

Submitted to the *Bernoulli*

arXiv: [math.PR/0000000](https://arxiv.org/abs/math.PR/0000000)

Optimal tuning of the Hybrid Monte-Carlo Algorithm

ALEXANDROS BESKOS¹ NATESH PILLAI² GARETH ROBERTS³ JESUS-MARIA SANZ-SERNA⁴ and ANDREW STUART⁵

¹*Department of Statistical Science, UCL, Gower Street, London, WC1E 6BT, UK*
E-mail: alex@stats.ucl.ac.uk

²*Department of Statistics, Harvard University, Cambridge, MA 02138-2901, USA*
E-mail: pillai@fas.harvard.edu

³*Department of Statistics, University of Warwick, Coventry, CV4 7AL, UK*
E-mail: gareth.o.roberts@warwick.ac.uk

⁴*Departamento de Matematica Aplicada, Facultad de Ciencias, Universidad de Valladolid, Spain*
E-mail: sanzsern@mac.uva.es

⁵*Mathematics Institute, University of Warwick, Coventry, CV4 7AL, UK*
E-mail: a.m.stuart@warwick.ac.uk

We investigate the properties of the Hybrid Monte Carlo algorithm (HMC) in high dimensions. HMC develops a Markov chain reversible w.r.t. a given target distribution Π by using separable Hamiltonian dynamics with potential $-\log \Pi$. The additional momentum variables are chosen at random from the Boltzmann distribution and the continuous-time Hamiltonian dynamics are then discretised using the leapfrog scheme. The induced bias is removed via a Metropolis-Hastings accept/reject rule. In the simplified scenario of independent, identically distributed components, we prove that, to obtain an $\mathcal{O}(1)$ acceptance probability as the dimension d of the state space tends to ∞ , the leapfrog step-size h should be scaled as $h = l \times d^{-1/4}$. Therefore, in high dimensions, HMC requires $\mathcal{O}(d^{1/4})$ steps to traverse the state space. We also identify analytically the asymptotically optimal acceptance probability, which turns out to be 0.651 (to three decimal places). This is the choice which optimally balances the cost of generating a proposal, which *decreases* as l increases (because fewer steps are required to reach the desired final integration time), against the cost related to the average number of proposals required to obtain acceptance, which *increases* as l increases.

Keywords: Hamiltonian dynamics, high dimensions, optimal acceptance probability, leapfrog scheme, squared jumping distance.

1. Introduction

The Hybrid Monte Carlo (HMC) algorithm originates from the physics literature [15] where it was introduced as a fast method for simulating molecular dynamics. It has since become popular in a number of application areas including statistical physics [17, 18, 43,

23, 1], computational chemistry [22, 24, 30, 42, 45], data assimilation [2], geophysics [30] and neural networks [32, 46]. The algorithm has also been proposed as a generic tool for Bayesian statistical inference [31, 12, 16].

Many practitioners believe that HMC improves on traditional Markov Chain Monte Carlo (MCMC) algorithms. There are heuristic arguments to suggest why HMC might perform better, in particular based on the idea that it breaks down *random walk-like* behaviour intrinsic to many MCMC algorithms such as the Random-Walk Metropolis (RWM) algorithm. However there is very little theoretical understanding of this phenomenon (though see [14]). This lack of theoretical guidance concerning the choice of the free parameters for the algorithm partly accounts for its relative obscurity in statistical applications. The aim of this paper is to provide insight into the behavior of HMC in high dimensions and develop theoretical tools for improving the efficiency of the algorithm.

HMC uses the derivative of the target probability log-density to guide the Monte-Carlo trajectory towards areas of high probability. The standard RWM algorithm [29] proposes *local*, symmetric moves around the current position. In many cases (especially in high dimensions) the variance of the proposal must be small for the corresponding acceptance probability to be satisfactory. However smaller proposal variance leads to higher autocorrelations, and large computing time to explore the state space. In contrast, and as discussed in the following sections, HMC exploits the information on the derivative of the log density to deliver guided, *global* moves, with higher acceptance probability. Thus HMC has the potential to effectively decorrelate by exploiting Hamiltonian evolution, conferring a potential advantage over random walk based methods whose effective decorrelation time is determined by random walk behaviour.

HMC is closely related to the so-called Metropolis-adjusted Langevin algorithm (abbrev. MALA) [39] which uses the derivative of the log-density to propose steepest-ascent moves in the state space. MALA employs *Langevin* dynamics; the proposal is derived from an Euler discretisation of a Langevin stochastic differential equation that leaves the target density invariant. We note here that the statisticians' use of the term 'Langevin dynamics' refers to the dynamics of a first order equation which physicists normally term 'Brownian dynamics'; this model is derived from the second order Langevin equation in the over-damped limit. Indeed the idea of using such dynamics as a proposal for Monte Carlo predates its appearance in the statistical literature [33, 40]. On the other hand, HMC uses *Hamiltonian* dynamics. The original variable q is seen as a 'location' variable and an auxiliary 'momentum' variable p is introduced; Hamilton's ordinary differential equations are used to generate moves in the enlarged (q, p) phase space. These moves preserve the total energy, a fact that implies, in probability terms, that they preserve the target density Π of the original q variable, provided that the initial momentum is chosen randomly from an appropriate Gaussian distribution. Although seemingly of different origin, MALA can be thought of as a 'localised' version of HMC in the case where Hamilton's equations are integrated for only one time-step before the accept/reject mechanism is applied [28]. We will return to this point below.

In practice, continuous-time Hamiltonian dynamics are discretised by means of a numerical scheme; the popular *Störmer-Verlet* or *leapfrog* scheme [19, 25, 41, 44] is currently the scheme of choice. This integrator does not conserve energy exactly and the induced

bias is corrected via a Metropolis-Hastings accept/reject rule. In this way, HMC develops a Markov chain reversible w.r.t. Π , whose transitions incorporate information on Π in a natural way.

In this paper we will investigate the properties of HMC in high dimensions and, in such a context, offer some guidance over the *optimal* specification of the free parameters of the algorithm. We assume that we wish to sample from a density Π on \mathbb{R}^N with

$$\Pi(Q) = \exp(-\mathcal{V}(Q)) , \quad (1.1)$$

for $\mathcal{V} : \mathbb{R}^N \rightarrow \mathbb{R}$. We study the simplified scenario where $\Pi(Q)$ consists of $d \gg 1$ independent identically distributed (iid) vector components,

$$\Pi(Q) = \exp\left(-\sum_{i=1}^d V(q_i)\right) , \quad V : \mathbb{R}^m \rightarrow \mathbb{R} ; \quad N = m \times d . \quad (1.2)$$

For the leapfrog integrator, we show analytically that, under suitable hypotheses on V and as $d \rightarrow \infty$, HMC requires $\mathcal{O}(d^{1/4})$ steps to traverse the state space, and furthermore, we identify the associated optimal acceptance probability.

To be more precise, if h is the step-size employed in the leapfrog integrator, then we show that the choice

$$\text{HMC} : \quad h = l \cdot d^{-1/4} \quad (1.3)$$

leads to an average acceptance probability which is of $\mathcal{O}(1)$ as $d \rightarrow \infty$: Theorem 3.6. This implies that $\mathcal{O}(d^{1/4})$ steps are required for HMC to make $\mathcal{O}(1)$ moves in state space. Furthermore we provide a result of perhaps greater practical relevance. We prove that, for the leapfrog integrator and as $d \rightarrow \infty$, the asymptotically *optimal* algorithm corresponds to a well-defined value of the acceptance probability, *independent of the particular target* Π in (1.2). This value is (to three decimal places) 0.651: Theorems 4.1 and 4.2. Thus, when applying HMC in high dimensions, one should try to tune the free algorithmic parameters to obtain an acceptance probability close to that value. We give the precise definition of optimality when stating the theorems but, roughly, it is determined by the choice of l which balances the cost of generating a proposal, which *decreases* as l increases (because fewer steps are required to reach the desired final integration time), against the cost related to the average number of proposals required to obtain acceptance, which *increases* as l increases.

The scaling $\mathcal{O}(d^{1/4})$ to make $\mathcal{O}(1)$ moves in state space contrasts favorably with the corresponding scalings $\mathcal{O}(d)$ and $\mathcal{O}(d^{1/3})$ required in a similar context by RWM and MALA respectively (see the discussion below). Furthermore, the full analysis provided in this paper for the leapfrog scheme may be easily extended to high-order, volume-preserving, reversible integrators. For such an integrator the corresponding scaling would be $\mathcal{O}(d^{1/(2\nu)})$, where ν (an integer) represents the order of the method. For the standard HMC algorithm, previous works have already established the relevance of the choice $h = \mathcal{O}(d^{-1/4})$ (by heuristic arguments, see [18]) and an optimal acceptance probability of around 0.7 (by numerical experiments, see [12]). Our analytic study of the scaling issues in HMC was prompted by these two papers.

We end this discussion with a transparent disclaimer about the range of validity of our optimal scaling results. Our work contains two central assumptions: (i) we work in the setting of an iid target measure; (ii) this iid target is defined via a single potential V which (see below) is assumed to grow no faster than quadratically at infinity so that the tails of the distribution are no lighter than Gaussian. Regarding (i), we expect that our results will extend to some problems with a non-product structure, provided that the resulting measure is ‘close’ to iid. Examples of such problems are contained in the work of Bédard [4, 6, 5], and the papers [9, 27, 35]. These last three papers show that optimal scaling results for RWM and MALA type algorithms extend directly to target measures which have a density with respect to a Gaussian, uniformly as dimension $d \rightarrow \infty$. Since a Gaussian measure is iid when represented in appropriate coordinates such measures are indeed close to the iid case as almost sure properties of the Gaussian measure are inherited by the target measure. However, for all optimal scaling analyses of RWM, MALA and HMC, the extent and manner in which the ‘close to iid’ assumption can be violated, and yet the same optimality criteria apply, remains an open and interesting research question. Regarding (ii) we note that integration of Hamiltonian systems with super-quadratic potentials (more precisely superlinear forces) typically requires adaptive time-step integration [41] and that an open and interesting research direction concerns the generalization of HMC algorithms to this situation. We discuss these issues related to possible relaxation of our key assumptions also in the conclusions section.

The paper is organized as follows. Section 2 presents the HMC method and reviews the literature concerning scaling issues for the RWM and MALA algorithms. Section 3 studies the asymptotic behaviour of HMC as the dimensionality grows, $d \rightarrow \infty$, including the key Theorem 3.6. The optimal tuning of HMC is discussed in Section 4, including the key Theorems 4.1 and 4.2. Sections 5 and 6 are technical. The first of them contains the derivation of the required numerical analysis estimates on the leapfrog integrator, with careful attention paid to the dependence of constants in error estimates on the initial condition; estimates of this kind are not available in the literature and may be of independent interest. Section 6 gathers the probabilistic proofs. We finish with some conclusions and discussion in Section 7.

2. Hybrid Monte Carlo (HMC)

The Hybrid Monte Carlo method is described from a statisticians perspective in [26]. Here we provide a precise definition of the algorithm, recalling several important concepts from the theory of Hamiltonian dynamics, such as volume-preservation, Liouville equation and reversible integration;¹ rather than repeat these classical definitions here, we refer the reader to the text [41].

¹‘reversible’ here has a different meaning from that employed in the study of Markov chains.

2.1. Hamiltonian Dynamics

Consider the Hamiltonian function:

$$\mathcal{H}(Q, P) = \frac{1}{2} \langle P, \mathcal{M}^{-1} P \rangle + \mathcal{V}(Q) ,$$

on \mathbb{R}^{2N} , where \mathcal{M} is a symmetric positive definite matrix (the ‘mass’ matrix). One should think of Q as the *location* argument and $\mathcal{V}(Q)$ as the potential energy of the system; P as the *momenta*, and $(1/2)\langle P, \mathcal{M}^{-1} P \rangle$ as the kinetic energy. Thus $\mathcal{H}(Q, P)$ gives the total *energy*: the sum of the potential and the kinetic energy. The Hamiltonian dynamics associated with \mathcal{H} are governed by

$$\frac{dQ}{dt} = \mathcal{M}^{-1} P, \quad \frac{dP}{dt} = -\nabla \mathcal{V}(Q) , \quad (2.1)$$

a system of ordinary differential equations whose solution flow Φ_t defined by

$$(Q(t), P(t)) = \Phi_t(Q(0), P(0))$$

possesses some key properties relevant to HMC:

- **1. Conservation of Energy:** The change in the potential becomes kinetic energy; *i.e.*, $\mathcal{H} \circ \Phi_t = \mathcal{H}$, for all $t > 0$, that is $\mathcal{H}(\Phi_t(Q(0), P(0))) = \mathcal{H}(Q(0), P(0))$, for all $t > 0$ and all initial conditions $(Q(0), P(0))$.
- **2. Conservation of Volume:** The volume element $dP dQ$ of the phase space is conserved under the mapping Φ_t .
- **3. Time Reversibility:** If \mathcal{S} denotes the symmetry operator:

$$\mathcal{S}(Q, P) = (Q, -P)$$

then $\mathcal{H} \circ \mathcal{S} = \mathcal{H}$ and

$$\mathcal{S} \circ (\Phi_t)^{-1} \circ \mathcal{S} = \Phi_t . \quad (2.2)$$

Thus, changing the sign of the initial velocity, evolving backwards in time, and changing the sign of the final velocity reproduces the forward evolution.

From the Liouville equation for (2.1) it follows that, if the initial conditions are distributed according to a probability measure with Lebesgue density depending only on $\mathcal{H}(Q, P)$, then this probability measure is preserved by the Hamiltonian flow Φ_t . In particular, if the initial conditions $(Q(0), P(0))$ of (2.1) are distributed with a density (proportional to, since we omit the normalising constant for the Gaussian part)

$$\exp(-\mathcal{H}(Q, P)) = \exp(-(1/2)\langle P, \mathcal{M}^{-1} P \rangle) \exp(-\mathcal{V}(Q)) ,$$

then, for all $t > 0$, the marginal density of $Q(t)$ will also be $\exp(-\mathcal{V}(Q))$. This suggests that integration of equations (2.1) might form the basis for an exploration of the target density $\exp(-\mathcal{V}(Q))$.

2.2. The HMC Algorithm

To formulate a practical algorithm, the continuous-time dynamics (2.1) must be discretised. The most popular *explicit* method is the Störmer-Verlet or leapfrog scheme (see [19, 25, 41] and the references therein) defined as follows. Assume a current state (Q_0, P_0) ; then, after one step of length $h > 0$ the system (2.1) will be at a state (Q_h, P_h) defined by the three-stage procedure:

$$P_{h/2} = P_0 - \frac{h}{2} \nabla \mathcal{V}(Q_0) ; \quad (2.3a)$$

$$Q_h = Q_0 + h \mathcal{M}^{-1} P_{h/2} ; \quad (2.3b)$$

$$P_h = P_{h/2} - \frac{h}{2} \nabla \mathcal{V}(Q_h) . \quad (2.3c)$$

The scheme gives rise to a map:

$$\Psi_h : (Q_0, P_0) \mapsto (Q_h, P_h)$$

which approximates the flow Φ_h . The solution at time T is approximated by taking $\lfloor \frac{T}{h} \rfloor$ leapfrog steps:

$$(Q(T), P(T)) = \Phi_T((Q(0), P(0))) \approx \Psi_h^{\lfloor \frac{T}{h} \rfloor}((Q(0), P(0))) .$$

Note that this is a *deterministic* computation. The map

$$\Psi_h^{(T)} := \Psi_h^{\lfloor \frac{T}{h} \rfloor}$$

may be shown to be volume preserving and time reversible (see [19, 25, 41]) but it does not exactly conserve energy. As a consequence the leapfrog algorithm does not share the property of equations (2.1) following from the Liouville equation, namely that the probability density function proportional to $\exp(-\mathcal{H}(Q, P))$ is preserved. In order to restore this property an accept/reject step must be added. The work in [31] provides a clear derivation of the required acceptance criterion.

We can now describe the complete HMC algorithm. Let the current state be Q . The next state for the HMC Markov chain is determined by the dynamics described in Table 1.

Due to the time reversibility and volume conservation properties of the integrator map $\Psi_h^{(T)}$, the algorithm in Table 1 defines (see [15, 31]) a Markov chain reversible w.r.t $\Pi(Q)$; sampling this chain up to equilibrium will provide correlated samples Q^n from $\Pi(Q)$. We note that the momentum P is merely an auxiliary variable and that the user of the algorithm is free to choose h , T and the mass matrix \mathcal{M} . In this paper we concentrate on the optimal choice of h , for high dimensional targets.

2.3. Connection with other Metropolis-Hastings Algorithms

Earlier research has studied the optimal tuning of other Metropolis-Hastings algorithms, namely the Random-Walk Metropolis (RWM) and the Metropolis-adjusted Langevin

HMC(Q):

(i) Sample a momentum $P \sim N(0, \mathcal{M})$.

(ii) Accept the proposed update Q' defined via $(Q', P') = \Psi_h^{(T)}(Q, P)$ w.p.:

$$a((Q, P), (Q', P')) := 1 \wedge \exp\{\mathcal{H}(Q, P) - \mathcal{H}(Q', P')\} .$$

Table 1. The Markov transition for the Hybrid Monte-Carlo algorithm. Iterative application for a given starting location Q^0 , will yield a Markov chain Q^0, Q^1, \dots

algorithm (MALA). In contrast with HMC, whose proposals involve a deterministic element, those algorithms use updates that are purely stochastic. For the target density $\Pi(Q)$ in (1.1), RWM is specified through the proposed update

$$Q' = Q + \sqrt{h} Z ,$$

with $Z \sim N(0, I)$ (this simple case suffices for our exposition, but note that Z may be allowed to have an arbitrary mean zero distribution), while MALA is determined through the proposal

$$Q' = Q + \frac{h}{2} \nabla \log \Pi(Q) + \sqrt{h} Z .$$

The density Π is invariant for both algorithms when the proposals are accepted with probability

$$a(Q, Q') = 1 \wedge \frac{\Pi(Q')T(Q', Q)}{\Pi(Q)T(Q, Q')} ,$$

where

$$T(x, y) = \mathbb{P}[Q' \in dy \mid Q = x] / dy$$

is the transition density of the proposed update (note that for RWM the symmetry of the proposal implies $T(Q, Q') = T(Q', Q)$).

The proposal distribution for MALA corresponds to the Euler discretization of the stochastic differential equation (SDE)

$$dQ = \frac{1}{2} \nabla \log \Pi(Q) dt + dW ,$$

for which Π is an invariant density (here W denotes a standard multivariate Brownian motion with covariance I). One can easily check that HMC and MALA are connected because HMC reduces to MALA when $T \equiv h$, *i.e.*, when the algorithm makes only a single leapfrog step at each transition of the chain.

Assume now that RWM and MALA are applied with the scalings

$$\text{RWM : } h = l \cdot d^{-1}, \quad \text{MALA : } h = l \cdot d^{-1/3}, \quad (2.4)$$

for some constant $l > 0$, in the simplified scenario where the target Π has the iid structure (1.2) with $m = 1$. The papers [36], [37] prove that, as $d \rightarrow \infty$ and under regularity conditions on V (the function V must be seven times differentiable², with all derivatives having polynomial growth bounds, and all moments of $\exp(-V)$ must be finite), the acceptance probability approaches a nontrivial value:

$$\mathbb{E}[a(Q, Q')] \rightarrow a(l) \in (0, 1)$$

(the limit $a(l)$ is different for each of the two algorithms). Furthermore, if q_1^0, q_1^1, \dots denotes the projection of the trajectory Q^0, Q^1, \dots onto its first coordinate, in the above scenario it is possible to show ([36], [37]) the convergence of the continuous-time interpolation

$$\text{RWM : } t \mapsto q_1^{\lfloor t \cdot d \rfloor}, \quad \text{MALA : } t \mapsto q_1^{\lfloor t \cdot d^{1/3} \rfloor} \quad (2.5)$$

($\lfloor x \rfloor$ denoting the integer part of $x \in \mathbb{R}$) to the diffusion process governed by the SDE

$$dq = -\frac{1}{2} l a(l) V'(q) dt + \sqrt{l a(l)} dw, \quad (2.6)$$

(w represents a standard Brownian motion). In view of (2.4), (2.5) and (2.6) we deduce that the RWM and MALA algorithms cost $\mathcal{O}(d^2)$ and $\mathcal{O}(d^{4/3})$ respectively to explore the invariant measure in stationarity, for product measures where the cost of each step of the algorithm is $\mathcal{O}(d)$ (since, recall, m is fixed and $d \rightarrow \infty$). Furthermore, as the product $l a(l)$ determines the *speed* of the limiting diffusion the state space will be explored faster for the choice l_{opt} of l that maximises $l a(l)$. While l_{opt} depends on the target distribution, it turns out that the optimal acceptance probability $a(l_{opt})$ is independent of V . In fact, with three decimal places, one finds:

$$\text{RWM : } a(l_{opt}) = 0.234, \quad \text{MALA : } a(l_{opt}) = 0.574.$$

Asymptotically as $d \rightarrow \infty$, this analysis identifies algorithms that may be regarded as *uniformly* optimal, because, as discussed in [38], ergodic averages of trajectories corresponding to $l = l_{opt}$ provide optimal estimation of expectations $\mathbb{E}[f(q)]$, $q \sim \exp(-V)$, irrespectively of the choice of the (regular) function f . These investigations of the optimal tuning of RWM and MALA have been subsequently extended in [9] and [10] to non-product target distributions.

For HMC we show that the scaling (1.3) leads to an average acceptance probability of $\mathcal{O}(1)$ and hence to a cost of $\mathcal{O}(d^{5/4})$ to make the $\mathcal{O}(1)$ moves necessary to explore the (product) invariant measure. However, in contrast to RWM and MALA, we are not able to provide a simple description of the limiting dynamics of a single coordinate of the Markov chain. Consequently optimality is harder to define.

²this is mostly a technical requirement which may be relaxed.

3. Hybrid Monte Carlo in the Limit $d \rightarrow \infty$.

The primary aim of this section is to prove Theorem 3.6 concerning the scaling of the step-size h in HMC. We also provide some insight into the limiting behaviour of the resulting Markov chain, under this scaling, in Propositions 3.8 and 3.9.

3.1. HMC in the iid Scenario

We now study the asymptotic behaviour of the HMC algorithm in the iid scenario (1.2), when the number d of ‘particles’ goes to infinity. Due to such a scenario for our target, a global $m \times d$ -dimensional implementation of the Hamiltonian dynamics (2.1) or indeed of the practical leapfrog scheme (2.3), can now be decomposed into d independent implementations along each of the identical m -dimensional constituent components (assuming that the auxiliary variable P is also chosen to have a similar iid structure). We will exploit this simplified structure in our analysis.

We write $Q = (q_i)_{i=1}^d$ and $P = (p_i)_{i=1}^d$ to distinguish the individual components, and use the following notation for the combination location/momentum:

$$X = (x_i)_{i=1}^d; \quad x_i := (q_i, p_i) \in \mathbb{R}^{2m}.$$

We denote by \mathcal{P}_q and \mathcal{P}_p the projections onto the position and momentum components of x , i.e. $\mathcal{P}_q(q, p) = q$, $\mathcal{P}_p(q, p) = p$. We have:

$$\mathcal{H}(Q, P) = \sum_{i=1}^d H(q_i, p_i); \quad H(q, p) := \frac{1}{2} \langle p, M^{-1}p \rangle + V(q) - \log(c),$$

where M is a $m \times m$ symmetric, positive definite matrix. Also, $c > 0$ is the normalising constant for the Gaussian part, i.e. $c^{-1} = \int e^{-\frac{1}{2} \langle p, M^{-1}p \rangle} dp$. We have only included it here to avoid repeatedly using a normalising constant in the mathematical expressions for expectations used below. Of course, HMC only uses differences in the energy $H(q, p)$ or it’s derivatives, so normalising constants under the distribution of p or q are not required by the algorithm. The Hamiltonian differential equations for a single (m -dimensional) particle are

$$\frac{dq}{dt} = M^{-1}p, \quad \frac{dp}{dt} = -\nabla V(q), \quad (3.1)$$

where $V : \mathbb{R}^m \rightarrow \mathbb{R}$. We denote the corresponding flow by φ_t and the leapfrog solution operator over one h -step by ψ_h . Thus the acceptance probability for the evolution of the d particles is given by (see Table 1):

$$a(X, Y) = 1 \wedge \exp \left(\sum_{i=1}^d [H(x_i) - H(\psi_h^{(T)}(x_i))] \right) \quad (3.2)$$

with $Y = (y_i)_{i=1}^d = \Psi_h^{(T)}(X)$ denoting the HMC proposal.

As mentioned above, due to the iid scenario, the leapfrog scheme (2.3) disentangles into independent implementations for each of the d particles (q_i, p_i) , with the different particles being only connected through the accept/reject decision (3.2).

3.2. Energy Increments

Our first aim is to estimate (in an analytical sense) the exponent in the right-hand side of (3.2). Since the d particles play the same role, it is sufficient to study a single term $H(x_i) - H(\psi_h^{(T)}(x_i))$. We set

$$\Delta(x, h) := H(\psi_h^{(T)}(x)) - H(\varphi_T(x)) = H(\psi_h^{(T)}(x)) - H(x) . \quad (3.3)$$

This is the energy change, due to the leapfrog scheme, over $0 \leq t \leq T$, with step-size h and initial condition x , which by conservation of energy under the true dynamics, is simply the energy error at time T . We will study the first and second moments:

$$\begin{aligned} \mu(h) &:= \mathbb{E} [\Delta(x, h)] = \int_{\mathbb{R}^{2m}} \Delta(x, h) e^{-H(x)} dx , \\ s^2(h) &:= \mathbb{E} [\Delta(x, h)^2] , \end{aligned}$$

and the corresponding variance

$$\sigma^2(h) = s^2(h) - \mu^2(h) .$$

If the integrator were exactly energy-preserving, one would have $\Delta \equiv 0$ and all proposals would be accepted. However it is well known that the size of $\Delta(x, h)$ is in general no better than the size of the integration error $\psi_h^{(T)}(x) - \varphi_T(x)$, *i.e.* $\mathcal{O}(h^2)$. In fact, under natural smoothness assumptions on V the following condition holds (see Section 5 for a proof):

Condition 3.1. *There exist functions $\alpha(x)$, $\rho(x, h)$ such that*

$$\Delta(x, h) = h^2 \alpha(x) + h^2 \rho(x, h) \quad (3.4)$$

with $\lim_{h \rightarrow 0} \rho(x, h) = 0$.

Furthermore in the proofs of the theorems below we shall use an additional condition to control the variation of Δ as a function of x . This condition will be shown in Section 5 to hold under suitable assumptions on the growth of V and its derivatives.

Condition 3.2. *There exists a function $D : \mathbb{R}^{2m} \rightarrow \mathbb{R}$ such that*

$$\sup_{0 \leq h \leq 1} \frac{|\Delta(x, h)|^2}{h^4} \leq D(x) ,$$

with

$$\int_{\mathbb{R}^{2m}} D(x) e^{-H(x)} dx < \infty .$$

Key to the proof of Theorem 3.6 is the fact that the average energy increment scales as $\mathcal{O}(h^4)$. We show this in Proposition 3.4 using the following simple lemma that holds for general volume preserving, time reversible integrators:

Lemma 3.3. *Let $\psi_h^{(T)}$ be any volume preserving, time reversible numerical integrator of the Hamiltonian equations (3.1) and $\Delta(x, h) : \mathbb{R}^{2m} \times \mathbb{R}_+ \rightarrow \mathbb{R}$ be as in (3.3). If $g : \mathbb{R} \rightarrow \mathbb{R}$ is an odd function then:*

$$\int_{\mathbb{R}^{2m}} g(\Delta(x, h)) e^{-H(x)} dx = - \int_{\mathbb{R}^{2m}} g(\Delta(x, h)) e^{-H(\psi_h^{(T)}(x))} dx$$

provided at least one of the integrals above exist. If g is an even function, then:

$$\int_{\mathbb{R}^{2m}} g(\Delta(x, h)) e^{-H(x)} dx = \int_{\mathbb{R}^{2m}} g(\Delta(x, h)) e^{-H(\psi_h^{(T)}(x))} dx ,$$

provided at least one of the integrals above exist.

Proof. See Section 6. □

Applying this lemma with $g(u) = u$, we obtain

$$\mu(h) = - \int_{\mathbb{R}^{2m}} \Delta(x, h) e^{-H(\psi_h^{(T)}(x))} dx ,$$

which implies that

$$2\mu(h) = \int_{\mathbb{R}^{2m}} \Delta(x, h) [1 - \exp(-\Delta(x, h))] e^{-H(x)} dx . \quad (3.5)$$

We now use first the inequality $|e^u - 1| \leq |u|(e^u + 1)$ and then Lemma 3.3 with $g(u) = u^2$ to conclude that

$$\begin{aligned} |2\mu(h)| &\leq \int_{\mathbb{R}^{2m}} |\Delta(x, h)|^2 e^{-H(\psi_h^{(T)}(x))} dx + \int_{\mathbb{R}^{2m}} |\Delta(x, h)|^2 e^{-H(x)} dx \\ &\leq 2 \int_{\mathbb{R}^{2m}} |\Delta(x, h)|^2 e^{-H(x)} dx = 2s^2(h) . \end{aligned} \quad (3.6)$$

The bound in (3.6) is important: it shows that the average of $\Delta(x, h)$ is actually of the order of (the average of) $\Delta(x, h)^2$. Since for the second-order leapfrog scheme $\Delta(x, h) = \mathcal{O}(h^2)$, we see from (3.6) that we may expect the average $\mu(h)$ to actually behave as $\mathcal{O}(h^4)$. This is made precise in the following theorem.

Proposition 3.4. *If the potential V is such that Conditions 3.1 and 3.2 hold for the leapfrog integrator $\psi_h^{(T)}$, then*

$$\lim_{h \rightarrow 0} \frac{\mu(h)}{h^4} = \mu , \quad \lim_{h \rightarrow 0} \frac{\sigma^2(h)}{h^4} = \Sigma ,$$

for the constants:

$$\Sigma = \int_{R^{2m}} \alpha^2(x) e^{-H(x)} dx ; \quad \mu = \Sigma/2 .$$

Proof. See Section 6. □

Next, we perform explicit calculations for the example of the harmonic oscillator and verify (for this case) the conclusions of Proposition 3.4.

Example 3.5 (Harmonic Oscillator). *Consider the Hamiltonian*

$$H(q, p) = \frac{1}{2}p^2 + \frac{1}{2}q^2$$

that gives rise to the system

$$\begin{pmatrix} dq/dt \\ dp/dt \end{pmatrix} = \begin{pmatrix} p \\ -q \end{pmatrix} ,$$

with solutions

$$\begin{pmatrix} q(t) \\ p(t) \end{pmatrix} = \begin{pmatrix} \cos(t) & \sin(t) \\ -\sin(t) & \cos(t) \end{pmatrix} \begin{pmatrix} q(0) \\ p(0) \end{pmatrix} .$$

In this case, the leapfrog integration can be written as:

$$\psi_h = \psi_h(q, p) = \begin{pmatrix} 1 - h^2/2 & h \\ -h + h^3/4 & 1 - h^2/2 \end{pmatrix} \begin{pmatrix} q \\ p \end{pmatrix} = \Xi \begin{pmatrix} q \\ p \end{pmatrix} ,$$

and, accordingly, the numerical solution after $\lfloor \frac{1}{h} \rfloor$ steps is given by:

$$\psi_h^{(1)}(q, p) = \Xi^{\lfloor \frac{1}{h} \rfloor} \begin{pmatrix} q \\ p \end{pmatrix} .$$

Diagonalizing Ξ and exponentiating yields:

$$\Xi^n = \begin{pmatrix} \cos(\theta n) & \frac{1}{\sqrt{1-h^2/4}} \sin(\theta n) \\ -\sqrt{1-h^2/4} \sin(\theta n) & \cos(\theta n) \end{pmatrix}$$

where $\theta = \cos^{-1}(1 - h^2/2)$. Using, for instance, MATHEMATICA, one can now obtain the Taylor expansion:

$$\Delta(x, h) = H(\psi_h^{(1)}(x)) - H(x) = h^2\alpha(x) + h^4\beta(x) + \mathcal{O}(h^6)$$

where:

$$\alpha(q, p) = ((p^2 - q^2) \sin^2(1) + pq \sin(2)) / 8 ;$$

$$\beta(q, p) = \left(-q^2 \sin(2) + pq(2 \cos(2) + 3 \sin(2)) + p^2(3 - 3 \cos(2) + \sin(2)) \right) / 192 .$$

Notice that, in the stationary regime, q, p are standard normal variables. Therefore, the expectation of $\alpha(x)$ is 0. Tedious calculations give:

$$\text{Var}[\alpha(x)] = \frac{1}{16} \sin^2(1) , \quad \mathbb{E}[\beta(x)] = \frac{1}{32} \sin^2(1) ,$$

in agreement with Proposition 3.4.

3.3. Expected Acceptance Probability

We are now in a position to identify the scaling for h that gives non-trivial acceptance probability as $d \rightarrow \infty$.

Theorem 3.6. *Assume that the potential V is such that the leapfrog integrator $\psi_h^{(T)}$ satisfies Conditions 3.1 and 3.2 and that*

$$h = l \cdot d^{-1/4} , \quad (3.7)$$

for a constant $l > 0$. Then in stationarity, i.e. for $X \sim \exp(-\mathcal{H})$ and $Y = \Psi_h^{(T)}(X)$,

$$\lim_{d \rightarrow \infty} \mathbb{E} [a(X, Y)] = 2 \Phi(-l^2 \sqrt{\Sigma}/2) =: a(l)$$

where the constant Σ is as defined in Proposition 3.4.

Proof. To grasp the main idea, note that the acceptance probability (3.2) is given by

$$a(X, Y) = 1 \wedge e^{R_d} ; \quad R_d = - \sum_{i=1}^d \Delta(x_i, h) . \quad (3.8)$$

Due to the simple structure of the target density and stationarity, the terms $\Delta(x_i, h)$ being added in (3.8) are iid random variables. Since the expectation and variance of $\Delta(x, h)$ are both $\mathcal{O}(h^4)$ and we have d terms, the natural scaling to obtain a distributional limit is given by (3.7). Then $R_d \approx N(-\frac{1}{2}l^4\Sigma, l^4\Sigma)$ and the desired result follows. See Section 6 for a detailed proof. \square

In Theorem 3.6 the limit acceptance probability arises from the use of the Central Limit Theorem. If Condition 3.2 is not satisfied and $\sigma^2(h) = \infty$, then a Gaussian limit is not guaranteed and it may be necessary to consider a different scaling to obtain a heavy tailed limiting distribution such as, say, a stable law.

The scaling (3.7) is a direct consequence of the fact that the leapfrog integrator possesses second order accuracy. Arguments similar to those used above prove that the use of a volume-preserving, symmetric ν -th order integrator would result in a scaling $h = \mathcal{O}(d^{-1/(2\nu)})$ (ν is an even integer) to obtain an acceptance probability of $\mathcal{O}(1)$.

3.4. Displacement of one Particle in a Transition

We now turn our attention to the displacement $q_1^{n+1} - q_1^n$ of a single particle in a transition $n \rightarrow n + 1$ of the chain. Note that clearly

$$q_1^{n+1} = I^n \cdot \mathcal{P}_q \psi_h^{(T)}(q_1^n, p_1^n) + (1 - I^n)q_1^n; \quad I^n = \mathbb{I}_{U^n \leq a(X^n, Y^n)} . \quad (3.9)$$

with U_n having a uniform distribution in $[0, 1]$. While Conditions 3.1 and 3.2 above refer to the error in energy, the proof of the next results requires a condition on the leapfrog integration error in the dynamic variables q and p . In Section 5 we describe conditions on V that guarantee the fulfillment of this condition.

Condition 3.7. *There exists a function $E : \mathbb{R}^{2m} \rightarrow \mathbb{R}$ such that*

$$\sup_{0 \leq h \leq 1} \frac{|\psi_h^{(T)}(x) - \varphi_T(x)|}{h^2} \leq E(x) ,$$

with

$$\int_{\mathbb{R}^{2m}} E(x)^4 e^{-H(x)} dx < \infty .$$

Under the scaling (3.7) and at stationarity, the second moment $\mathbb{E}[(q_1^{n+1} - q_1^n)^2]$ will also approach a nontrivial limit:

Proposition 3.8. *Assume that the hypotheses of Theorem 3.6 and Condition 3.7 hold and, furthermore, that the density $\exp(-V(q))$ possesses finite fourth moments. Then, in stationarity,*

$$\lim_{d \rightarrow \infty} \mathbb{E}[(q_1^{n+1} - q_1^n)^2] = C_J \cdot a(l)$$

where the value of the constant C_J is given by

$$C_J = \mathbb{E}[(\mathcal{P}_q \varphi_T(q, p) - q)^2] ; \quad (q, p) \sim \exp(-H(q, p)) .$$

Proof. See Section 6. □

Notice that the computational work required to integrate up to a fixed time T is inversely proportional to the parameter l . Thus Proposition 3.8 suggests that it is reasonable to choose a value for l that maximizes the quantity $a(l)l$. This choice of l is optimal in the sense that it seeks a middle ground between smaller values of l , which lead to a higher acceptance probability (and hence larger mean square jumps) but need more computational work, and large values of l which have a smaller acceptance probability (and hence smaller mean square jumps) but need less computational resources. In Section 4 we expand this idea, define a precise notion of optimality which encodes this trade-off, and derive the resulting optimal acceptance probability.

3.5. The Limit Dynamics

We now discuss the limiting dynamics of the Markov chain, under the same assumptions as in Proposition 3.8. For HMC (as for RWM or MALA) the marginal process $\{q_1^n\}_{n \geq 0}$ is not Markovian w.r.t. its own filtration since its dynamics depend on the current position of all d particles via the acceptance probability $a(X^n, Y^n)$ (see (3.9)). In the case of MALA and RWM, $\{q_1^n\}_{n \geq 0}$ is *asymptotically* Markovian: as $d \rightarrow \infty$ the effect of the rest of the particles gets averaged to a constant via the Strong Law of Large Numbers. This allows for the interpolants of (2.5) to converge to solutions of the SDE (2.6), which defines a Markov process. We will now argue that for HMC $\{q_1^n\}_{n \geq 0}$ cannot be expected

to be *asymptotically* Markovian. In order to simplify the exposition we will not present all the technicalities of the argument that follows.

It is well known (see for instance [44]) that, due to time reversibility and under suitable smoothness assumptions on V , the energy increments of the leapfrog integrator may be expanded in even powers of h as follows (cf. (3.4)):

$$\Delta(x, h) = h^2\alpha(x) + h^4\beta(x) + \mathcal{O}(h^6) .$$

Necessarily $\mathbb{E}[\alpha(x)] = 0$ since from Proposition 3.4 we know that $\mathbb{E}[\Delta(x, h)] = \mathcal{O}(h^4)$. Ignoring $\mathcal{O}(h^6)$ -terms, we can write:

$$a(X^n, Y^n) = 1 \wedge e^{R_{1,d}^n + R_{2,d}^n}$$

with

$$R_{1,d}^n = -h^2 \sum_{i=1}^d \{ \alpha(x_i^n) - \mathbb{E}[\alpha(x_i^n) | q_i^n] \} - h^4 \sum_{i=1}^d \beta(x_i^n) ,$$

$$R_{2,d}^n = -h^2 \sum_{i=1}^d \mathbb{E}[\alpha(x_i^n) | q_i^n] .$$

Under appropriate conditions, $R_{1,d}^n$ converges, as $d \rightarrow \infty$, to a Gaussian limit independent of the σ -algebra $\sigma(q_1^n, q_2^n, \dots)$. To see that, note that, due to the Strong Law of Large Numbers and since $h^4 = l^4/d$, the second sum in $R_{1,d}^n$ converges a.s. to a constant. *Conditionally* on $\sigma(q_1^n, q_2^n, \dots)$, the distributional limit of the first term in $R_{1,d}^n$ is Gaussian with zero mean and a variance determined by the almost surely constant limit of $h^4 \sum_{i=1}^d \{ \alpha(x_i^n) - \mathbb{E}[\alpha(x_i^n) | q_i^n] \}^2$; this follows from the Martingale Central Limit Theorem (see e.g. Theorem 3.2 of [21]). On the other hand, the limit distribution of $R_{2,d}^n$ is Gaussian with zero mean but, in general, cannot be asymptotically independent of $\sigma(q_1^n, q_2^n, \dots)$. Instead, it seems that $R_{2,d}^n$ will result to a quantity appearing in the acceptance probability that is non-trivial as $d \rightarrow \infty$ and makes it impossible for having a Markovian limit for the trajectory of q_1 . In the case of RWM or MALA, the conditional expectations that play the role played here by $\mathbb{E}[\alpha(x_i^n) | q_i^n]$ are identically zero (see the expansions for the acceptance probability in [36] and [37]) and this implies that the corresponding acceptance probabilities are asymptotically independent from $\sigma(q_1^n, q_2^n, \dots)$ and that the marginal processes $\{q_1^n\}_{n \geq 0}$ are asymptotically Markovian.

The last result in this section provides insight into the limit dynamics of $\{q_1^n\}_{n \geq 0}$:

Proposition 3.9. *Let $Q^n \sim \Pi(Q)$, define*

$$q_1^{n+1} = l^n \cdot \mathcal{P}_q \varphi_T(q_1^n, p_1^n) + (1 - l^n)q_1^n; \quad l^n = \mathbb{I}_{U^n \leq a(l)} ,$$

and consider q_1^{n+1} in (3.9). Then, under the hypotheses of Proposition 3.8, as $d \rightarrow \infty$:

$$(q_1^n, q_1^{n+1}) \xrightarrow{\mathcal{L}} (q_1^n, \mathbf{q}_1^{n+1}) .$$

Proof. See Section 6. □

This proposition provides a simple description of the asymptotic behaviour of the one-transition dynamics of the marginal trajectories of HMC. As $d \rightarrow \infty$, with probability $a(l)$, the HMC particle moves under the *correct* Hamiltonian dynamics. However, the deviation from the true Hamiltonian dynamics, due to the energy errors accumulated from the leapfrog integration of all d particles, gives rise to the alternative event of staying at the current position q^n , with probability $1 - a(l)$.

4. Optimal Tuning of HMC

In the previous section we addressed the question of how to scale the step-size in the leapfrog integration in terms of the dimension d , leading to Theorem 3.6. In this section we refine this analysis and study the choice of constant l in (3.7). Regardless of the metrics used to measure the efficiency of the algorithm, a good choice of l in (3.7) has to balance the amount of work needed to simulate a full T -leg (interval of length T) of the Hamiltonian dynamics and the probability of accepting the resulting proposal. Increasing l decreases the acceptance probability but also decreases the computational cost of each T -leg integration; decreasing l will yield the opposite effects, suggesting an optimal value of l . In this section we present an analysis that avoids the complex calculations typically associated with the estimation of mixing times of Markov chains, but still provides useful guidance regarding the choice of l . We provide two alternative ways of doing this, summarized in Theorems 4.1 and Theorem 4.2.

4.1. Asymptotically Optimal Acceptance Probability

Clearly, the number of leapfrog steps of length h needed to compute a proposal is $\lceil T/h \rceil$. Furthermore, at each step of the chain, it is necessary to evaluate $a(X, Y)$ and sample P . Thus the computing time for a single proposal will be

$$C_{l,d} := \left\lceil \frac{T d^{1/4}}{l} \right\rceil \cdot d \cdot C_{LF} + d \cdot C_O, \quad (4.1)$$

for some constants C_{LF} , C_O that measure, for one particle, the leapfrog costs and the overheads. Let $E_{l,d}$ denote the expected computing time until the first accepted T -leg, in stationarity. Recall that Q denotes the vector of positions within the Hamiltonian model so that $X = (Q, P)$. If N denotes the number of proposals until (and including) the first to be accepted, then

$$E_{l,d} = C_{l,d} \mathbb{E}[N] = C_{l,d} \mathbb{E}[\mathbb{E}[N | Q]] = C_{l,d} \mathbb{E} \left[\frac{1}{\mathbb{E}[a(X, Y) | Q]} \right].$$

Here we have used the fact that, given the locations Q , the number of proposed T -legs follows a geometric distribution with probability of success $\mathbb{E}[a(X, Y) | Q]$. Jensen's

inequality yields

$$E_{l,d} \geq \frac{C_{l,d}}{\mathbb{E}[a(X,Y)]} =: E_{l,d}^* , \quad (4.2)$$

and, from (4.1) and Theorem 3.6, we conclude that:

$$\lim_{d \rightarrow \infty} d^{-5/4} \times E_{l,d}^* = \frac{T C_{LF}}{a(l) l} .$$

A sensible choice for l is that which minimizes the asymptotic cost $E_{l,d}^*$, that is:

$$l_{opt} = \arg \max_{l > 0} \text{eff}(l); \quad \text{eff}(l) := a(l) l .$$

The value of l_{opt} will in general depend on the specific target distribution under consideration. However, by expressing eff as a function of $a = a(l)$, we may write

$$\text{eff} = \left(\frac{\sqrt{2}}{\Sigma^{\frac{1}{4}}} \right) \cdot a \cdot \left(\Phi^{-1} \left(1 - \frac{a}{2} \right) \right)^{\frac{1}{2}} \quad (4.3)$$

and this equality makes it apparent that $a(l_{opt})$ *does not vary with the selected target*. Fig.1 illustrates the mapping $a \mapsto \text{eff}(a)$; different choices of target distribution only change the vertical scale. In summary, we have:

Theorem 4.1. *Under the hypotheses of Theorem 3.6 and as $d \rightarrow \infty$, the measure of cost $E_{l,d}^*$ defined in (4.2) is minimised for the choice l_{opt} of l that leads to the value of $a = a(l)$ that maximises (4.3). Rounded to 3 decimal places the (target independent) optimal value of the limit probability a is*

$$a(l_{opt}) = 0.651 .$$

The optimal value identified in the preceding theorem is based on the quantity $E_{l,d}^*$ that underestimates the expected number of proposals. It may be assumed that the practical optimal average acceptance probability is in fact *greater than* or equal to 0.651. In the next subsection we use an alternative measure of efficiency: the expected squared jumping distance. Consideration of this alternative metric will also lead to the same asymptotically optimal acceptance probability of precisely 0.651 as did the minimisation of $E_{l,d}^*$. This suggests that, as $d \rightarrow \infty$, the consequences of the fact that $E_{l,d}^*$ underestimates $E_{l,d}$ become negligible; proving analytically such a conjecture seems hard given our current understanding of the limiting HMC dynamics.

4.2. Squared Jumping Distance

We now consider the chain Q^0, Q^1, \dots in stationarity (*i.e.* $Q^0 \sim \Pi(Q)$) and account for the computing cost $C_{l,d}$ in (4.1) by introducing the continuous-time process $Q^{N(t)}$, where $\{N(t); t \geq 0\}$ denotes a Poisson process, independent of the HMC Markov chain,

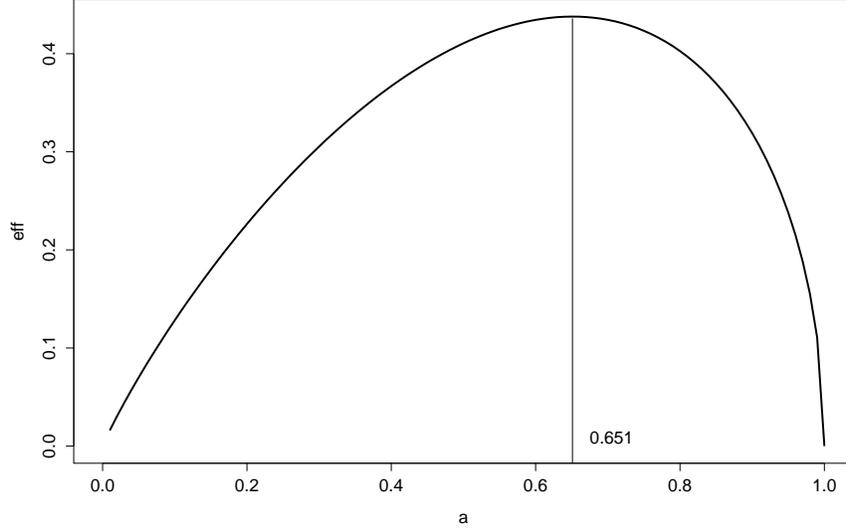


Figure 1. The efficiency function $\text{eff} = \text{eff}(a)$.

of intensity $\lambda_d = 1/C_{l,d}$. If $q_d(t) := q_1^{N(t)}$ denotes the projection of $Q^{N(t)}$ onto the first particle and $\delta > 0$ is a time increment, we measure the efficiency of HMC algorithms by using the expected squared jump distance:

$$\mathcal{SJD}_d(\delta) = \mathbb{E}[(q_d(t + \delta) - q_d(t))^2] .$$

This measure of efficiency is a fairly standard one: see [34, 9] for example.

In the HMC algorithm the computational time (cost) expended between successive steps of the Markov chain is essentially fixed and equal to $C_{l,d}$. Using an auxiliary Poisson process instead with mean interarrival time equal to $C_{l,d}$ is merely a device that allows for the definition of processes (over the different choices of l) that take the computational time per step (that changes with l) under consideration in a reasonable manner and can be meaningfully compared via an easy to calculate measure such as $\mathcal{SJD}_d(\delta)$.

The following result shows that $\mathcal{SJD}_d(\delta)$ is indeed asymptotically maximized by maximizing $a(l)l$:

Theorem 4.2. *Under the hypotheses of Proposition 3.8:*

$$\lim_{d \rightarrow \infty} d^{5/4} \times \mathcal{SJD}_d = \frac{C_J \delta}{T C_{LF}} \times a(l)l .$$

Proof. See Section 6. □

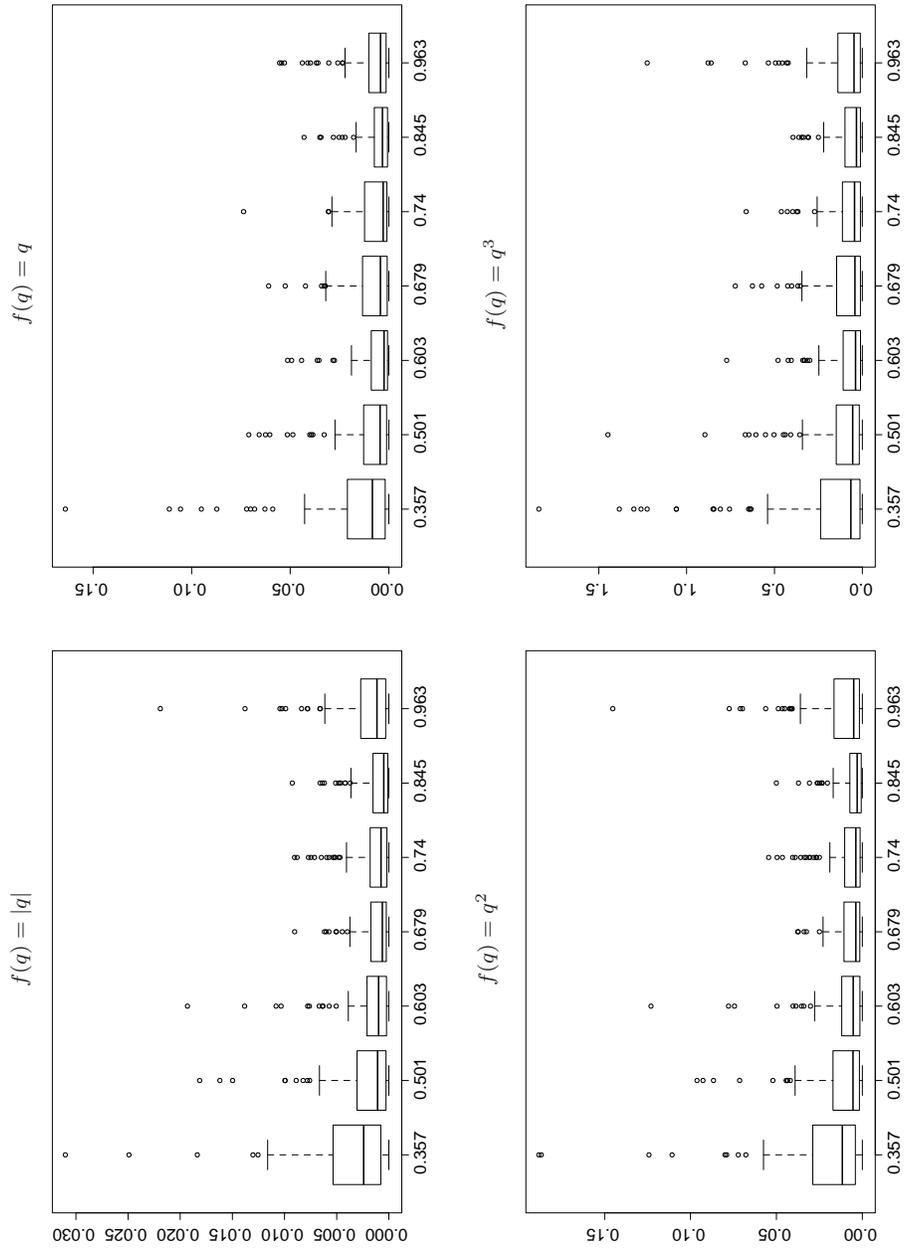


Figure 2. Boxplots of Squared Errors (SEs) from Monte-Carlo averages of HMC with $T = 1$ for 7 different selections of number of leapfrog steps or step-sizes h (corresponding to the different boxplots in each panel); the number of leapfrog steps used in the 7 scenaria were (6, 7, 8, 9, 10, 13, 27). We ran HMC 120 times; every run was allowed a computing time of 30s. Each boxplot corresponds to the 120 SEs in estimating $\mathbb{E}[f(q)]$, for a particular h and $f(\cdot)$. Written at the bottom of each boxplots is the median of the 120 empirical average acceptance probabilities for the corresponding h . (Notice that these medians change in a non-linear fashion from one boxplot to the next.)
 imsart-bj ver. 2009/08/13 file: hmc.tex date: October 12, 2011

4.3. Optimal Acceptance Probability in Practice

As $d \rightarrow \infty$, the computing time required for a proposal scales as $1/l$ (see (4.1)) and the number of proposals that may be performed in a given amount of time scales as l . Inspection of (4.1) reveals however that selecting a big value of l gives the full benefit of a proportional increase of the number of proposals only asymptotically, and at the slow rate of $\mathcal{O}(d^{-1/4})$. On the other hand, the average acceptance probability converges at the faster rate $\mathcal{O}(d^{-1/2})$ (this is an application of Stein's method, see e.g. [3]). These considerations suggest that unless $d^{-1/4}$ is very small the algorithm will tend to benefit from average acceptance probabilities higher than 0.651.

Fig.2 shows the results of a numerical study on HMC. The target distribution is a product of $d = 10^5$ standard Gaussian densities $N(0,1)$. We have applied HMC with different choices of the step-size h but the same length of time horizon $T = 1$ and, in all cases, allowed the algorithm to run during a computational time t_{comp} of 30 seconds. We used Monte-Carlo averages of the output

$$\hat{f} = \frac{1}{N_{t_{comp}}} \sum_{n=1}^{N_{t_{comp}}} f(q_1^n)$$

to estimate, for different choices of f , the expectation $\mathbb{E}[f] = \mathbb{E}[f(q)]$, $q \sim N(0,1)$; here $N_{t_{comp}}$ denotes the number of T -legs carried out within the allowed time t_{comp} . For each choice of h we ran the HMC algorithm 120 times.

Each of the four panels in Fig.2 corresponds to a different choice of $f(\cdot)$. In each of the panels, the various boxplots correspond to choices of h ; at the bottom of each boxplot we have written the median of the 120 empirical average acceptance probabilities. The boxplots themselves use the 120 realizations of the squared distances: $(\hat{f} - \mathbb{E}[f])^2$. The shape of the boxplots endorses the point made above, that the optimal acceptance probability for large (but finite) d is larger than the asymptotically optimal value of 0.651.

5. Estimates for the Leapfrog Algorithm

In this section we identify analytical hypotheses on V under which Conditions 3.1, 3.2 and 3.7 in Section 3 hold.

We set $f := -\nabla V$ (the 'force') and denote by $f'(q) := f^{(1)}(q), f^{(2)}(q), \dots$ the successive Fréchet derivatives of f at q . Thus, at a fixed q , $f^{(k)}(q)$ is a multilinear operator from $(\mathbb{R}^m)^{k+1}$ to \mathbb{R} . For the rest of this section we will use the following assumptions on V :

Assumptions 5.1. *The function $V : \mathbb{R}^m \rightarrow \mathbb{R}$ satisfies:*

- (i) $V \in C^4(\mathbb{R}^m \rightarrow \mathbb{R}_+)$.
- (ii) $f', f^{(2)}, f^{(3)}$ are uniformly bounded by a constant B .

These assumptions imply that the potential $V(q)$ can grow at most quadratically at infinity as $|q| \rightarrow \infty$. (If the growth of V is more than quadratic, then the leapfrog algorithm as applied with a constant value of h throughout the phase space is in fact unstable whenever the initial condition is large.) The case where V takes negative values but is bounded from below can be reduced to the case $V \geq 0$ by adding a suitable constant to V . In terms of the target measure this just involves changing the normalization constant and hence is irrelevant in the HMC algorithm.

5.1. Preliminaries

Differentiating (3.1) with respect to t , we find successively:

$$\begin{aligned} \dot{p}(t) &= f'(q(t))M^{-1}p(t) , \\ \dot{q}(t) &= M^{-1}f(q(t)) , \\ \ddot{p}(t) &= f^{(2)}(q(t))(M^{-1}p(t), M^{-1}p(t)) + f'(q(t))M^{-1}f(q(t)) , \\ \ddot{q}(t) &= M^{-1}f'(q(t))M^{-1}p(t) , \\ \dddot{p}(t) &= f^{(3)}(q(t))(M^{-1}p(t), M^{-1}p(t), M^{-1}p(t)) + \\ &\quad 3f^{(2)}(q(t))(M^{-1}f(q(t)), M^{-1}p(t)) + f'(q(t))M^{-1}f'(q(t))M^{-1}f(q(t)) , \\ \dddot{q}(t) &= M^{-1}f^{(2)}(q(t))(M^{-1}p(t), M^{-1}p(t)) + M^{-1}f'(q(t))M^{-1}f(q(t)) . \end{aligned}$$

In this section letter K will denote a generic constant which may vary from one appearance to the next, but will depend only on $B, T, \|M\|, \|M^{-1}\|$. From the above equations for the derivatives and using the assumptions on V , we obtain the following bounds:

$$\begin{aligned} |\dot{p}(t)| &\leq |f(q(t))| , & |\dot{q}(t)| &\leq K |p(t)| , \\ |\ddot{p}(t)| &\leq K |p(t)| , & |\ddot{q}(t)| &\leq K |f(q(t))| , \\ |\dddot{p}(t)| &\leq K (|p(t)|^2 + |f(q(t))|) , & |\dddot{q}(t)| &\leq K |p(t)| , \\ |\dddot{p}(t)| &\leq K (|p(t)|^3 + |p(t)||f(q(t))| + |f(q(t))|) , & |\dddot{q}(t)| &\leq K (|p(t)|^2 + |f(q(t))|) . \end{aligned} \tag{5.1}$$

5.2. Asymptotic Expansion for the Leapfrog Solution

In previous sections we have used a subscript to denote the different particles comprising our state space. Here we consider leapfrog integration of a single particle and use the subscript to denote the time-level in this integration. The leapfrog scheme can then be compactly written as

$$q_{n+1} = q_n + hM^{-1}p_n + \frac{h^2}{2}M^{-1}f(q_n) , \tag{5.2}$$

$$p_{n+1} = p_n + \frac{h}{2}f(q_n) + \frac{h}{2}f(q_n + hM^{-1}p_n + \frac{h^2}{2}M^{-1}f(q_n)) . \tag{5.3}$$

We define the truncation error in the usual way:

$$\begin{aligned} -\tau_n^{(q)} &:= q(t_{n+1}) - \left(q(t_n) + hM^{-1}p(t_n) + \frac{h^2}{2}M^{-1}f(q(t_n)) \right), \\ -\tau_n^{(p)} &:= p(t_{n+1}) - \left(p(t_n) + \frac{h}{2}f(q_n) + \frac{h}{2}f(q(t_n) + hM^{-1}p(t_n) + \frac{h^2}{2}M^{-1}f(q(t_n))) \right), \end{aligned}$$

where we have set $t_n = nh \in [0, T]$. Expanding and using standard truncation error analysis (see [19], for example) we obtain:

$$\begin{aligned} \tau_n^{(q)} &= \frac{1}{6}h^3 \ddot{q}(t_n) + h^4 \mathcal{O}(\|\ddot{q}(\cdot)\|_\infty), \\ \tau_n^{(p)} &= -\frac{1}{12}h^3 \ddot{p}(t_n) + h^4 \mathcal{O}(\|\ddot{p}(\cdot)\|_\infty) + h \mathcal{O}(\tau_n^{(q)}), \end{aligned}$$

where, for arbitrary function g :

$$\|g(\cdot)\|_\infty := \sup_{0 \leq t \leq T} |g(t)|.$$

In view of these estimates, $\frac{1}{6}h^3 \ddot{q}(t_n)$ and $-\frac{1}{12}h^3 \ddot{p}(t_n)$ are the leading terms in the asymptotic expansion of the truncation error. Standard results (see, for instance, [20], Section II.8) show that the numerical solution possesses an asymptotic expansion:

$$\begin{aligned} q_n &= q(t_n) + h^2v(t_n) + \mathcal{O}(h^3), \\ p_n &= p(t_n) + h^2u(t_n) + \mathcal{O}(h^3), \end{aligned} \tag{5.4}$$

where functions $u(\cdot)$ and $v(\cdot)$ are the solutions, with initial condition $u(0) = v(0) = 0$, of the *variational* system

$$\begin{pmatrix} \dot{u}(t) \\ \dot{v}(t) \end{pmatrix} = \begin{pmatrix} 0 & M^{-1}f'(q(t)) \\ I & 0 \end{pmatrix} \begin{pmatrix} u(t) \\ v(t) \end{pmatrix} + \begin{pmatrix} \frac{1}{12}\ddot{p}(t) \\ -\frac{1}{6}\ddot{q}(t) \end{pmatrix}. \tag{5.5}$$

Remark 5.2. Notice here that $u(\cdot), v(\cdot)$ depend on the initial conditions $(q(0), p(0))$ via $(q(\cdot), p(\cdot))$ but this dependence is not reflected in the notation. One should keep in mind that most of the norms appearing in the sequel are functions of $(q(0), p(0))$.

Applying Gronwall's lemma and using the estimates (5.1), we obtain the bound:

$$\|u(\cdot)\|_\infty + \|v(\cdot)\|_\infty \leq K (\|p(\cdot)\|_\infty^2 + \|f(q(\cdot))\|_\infty) \tag{5.6}$$

and, by differentiating (5.5) with respect to t , expressing \dot{u}, \dot{v} in terms of u, v , and using (5.1) again, we obtain in turn:

$$\|\ddot{u}(\cdot)\|_\infty \leq K (\|p(\cdot)\|_\infty^3 + \|p(\cdot)\|_\infty \|f(q(\cdot))\|_\infty + \|f(q(\cdot))\|_\infty), \tag{5.7}$$

$$\|\ddot{v}(\cdot)\|_\infty \leq K (\|p(\cdot)\|_\infty^2 + \|f(q(\cdot))\|_\infty). \tag{5.8}$$

5.3. Estimates for the Global Error

With the leading coefficients u, v of the global errors $q_n - q(t_n), p_n - p(t_n)$ estimated in (5.6), our task now is to obtain an explicit bound for the constants implied in the $\mathcal{O}(h^3)$ remainder in (5.4). To this end, we define the quantities

$$\begin{aligned} z_n &:= q(t_n) + h^2 v(t_n) , \\ w_n &:= p(t_n) + h^2 u(t_n) , \end{aligned}$$

and denote by $\tau_n^{(q)*}, \tau_n^{(p)*}$ the residuals they generate when substituted in (5.2), (5.3) respectively, *i.e.*,

$$\begin{aligned} -\tau_n^{(q)*} &= z_{n+1} - z_n - h M^{-1} w_n - \frac{h^2}{2} M^{-1} f(z_n) , \\ -\tau_n^{(p)*} &= w_{n+1} - w_n - \frac{h}{2} f(z_n) - \frac{h}{2} f\left(z_n + h M^{-1} w_n + \frac{h^2}{2} M^{-1} f(z_n)\right) . \end{aligned}$$

Since the leapfrog scheme is stable, standard numerical analysis techniques [20] show that it is possible to estimate the global errors in terms of the local residuals (truncation errors). This gives

$$\max_{0 \leq t_n \leq T} (|q_n - z_n| + |p_n - w_n|) \leq \frac{C}{h} \max_{0 \leq t_n \leq T} (|\tau_n^{(q)*}| + |\tau_n^{(p)*}|) \quad (5.9)$$

with the constant C depending only on T and Lipschitz constant of the map $(q_n, p_n) \mapsto (q_{n+1}, p_{n+1})$, which in turn depends on $\|M^{-1}\|$ and the bound for f' . The stability bound (5.9) is the basis of the proof of the following estimation of the global error:

Proposition 5.3. *If the potential V satisfies Assumptions 5.1, then for $0 \leq t_n \leq T$,*

$$\begin{aligned} |p_n - (p(t_n) + h^2 u(t_n))| &\leq K h^3 (\|p(\cdot)\|_\infty^4 + \|f(q(\cdot))\|_\infty^2 + 1) , \\ |q_n - (q(t_n) + h^2 v(t_n))| &\leq K h^3 (\|p(\cdot)\|_\infty^4 + \|f(q(\cdot))\|_\infty^2 + 1) . \end{aligned}$$

Proof. Our task is reduced to estimating $\tau_n^{(q)*}, \tau_n^{(p)*}$. We only present the estimation for $\tau_n^{(p)*}$, since the computations for $\tau_n^{(q)*}$ are similar but simpler.

After regrouping the terms in $\tau_n^{(p)*}$ we find that

$$\begin{aligned}
-\tau_n^{(p)*} &= \underbrace{p(t_{n+1}) - p(t_n) - \frac{h}{2}f(q(t_n)) - \frac{h}{2}f(q(t_{n+1})) + \frac{h^3}{12}\ddot{p}(t)}_{I_1} \\
&\quad + \underbrace{h^2 \left(u(t_{n+1}) - u(t_n) - hf'(q(t_n))v(t_n) - \frac{h}{12}\ddot{p}(t) \right)}_{I_2} \\
&\quad + \underbrace{\frac{h}{2} \left(f(q(t_n)) - f(z_n) + h^2f'(q(t_n))v(t_n) \right)}_{I_3} \\
&\quad - \underbrace{\frac{h}{2} \left(f(z_n + hM^{-1}w_n + \frac{h^2}{2}M^{-1}f(q(t_n))) - f(q(t_{n+1})) - h^2f'(q(t_n))v(t_n) \right)}_{I_4} \\
&\quad + \underbrace{\frac{h}{2} \left(f(z_n + hM^{-1}w_n + \frac{h^2}{2}M^{-1}f(q(t_n))) - f(z_n + hM^{-1}w_n + \frac{h^2}{2}M^{-1}f(z_n)) \right)}_{I_5}
\end{aligned}$$

Now we estimate the above five terms separately.

I_1 : We note that

$$p(t_{n+1}) - p(t_n) - \frac{h}{2}f(q(t_n)) - \frac{h}{2}f(q(t_{n+1})) = p(t_{n+1}) - p(t_n) - \frac{h}{2}\dot{p}(t_{n+1}) - \frac{h}{2}\dot{p}(t_n).$$

and by using the estimates in (5.1) it follows that

$$|I_1| \leq K h^4 (\|p(\cdot)\|_\infty + \|p(\cdot)\|_\infty \|f(q(\cdot))\|_\infty + \|f(q(\cdot))\|_\infty).$$

I_2 : Here we write $I_2 = h^2(u(t_{n+1}) - u(t_n) - h\dot{u}(t_n))$ so that by (5.7)

$$|I_2| \leq K h^4 (\|p(\cdot)\|_\infty^3 + \|p(\cdot)\|_\infty \|f(q(\cdot))\|_\infty + \|f(q(\cdot))\|_\infty).$$

I_3 : This term is estimated, after Taylor expanding $f(z_n)$ near $f(q(t_n))$, by

$$|I_3| \leq K h^5 (\|p(\cdot)\|_\infty + \|f(q(\cdot))\|_\infty)^2.$$

I_4 : We rewrite this as

$$\frac{h}{2} \left(f(q(t_{n+1}) + \tau_n^{(q)} + h^2v(t_n) + h^3M^{-1}v(t_n)) - f(q(t_{n+1})) - h^2f'(q(t_n))v(t_n) \right)$$

and Taylor expand around $f(q(t_n))$ to derive the bound:

$$|I_4| \leq K h^4 (\|p(\cdot)\|_\infty^4 + \|f(q(\cdot))\|_\infty^2).$$

I_5 : This term is easily estimated as:

$$|I_5| \leq K h^5 \|v(\cdot)\|_\infty \leq K h^5 (\|p(\cdot)\|_\infty^2 + \|f(q(\cdot))\|_\infty).$$

Combining all the above estimates, we have the bound

$$|\tau_n^{(p)*}| \leq K h^4 (\|p(\cdot)\|_\infty^4 + \|f(q(\cdot))\|_\infty^2).$$

A similar analysis for $\tau_n^{(q)*}$ yields the bound

$$|\tau_n^{(q)*}| \leq K h^4 (\|p(\cdot)\|_\infty^4 + \|f(q(\cdot))\|_\infty^2) .$$

The proof is completed by substituting the above estimates in (5.9). \square

We now use the estimates in Proposition 5.3 to derive the asymptotic expansion for the energy increment for the leapfrog scheme (cf. Condition 1).

Proposition 5.4. *Let potential V satisfy Assumptions 5.1. Then, for the leapfrog scheme, we get*

$$\Delta(x, h) = h^2 \alpha(x) + h^2 \rho(x, h) ,$$

with

$$\begin{aligned} \alpha(x) &= \langle M^{-1}p(T), u(T) \rangle - \langle f(q(T)), v(T) \rangle , \\ |\alpha(x)| &\leq K (\|p(\cdot)\|_\infty^3 + \|f(q(\cdot))\|_\infty^2 + 1) , \\ |\rho(x, h)| &\leq K h (\|p(\cdot)\|_\infty^8 + \|f(q(\cdot))\|_\infty^2 + 1) , \quad 0 < h \leq 1 , \end{aligned}$$

where $(q(\cdot), p(\cdot))$ denotes the solution of (3.1) with initial data $x \equiv (q(0), p(0))$ and $u(\cdot), v(\cdot)$ are the solutions of the corresponding variational system given in (5.5) with $u(0) = v(0) = 0$.

Proof. We only consider the case when T/h is an integer. The general case follows with minor adjustments. By Proposition 5.3,

$$\begin{aligned} \Delta(x, h) &= H(\psi_h^{(T)}(x)) - H(x) = H(\psi_h^{(T)}(x)) - H(\varphi_T(x)) = \\ &= \langle M^{-1}p(T), h^2 u(T) + h^3 R_1 \rangle + \frac{1}{2} \langle M^{-1}(h^2 u(T) + h^3 R_1), (h^2 u(T) + h^3 R_1) \rangle \\ &\quad + V(q(T) + h^2 v(T) + h^3 R_2) - V(q(T)) , \end{aligned}$$

where R_1, R_2 are remainders with

$$|R_1| + |R_2| \leq K (\|p(\cdot)\|_\infty^4 + \|f(q(\cdot))\|_\infty^2 + 1) .$$

By Taylor expanding $V(\cdot)$ around $q(T)$ we obtain,

$$\Delta(x, h) = h^2 (\langle M^{-1}p(T), u(T) \rangle - \langle f(q(T)), v(T) \rangle) + \rho(x, h) ,$$

with

$$|\rho(x, h)| \leq K h^3 (\|p(\cdot)\|_\infty^8 + \|f(q(\cdot))\|_\infty^2 + 1)$$

for $0 \leq h \leq 1$. From the bound (5.6) it follows that

$$\begin{aligned} |\alpha(x)| &\leq K (\|p(\cdot)\|_\infty \|u(\cdot)\|_\infty + \|f(q(\cdot))\|_\infty \|v(\cdot)\|_\infty) \\ &\leq K (\|p(\cdot)\|_\infty^3 + \|f(\cdot)\|_\infty^2 + 1) \end{aligned}$$

and the theorem is proved. \square

Our analysis is completed by estimating the quantities $\|p(\cdot)\|_\infty$ and $\|q(\cdot)\|_\infty$, that feature in the preceding theorems, in terms of the initial data $(q(0), p(0))$. We obtain these estimates for two families of potentials which include most of the interesting/useful target distributions. The corresponding estimates for other potentials may be obtained using similar methods.

Proposition 5.5. *Let potential V satisfy Assumptions 5.1. If V satisfies, in addition, either of the following conditions:*

(1) *f is bounded and*

$$\int_{\mathbb{R}^m} |V(q)|^8 e^{-V(q)} dq < \infty ; \quad (5.10)$$

(2) *there exist constants $C_1, C_2 > 0$ and $0 < \gamma \leq 1$ such that for all $|q| \geq C_2$, we have $V(q) \geq C_1|q|^\gamma$;*

then Conditions 3.1, 3.2 and 3.7 all hold.

Proof. We only present the treatment of Conditions 3.1 and 3.2. The derivation of Condition 3.7 is similar and simpler.

From Proposition 5.4 we observe that function $D(x)$ in Condition 3.2 may be taken to be

$$D(x) = K \left(\|p(\cdot)\|_\infty^{16} + \|f(q(\cdot))\|_\infty^4 + 1 \right) .$$

Thus, to prove integrability of $D(\cdot)$ we need to estimate $\|p(\cdot)\|_\infty$ and $\|f(q(\cdot))\|_\infty$. Estimating $\|p(\cdot)\|_\infty$ is easier. Indeed, by conservation of energy,

$$\frac{1}{2} \langle p(t), M^{-1}p(t) \rangle \leq \frac{1}{2} \langle p(0), M^{-1}p(0) \rangle + V(q(0)) ,$$

which implies

$$|p(t)|^{16} \leq K \left(|p(0)|^{16} + |V(q(0))|^8 \right) . \quad (5.11)$$

Now, we prove integrability of $D(\cdot)$ under each of the two stated hypothesis.

Under hypothesis (1): Suppose f is bounded. In this case we obtain that $|D(x)| \leq K(\|p(\cdot)\|_\infty^{16} + 1)$, therefore it is enough to estimate $\|p(\cdot)\|_\infty$. Since the Gaussian distribution has all moments, integrability of D follows from (5.10) and (5.11).

Under hypothesis (2): Using the stated hypothesis on $V(q)$ we obtain

$$C_1|q(t)|^\gamma \leq V(q(t)) \leq \frac{1}{2} \langle p(0), M^{-1}p(0) \rangle + V(q(0)) ,$$

which implies that:

$$|q(t)| \leq K \left(|p(0)|^{\frac{2}{\gamma}} + |V(q(0))|^{\frac{1}{\gamma}} \right) .$$

By Assumptions 5.1(i), $|f(q(t))| \leq K(1 + |q(t)|)$ and arguing as above and using the bound (5.11), integrability of D follows if we show that

$$\int_{\mathbb{R}^m} |V(q)|^\delta e^{-V(q)} dq < \infty , \quad \delta = \max\left(8, \frac{4}{\gamma}\right) .$$

Since $|V(q)| \leq K(1 + |q|^2)$,

$$\int_{\mathbb{R}^m} |V(q)|^\delta e^{-V(q)} dq \leq K \int_{\mathbb{R}^m} (1 + |q|^{2\delta}) e^{-B|q|^\gamma} dq < \infty$$

and we are done. \square

6. Proofs of Probabilistic Results

Proof of Lemma 3.3. The volume preservation property of $\psi_h^{(T)}(\cdot)$ implies that the associated Jacobian is unit. Thus, setting $x = (\psi_h^{(T)})^{-1}(y)$ we get:

$$\begin{aligned} \int_{\mathbb{R}^{2m}} g(\Delta(x, h)) e^{-H(x)} dx &= \int_{\mathbb{R}^{2m}} g(H(\psi_h^{(T)}(x)) - H(x)) e^{-H(x)} dx \\ &= \int_{\mathbb{R}^{2m}} g[H(y) - H((\psi_h^{(T)})^{-1}(y))] e^{-H((\psi_h^{(T)})^{-1}(y))} dy . \end{aligned}$$

Following the definition of time reversibility in (2.2), we have:

$$S \circ \psi_h^{(T)} = (\psi_h^{(T)})^{-1} \circ S$$

for the symmetry operator S such that $S(q, p) = (q, -p)$. Using now the volume preserving transformation $y = Sz$ and continuing from above, we get:

$$\begin{aligned} \int_{\mathbb{R}^{2m}} g(\Delta(x, h)) e^{-H(x)} dx &= \int_{\mathbb{R}^{2m}} g(H(Sz) - H((\psi_h^{(T)})^{-1}(Sz))) e^{-H((\psi_h^{(T)})^{-1}(Sz))} dz \\ &= \int_{\mathbb{R}^{2m}} g(H(Sz) - H(S\psi_h^{(T)}(z))) e^{-H(S\psi_h^{(T)}(z))} dz \\ &= \int_{\mathbb{R}^{2m}} g(H(z) - H(\psi_h^{(T)}(z))) e^{-H(\psi_h^{(T)}(z))} dz , \end{aligned}$$

where in the last equation we have used the identity $H(Sz) = H(z)$. \square

Proof of Proposition 3.4. We will first find the limit of $\sigma^2(h)/h^4$. Conditions 3.1 and 3.2 imply that:

$$\frac{\Delta^2(x, h)}{h^4} = \alpha^2(x) + \rho^2(x, h) + 2\rho(x, h)\alpha(x) \leq D(x)$$

and since, for fixed x , $\Delta^2(x, h)/h^4 \rightarrow \alpha^2(x)$, the dominated convergence theorem shows:

$$\lim_{h \rightarrow 0} \frac{s^2(h)}{h^4} = \int_{\mathbb{R}^{2m}} \alpha^2(x) e^{-H(x)} dx = \Sigma .$$

Now, (3.6) implies that:

$$\lim_{h \rightarrow 0} \frac{\mu^2(h)}{h^4} = 0, \quad (6.1)$$

and the required limit for $\sigma^2(h)/h^4$ follows directly. Then, from (3.5) we obtain

$$\begin{aligned} \frac{2\mu(h) - \sigma^2(h)}{h^4} = & \\ & - \int_{\mathbb{R}^{2m}} \frac{\Delta(x, h)}{h^2} \frac{[\exp(-\Delta(x, h)) - 1 + \Delta(x, h)]}{h^2} e^{-H(x)} dx + \frac{\mu^2(h)}{h^4}. \end{aligned}$$

Since for any fixed x , Conditions 3.1 and 3.2 imply that $\Delta(x, h) \rightarrow 0$ as $h \rightarrow 0$ and $\Delta^2(x, h) = \mathcal{O}(h^4)$, we have the pointwise limit

$$\lim_{h \rightarrow 0} \frac{\exp(-\Delta(x, h)) - 1 + \Delta(x, h)}{h^2} = 0.$$

Using inequality $|u|e^u - 1 - u \leq |u|^2(e^u + 2)$, we deduce that for all sufficiently small h ,

$$\begin{aligned} & \int_{\mathbb{R}^{2m}} \frac{|\Delta(x, h)|}{h^2} \frac{|\exp(-\Delta(x, h)) - 1 + \Delta(x, h)|}{h^2} e^{-H(x)} dx \\ & \leq \int_{\mathbb{R}^{2m}} \frac{|\Delta^2(x, h)|}{h^4} \exp(-\Delta(x, h)) e^{-H(x)} dx + 2 \int_{\mathbb{R}^{2m}} \frac{|\Delta^2(x, h)|}{h^4} e^{-H(x)} dx \\ & \leq 3 \int_{\mathbb{R}^{2m}} D(x) e^{-H(x)} dx < \infty, \end{aligned}$$

where the last line follows from applying Lemma 3.3 with $\varphi(x) = x^2$ and Condition 3.2. So, the dominated convergence theorem yields

$$\lim_{h \rightarrow 0} \frac{2\mu(h) - \sigma^2(h)}{h^4} = 0.$$

This completes the proof of the proposition. \square

Proof of Theorem 3.6. We continue from (3.8). In view of the scaling $h = l \cdot d^{-1/4}$ we obtain, after using Proposition 3.4:

$$\mathbb{E}[R_d] = -d \cdot \mu(h) \rightarrow -\frac{l^4 \sigma}{2}$$

and

$$\text{Var}[R_d] = d \cdot \sigma^2(h) \rightarrow l^4 \Sigma.$$

The Lindeberg condition is easily seen to hold and therefore:

$$R_d \xrightarrow{\mathcal{L}} R_\infty := N\left(-\frac{l^4 \Sigma}{2}, l^4 \Sigma\right).$$

From the boundedness of $u \mapsto 1 \wedge e^u$ we may write:

$$\mathbb{E}[a(X, Y)] \rightarrow \mathbb{E}[1 \wedge e^{R_\infty}] ,$$

where the last expectation can be found analytically (see e.g. [36]) to be:

$$\mathbb{E}[1 \wedge e^{R_\infty}] = 2\Phi(-l^2\sqrt{\Sigma}/2) .$$

This completes the proof. \square

Proof of Proposition 3.8. For simplicity, we will write just q^n , q^{n+1} and p^n instead of q_1^n , q_1^{n+1} , p_1^n respectively. Using (3.9), we get:

$$(q^{n+1} - q^n)^2 = I^n (\mathcal{P}_q \psi_h^{(T)}(q^n, p^n) - q^n)^2 .$$

We define:

$$a^-(X^n, Y^n) := 1 \wedge \exp\left\{-\sum_{i=2}^d \Delta(x_i^n, h)\right\}; \quad I^{n-} := \mathbb{1}_{U^n < a^-(X^n, Y^n)} , \quad (6.2)$$

and set

$$\xi^n = I^{n-} (\mathcal{P}_q \psi_h^{(T)}(q^n, p^n) - q^n)^2 .$$

Using the Lipschitz continuity of $u \mapsto \mathbb{1}_{U \leq 1 \wedge e^u}$ and the Cauchy-Schwartz inequality we get:

$$\mathbb{E}|(q^{n+1} - q^n)^2 - \xi^n| \leq |\Delta(x_1, h)|_{L_2} |(\mathcal{P}_q \psi_h^{(T)}(q^n, p^n) - q^n)^2|_{L_2}$$

Now, Conditions 3.1 and 3.2 imply that

$$|\Delta(x_1, h)|_{L_2} = \mathcal{O}(h^2) .$$

Also, from Condition 3.7 and the stated hypothesis on the density $\exp(-V)$, q^n and $\mathcal{P}_q \psi_h^{(T)}(q^n, p^n)$ have bounded fourth moments uniformly in h , so:

$$|(\mathcal{P}_q \psi_h^{(T)}(q^n, p^n) - q^n)^2|_{L_2} \leq C ,$$

for some constant $C > 0$. The last two statements imply that:

$$\mathbb{E}|(q^{n+1} - q^n)^2 - \xi^n| = \mathcal{O}(h^2) . \quad (6.3)$$

Exploiting the independence between I^{n-} and the first particle:

$$\mathbb{E}[\xi_n] = \mathbb{E}[a^-(X, Y)] \times \mathbb{E}[(\mathcal{P}_q \psi_h^{(T)}(q^n, p^n) - q^n)^2] \rightarrow a(l) \cdot \mathbb{E}[(\mathcal{P}_q \varphi_T(q^n, p^n) - q^n)^2] ,$$

where, for the first factor we used its limit from Theorem 3.6; for the second factor the limit is a consequence by Condition 3 and the dominated convergence theorem. Equation (6.3) completes the proof. \square

Proof of Proposition 3.9. Fix some $q_1^n \in \mathbb{R}^m$. We define $a^-(X^n, Y^n)$ and I^{n-} as in (6.2). For simplicity, we will write just q^n , q^{n+1} , \mathbf{q}^{n+1} and p^n instead of q_1^n , q_1^{n+1} , \mathbf{q}_1^{n+1} and p_1^n respectively.

We set

$$g^{n+1} = I^{n-} \cdot \mathcal{P}_q \varphi_T(q^n, p^n) + (1 - I^{n-}) q^n .$$

Adding and subtracting $I^n \cdot \mathcal{P}_q(\varphi_T(q^n, p^n))$ yields:

$$\begin{aligned} |q^{n+1} - g^{n+1}| &\leq |\mathcal{P}_q(\psi_h^{(T)}(q^n, p^n)) - \mathcal{P}_q(\varphi_T(q^n, p^n))| \\ &\quad + |I^{n-} - I^n| (|\mathcal{P}_q(\varphi_T(q^n, p^n))| + |q^n|) . \end{aligned} \quad (6.4)$$

Using the Lipschitz continuity (with constant 1) of $u \mapsto \mathbb{I}_{U \leq 1 \wedge \exp(u)}$:

$$|I^{n-} - I^n| \leq |\Delta(x_1, h)| . \quad (6.5)$$

Now, Condition 3.7 implies that the first term on the right-hand side of (6.4) vanishes w.p.1 and Condition 3.1 implies (via (6.5)) that also the second term vanishes w.p.1. Therefore, as $d \rightarrow \infty$:

$$q^{n+1} - g^{n+1} \rightarrow 0, \text{ a.s. .}$$

Theorem 3.6 immediately implies that $I^{n-} \xrightarrow{\mathcal{L}} I^n$, thus:

$$g^{n+1} \xrightarrow{\mathcal{L}} \mathbf{q}^{n+1} .$$

From these two limits, we have $q^{n+1} \xrightarrow{\mathcal{L}} \mathbf{q}^{n+1}$, and this completes the proof. \square

Proof of Theorem 4.2. To simplify the notation we again drop the subscript 1. Conditionally on the trajectory q^0, q^1, \dots we get:

$$(q(t+\delta) - q(t))^2 = \begin{cases} 0, & \text{w.p. } 1 - \lambda_d \delta + \mathcal{O}((\lambda_d \delta)^2), \\ (q^{N(t)+1} - q^{N(t)})^2, & \text{w.p. } \lambda_d \delta + \mathcal{O}((\lambda_d \delta)^2), \\ (q^{N(t)+1+j} - q^{N(t)})^2, j \geq 1, & \text{w.p. } \mathcal{O}((\lambda_d \delta)^{j+1}). \end{cases}$$

Therefore,

$$\begin{aligned} \mathcal{SJD}_d &= \mathbb{E}[(q^{N(t)+1} - q^{N(t)})^2] (\lambda_d \delta + \mathcal{O}((\lambda_d \delta)^2)) \\ &\quad + \sum_{j \geq 1} \mathbb{E}[(q^{N(t)+1+j} - q^{N(t)})^2] \mathcal{O}((\lambda_d \delta)^{j+1}) . \end{aligned} \quad (6.6)$$

Note now that:

$$\begin{aligned} \mathbb{E}[(q^{N(t)+1+j} - q^{N(t)})^2] &\leq \left(\sum_{k=1}^{j+1} |q^{N(t)+k} - q^{N(t)+k-1}|_{L_2} \right)^2 \\ &= (j+1)^2 \mathbb{E}[(q^{n+1} - q^n)^2] , \end{aligned}$$

since we have assumed stationarity. From (4.1):

$$\lambda_d = d^{-5/4} \frac{l}{TC_{LF}} + \mathcal{O}(d^{-6/4}) .$$

and, from Proposition 3.8, $\mathbb{E}[(q^{n+1} - q^n)^2] = \mathcal{O}(1)$. Therefore,

$$d^{5/4} \times \sum_{j \geq 1} \mathbb{E}[(q^{N(t)+1+j} - q^{N(t)})^2] \mathcal{O}((\lambda_d \delta)^{j+1})$$

is of the same order in d as

$$\lambda_d^2 \cdot d^{5/4} \times \sum_{j \geq 1} (j+1)^2 \mathcal{O}(\lambda_d^{j-1}) ,$$

thus:

$$d^{5/4} \times \sum_{j \geq 1} \mathbb{E}[(q^{N(t)+1+j} - q^{N(t)})^2] \mathcal{O}((\lambda_d \delta)^{j+1}) = \mathcal{O}(\lambda_d) .$$

Using this result, and continuing from (6.6), Proposition 3.8 provides the required statement. \square

7. Conclusions

The HMC methodology provides a promising framework for the study of sampling problems, especially in high dimensions. There are a number of directions in which the research direction taken in this paper could be developed further, and a number of observations to be made concerning optimal tuning of MCMC methods in general. We conclude by listing some of these issues.

- The overall optimization involves tuning *three* free parameters (h, T, M); since M is a symmetric matrix, the number of tuning parameters is $2 + m(m+1)/2$. In this paper, we have fixed M and T and illustrated that the choice $h = l d^{-1/4}$ provides non-vanishing $\mathcal{O}(1)$ acceptance probabilities as $d \rightarrow \infty$. We then focussed on optimizing the HMC algorithm over choice of l . The natural next step would be to study the algorithm for various choices of the mass matrix M and the integration time T .
- There is interesting recent computational work [16] concerning exploration of state space by means of nonseparable Hamiltonian dynamics; this work opens up several theoretical research directions.
- The issue of irreducibility for the transition kernel of HMC is subtle, and requires further investigation, as certain exceptional cases can lead to nonergodic behaviour (see [11, 42] and the references therein).
- Our analysis of the HMC algorithm is conducted in stationarity. It is possible that different scaling analyses will be required to study the burn-in phase of the algorithm, as for the study of random walk type algorithms in [13].

- There is evidence that the limiting properties of MALA for high-dimensional target densities do not appear to depend critically on the tail behaviour of the target (see [37]). However in the present paper for HMC, we have considered densities that are no lighter than Gaussian at infinity. It would thus be interesting to extend the work to light-tailed densities. This links naturally to the question of using variable step size integration [41] for HMC since light tailed densities will lead to superlinear vector fields at infinity in (2.1). This also links to the work in [16] where non-separable Hamiltonians arise via introduction of a non-standard metric on phase space, related to the Fisher information. This metric introduces a rescaling of state space and this rescaling induces similar algorithmic properties to those induced by variable time-stepping.
- We have shown how to scale the HMC method to obtain $\mathcal{O}(1)$ acceptance probabilities as the dimension of the target product measure grows. We have also shown how to minimize a reasonable measure of computational cost, defined as the work needed to make an $\mathcal{O}(1)$ move in state space. However, in contrast to similar work for RWM and MALA ([36, 37]) where a scalar SDE governs, for large d , the evolution of a single component of the Markov chain, we have not identified a limiting Markov process which arises in the infinite dimensional limit of HMC. This remains an interesting and technically demanding challenge.
- The work concerning optimal scaling of RWM and MALA in [36, 37] and the identification of the optimal acceptance probabilities of 0.234 and 0.574 respectively, concerns target measures with an iid structure. However, recent work [27, 35] shows that, for measures which have density with respect to a Gaussian measure (in the limit $d \rightarrow \infty$), hence are not necessarily iid, the same optimal acceptance probabilities arise. It would be natural to try to extend the work concerning HMC contained in this paper to non iid target measures in a similar manner. Note also that for measures with this special structure these results on optimal scaling are mainly of *theoretical* interest because they extend known results out of the iid scenario. For the particular case of measures which have density with respect to a Gaussian, and from a more *practical* perspective, the RWM, MALA and HMC algorithms should all be modified to exploit this underlying Gaussian structure. This is the subject of the next bullet.
- We have concentrated on explicit integration by the leapfrog method. For measures which have density with respect to a Gaussian measure (in the limit $d \rightarrow \infty$) it is natural to use semi-implicit integrators which compute the linear dynamics implicitly, leading to exact statistics in the pure Gaussian case. This idea was first developed for the MALA algorithm [8] and for the RWM algorithm in [10] and leads to methods which explore state space in $\mathcal{O}(1)$ steps for measures with this special structure. The idea has recently been developed for HMC methods in [7] and the resulting algorithm shown to outperform the semi-implicit MALA algorithm for some problems arising in conditioned diffusions. Developing a theoretical understanding of this behaviour would be of interest. Note that the optimal acceptance probabilities 0.234, 0.574 and 0.651 will not necessarily apply for these semi-implicit proposals as the optimal proposal variance does not shrink to zero

as $d \rightarrow \infty$; as a result different mechanisms may come into play when determining optimality.

- It would be interesting to conduct simulation studies which investigate the robustness of optimal scaling results for RWM, MALA and HMC in scenarios in which the target is not iid or change of measure from Gaussian. Such simulation studies could help guide future theoretical results on optimal scaling.

Acknowledgements

We thank Sebastian Reich for drawing our attention to the paper [18] which sparked our initial interest in the scaling issue for HMC. Further thanks also to Gabriel Stoltz and Robert D. Skeel for stimulating discussions and useful comments. NP gratefully acknowledges the NSF grant DMS 1107070; JMS gratefully acknowledges the grant TM2010-18246-C03 by Ministerio de Ciencia e Innovacion, Spain; AS is grateful to EPSRC and ERC for financial support. Part of this work was done when NP was a postdoctoral member of CRiSM, University of Warwick, and visited JMS at University of Valladolid, so we thank both these institutions for their warm hospitality. Finally, we thank two referees for their comments that greatly improved the content and presentation of the paper.

References

- [1] E. Akhmatskaya, N. Bou-Rabee, and S. Reich. A comparison of generalized Hybrid Monte Carlo methods with and without momentum flip. *J. Comput. Phys.*, 228(6), 2009.
- [2] F.J. Alexander, G.L. Eyink, and J.M. Restrepo. Accelerated Monte Carlo for optimal estimation of time series. *J. Stat. Phys.*, 119(5-6):1331–1345, 2005.
- [3] A. D. Barbour and Louis H. Y. Chen, editors. *An introduction to Stein's method*, volume 4 of *Lecture Notes Series. Institute for Mathematical Sciences. National University of Singapore*. Singapore University Press, Singapore, 2005.
- [4] M. Bédard. Weak convergence of Metropolis algorithms for non-iid target distributions. *The Annals of Applied Probability*, 17(4):1222–1244, 2007.
- [5] M. Bedard. Efficient sampling using Metropolis algorithms: Applications of optimal scaling results. *Journal of Computational and Graphical Statistics*, 17(2):312–332, 2008.
- [6] M. Bédard. Optimal acceptance rates for Metropolis algorithms: Moving beyond 0.234. *Stochastic Processes and their Applications*, 118(12):2198–2222, 2008.
- [7] A. Beskos, F. Pinski, J.M. Sanz-Serna, and A.M. Stuart. Hybrid Monte-Carlo on Hilbert spaces. *Stochastic Processes and Applications*, 121:2201–2230, 2011.
- [8] A. Beskos, G. Roberts, A.M. Stuart, and J. Voss. An MCMC method for diffusion bridges. *Stochastics and Dynamics*, 8:319–350, 2008.
- [9] A. Beskos, G.O. Roberts, and A.M. Stuart. Optimal scalings for local Metropolis-Hastings chains on non-product targets in high dimensions. *Ann. Appl. Probab.*, 19(3):863–898, 2009.

- [10] A. Beskos and A.M. Stuart. MCMC methods for sampling function space. *Proceedings of the International Congress of Industrial and Applied Mathematicians*, 2009.
- [11] E. Cancès, F. Legoll, and G. Stoltz. Theoretical and numerical comparison of some sampling methods for molecular dynamics. *M2AN Math. Model. Numer. Anal.*, 41(2):351–389, 2007.
- [12] L. Chen, Z. Qin, and J. Liu. Exploring Hybrid Monte Carlo in Bayesian computation. In *Proceedings, ISBA*, 2000.
- [13] O.F. Christensen, G.O. Roberts, and J.S. Rosenthal. Scaling limits for the transient phase of local Metropolis–Hastings algorithms. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):253–268, 2005.
- [14] P. Diaconis, S. Holmes, and R.M. Neal. Analysis of a nonreversible Markov chain sampler. *Annals of Applied Probability*, 10(3):726–752, 2000.
- [15] S. Duane, A.D. Kennedy, B. Pendleton, and D. Roweth. Hybrid Monte Carlo. *Phys. Lett. B*, 195(2):216–222, 1987.
- [16] M. Girolami and B. Calderhead. Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *JRSS - series B*, 73, 2011.
- [17] R. Gupta, G.W. Kilcup, and S.R. Sharpe. Tuning the Hybrid Monte Carlo algorithm. *Phys. Rev. D*, 38(4):1278–1287, 1988.
- [18] S. Gupta, A. Irbäck, F. Karsch, and B. Petersson. The acceptance probability in the Hybrid Monte Carlo method. *Phys. Lett. B*, 242(3-4):437–443, 1990.
- [19] Ernst Hairer, Christian Lubich, and Gerhard Wanner. *Geometric numerical integration*, volume 31 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, second edition, 2006. Structure-preserving algorithms for ordinary differential equations.
- [20] Ernst Hairer, Syvert P. Nørsett, and Gerhard Wanner. *Solving ordinary differential equations. I*, volume 8 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 1987. Nonstiff problems.
- [21] P. Hall and C. C. Heyde. *Martingale limit theory and its application*. Academic Press Inc. [Harcourt Brace Jovanovich Publishers], New York, 1980. Probability and Mathematical Statistics.
- [22] U. H. E. Hansmann, Y. Okamoto, and F. Eisenmenger. Molecular dynamics, Langevin and Hybrid Monte Carlo simulations in a multicanonical ensemble. *Chemical Physics Letters*, 259(3-4):321 – 330, 1996.
- [23] M. Hasenbusch. Speeding up the Hybrid Monte Carlo algorithm for dynamical fermions. *Phys. Lett. B*, 519(1-2):177–182, 2001.
- [24] J.A. Izaguirre and S.S. Hampton. Shadow hybrid Monte Carlo: an efficient propagator in phase space of macromolecules. *J. Comp. Phys.*, 200:581–604, 2004.
- [25] Benedict Leimkuhler and Sebastian Reich. *Simulating Hamiltonian dynamics*, volume 14 of *Cambridge Monographs on Applied and Computational Mathematics*. Cambridge University Press, Cambridge, 2004.
- [26] J.S. Liu. *Monte Carlo strategies in scientific computing*. Springer Verlag, 2008.
- [27] J.C. Mattingly, N. Pillai, and A.M. Stuart. Diffusion limits of random walk Metropolis algorithms in high dimensions. *Annals of Applied Probability*, 2012.

- [28] B. Mehlig, D.W. Heermann, and B.M. Forrest. Exact Langevin algorithms. *Molecular Physics*, 76(6):1347–1357, 1992.
- [29] N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and E. Teller. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6):1087–1092, 1953.
- [30] L. Mohamed, M. Christie, and V. Demyanov. Comparison of stochastic sampling algorithms for uncertainty quantification. Technical report, Institute of Petroleum Engineering, Heriot-Watt University, Edinburgh, 2009. SPE Reservoir Simulation Symposium.
- [31] R.M. Neal. Probabilistic inference using Markov chain Monte Carlo methods. Technical report, Department of Computer Science, University of Toronto, 1993.
- [32] R.M. Neal. *Bayesian learning for Neural networks*. Springer-Verlag, 1996.
- [33] C.S. Pangali, M. Rao, and B.J. Berne. On a novel Monte Carlo scheme for simulating water and aqueous solutions. *Chem. Phys. Lett*, 55:413–417, 1978.
- [34] C. Pasarica and A. Gelman. Adaptively scaling the Metropolis algorithm using expected squared jumped distance. *Statistica Sinica*, 20:343–364, 2010.
- [35] N.S. Pillai, A.M. Stuart, and A.H. Thiery. Optimal scaling and diffusion limits for the Langevin algorithm in high dimensions. Submitted.
- [36] G. O. Roberts, A. Gelman, and W. R. Gilks. Weak convergence and optimal scaling of random walk Metropolis algorithms. *Ann. Appl. Probab.*, 7(1):110–120, 1997.
- [37] G.O. Roberts and J.S. Rosenthal. Optimal scaling of discrete approximations to Langevin diffusions. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 60(1):255–268, 1998.
- [38] G.O. Roberts and J.S. Rosenthal. Optimal scaling for various Metropolis-Hastings algorithms. *Statist. Sci.*, 16(4):351–367, 2001.
- [39] G.O. Roberts and R.L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, 2(4):341–363, 1996.
- [40] P.J. Rossky, Doll. J.D., and H.L. Friedman. Brownian dynamics as smart Monte Carlo simulation. *J. Chem. Phys.*, 69:4628–4633, 1978.
- [41] Jesus-Maria Sanz-Serna and Mari-Paz Calvo. *Numerical Hamiltonian problems*, volume 7 of *Applied Mathematics and Mathematical Computation*. Chapman & Hall, London, 1994.
- [42] Christof Schütte. Conformational dynamics: Modelling, theory, algorithm, and application of biomolecules. 1998. Habilitation Thesis, Dept. of Mathematics and Computer Science, Free University Berlin, Available at <http://proteomics-berlin.de/89/>.
- [43] J. C. Sexton and D. H. Weingarten. Hamiltonian evolution for the Hybrid Monte Carlo algorithm. *Nuclear Physics B*, 380(3):665–677, 1992.
- [44] R.D. Skeel. Integration schemes for molecular dynamics and related applications. In *The graduate student’s guide to numerical analysis ’98 (Leicester)*, volume 26 of *Springer Ser. Comput. Math.*, pages 119–176. Springer, Berlin, 1999.
- [45] M.E. Tuckerman, B.J. Berne, G.J. Martyna, and M.L. Klein. Efficient molecular dynamics and hybrid Monte Carlo algorithms for path integrals. *The Journal of Chemical Physics*, 99(4):2796–2808, 1993.
- [46] M. Zlochín and Y. Baram. Manifold stochastic dynamics for Bayesian learning. *Neural Computation*, 13(11):2549–2572, 2001.