

# Adjustable Randomized Time Integrators for Hamiltonian Monte Carlo

Yunaf Bor-Rabee

Rutgers University Camden

joint with Milo Marsden @ Stanford

## 1. MCMC

Aims to approximately sample from a target probability measure  $\mu$  by simulating a Markov chain

$$X_0, X_1, \dots \quad \text{s.t.} \quad \text{Law}(X_n) \approx \mu.$$

How many MCMC steps till  $\text{Law}(X_n) \approx \mu$ ?

## 2. Hamiltonian MCMC

Hamiltonian Monte Carlo (HMC) is a class of MCMC methods for sampling

possibly rough.

$$\mu(dx) \propto \exp(-U(x)) \lambda^d(dx)$$

where  $U: \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable.

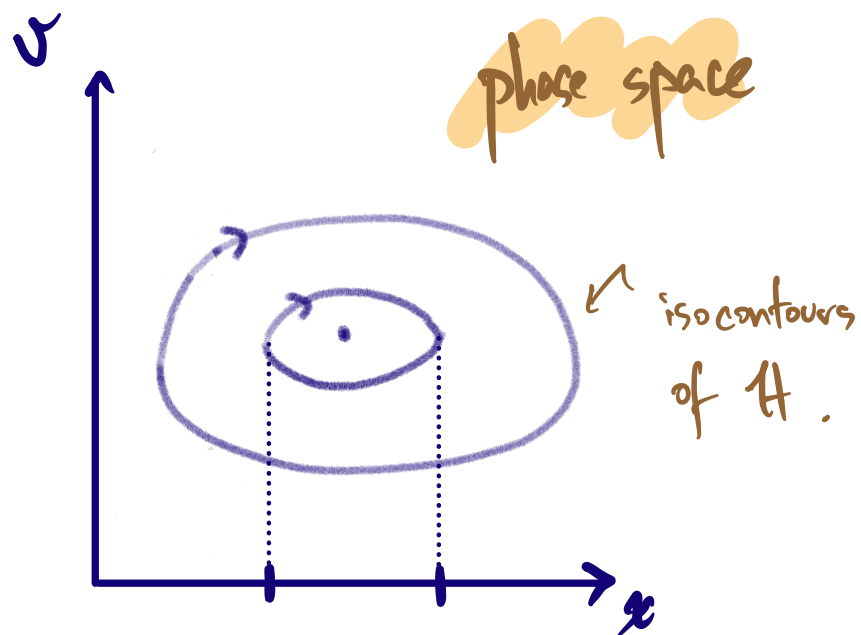
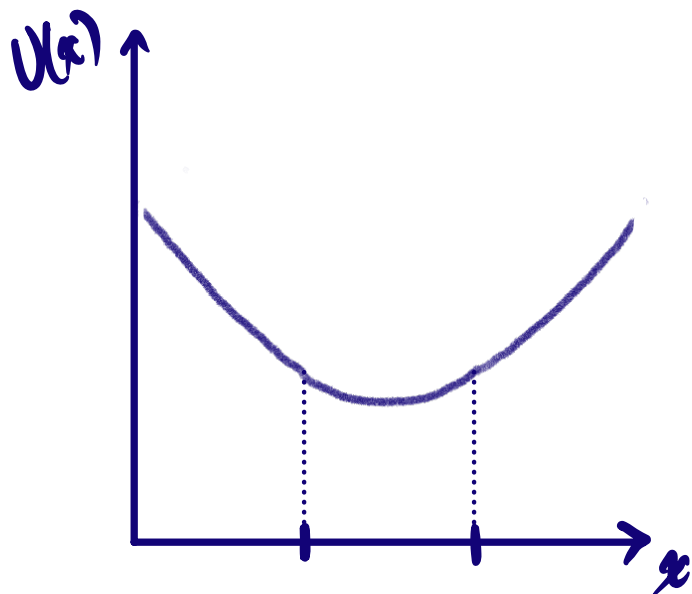
Motivated in part by success of momentum-based

algorithms for optimization: Polyak and Nesterov's methods.

HMC constructs a Markov chain "aimed" at

$$\mu_{\text{eq}}(dx dv) \propto e^{-H(x,v)} \lambda^d(dx) \lambda^d(dv)$$

where  $H(x,v) = \frac{1}{2}|v|^2 + U(x)$ .





HMC uses a fictitious Hamiltonian dynamics

$$\dot{q}_t = v_t \quad \dot{v}_t = -\nabla U(q_t)$$

Let  $\phi_t: \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$  be the corresponding flow.

energy conservation

$$(1) \quad H \circ \phi_t = H$$

volume conservation

$$(2) \quad \lambda^{2d} \circ \phi_t = \lambda^{2d}$$

However, Hamiltonian dynamics, by itself, preserves infinitely many invariant measures.

## Exact HMC Algorithm

Input : duration hyperparameter  $T > 0$

Output :  $X_0, X_1, \dots$

(Step 1) Draw  $\xi_k \sim \mathcal{N}(0, I_d)$ .

velocity randomization

(Step 2)  $X_{k+1} \leftarrow q_T(X_k, \xi_k)$

Transition Kernel  $\mathbb{T}_{\text{ex}}(x, \cdot) = \mathbb{E}[\delta_{q_T(x, \xi)}] \quad \xi \sim \mathcal{N}(0, I_d)$

Theorem (Chen and Vempala 2019)  $\leftarrow$  builds on Mangoubi / Smith 2017

Suppose  $U: \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies

$$(A1) \quad |\nabla U(x) - \nabla U(y)| \leq L |x - y|,$$

$\swarrow$  L-gradient  
Lipschitz

$$(A2) \quad (x - y) \cdot (\nabla U(x) - \nabla U(y)) \geq K |x - y|^2.$$

$\swarrow$  K-strong  
convexity

Suppose  $T > 0$  satisfies  $LT^2 \leq 1/8$   $\leftarrow$  No-U-TURN  
condition

For all distributions  $\nu, \eta \in \mathcal{P}(\mathbb{R}^{2d})$ ,

$$\boxed{W^P(\nu \overset{T}{\parallel}_\infty, \eta \overset{T}{\parallel}_\infty) \leq e^{-c} W^P(\nu, \eta)}$$

where  $c = KT^2/6$ .

$\nwarrow$   $L^P$ -Wasserstein contractivity

How many xHMC steps till  $\text{Law}(X_n) \approx \mu$ ?

Given accuracy  $\varepsilon > 0$ . Choose  $T \propto L^{-1/2}$ .

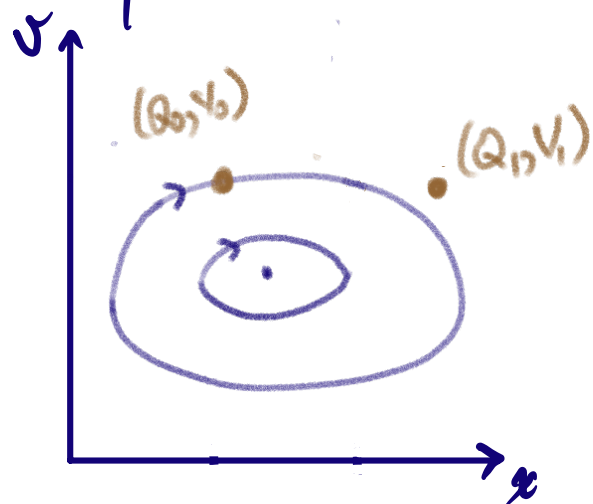
As long as  $m \geq \frac{L}{K} \log \left( \frac{\mathcal{W}^P(\mu, \nu)}{\varepsilon} \right)^+$ ,

$$\mathcal{W}^P(\nu \pi_{\text{ex}}^m, \mu) \leq \varepsilon.$$

$L/K$  is the condition # of  $V$ .

### 3. Unadjusted HMC with time integrator randomization

In practice, exact Hamiltonian flow is unavailable.



$$Q_1 = Q_0 + h V_0 + \frac{h^2}{2} F_0$$

$$V_1 = V_0 + h F_0$$

$$F_0 = -\nabla U(Q_0 + \frac{h}{2} V_0)$$

Idea: evaluate force at a random midpoint

$$F_0 = -\nabla U(Q_0 + U V_0) \text{ where } U \sim \text{Uniform}(0, h)$$

CS: Bou-Rabee, Schuh 2020  $\div$  Bou-Rabee, Eberle 2021

Unadjusted HMC replaces exact Hamilton flow w/  
flow of randomized midpoint method

$$\dot{Q}_t = V_t \quad \dot{V}_t = -\nabla U(Q_{\lfloor t \rfloor} + \gamma_t V_{\lfloor t \rfloor})$$

where  $\gamma_t := U_{\lfloor t \rfloor/h}$  and  $\{U_k\} \stackrel{i.i.d}{\sim} \text{Uniform}(0, h)$

Transition Kernel  $\Pi_u(x, \cdot) = E \left[ \delta_{Q_T(x, \xi)} \right] \quad \xi \sim \mathcal{N}(0, \mathbb{I})$

CS: Shen, Lee 2019 "RMM for Logconcave Sampling"  
Cao, Lu, Wang 2021 "Complexity of Randomized Algorithms..."

## Theorem (Bar-Rabee and Marsden 2022)

Suppose  $U: \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies

- (A0)  $U(0)=0$  and  $\nabla U(0)=0$
- (A1)  $L$ -gradient Lipschitz
- (A2)  $K$ -strongly convex

Suppose  $T > 0$  and  $h > 0$  satisfy  $LT^2 \leq 1/8$  and  $T/h \in \mathbb{Z}$ .

For all distributions  $\nu, \eta \in \mathcal{P}(\mathbb{R}^d)$ ,

$$\mathcal{W}^2(\nu \Pi_u, \eta \Pi_u) \leq e^{-c} \mathcal{W}^2(\nu, \eta)$$

where  $c = KT^2/6$ .

↑  $L^2$ -Wasserstein Contractivity

N.b.  $\nexists!$   $\tilde{\mu} = \tilde{\mu} \Pi_u$  but  $\tilde{\mu} \approx \mu$ .

## Idea of Proof

The proof is based on synchronously coupling two copies of  $\Pi_u$  and applying

$$|Q_T(x, v) - Q_T(y, v)|^2 \leq \left(1 - \frac{KT^2}{3}\right) |x - y|^2$$

which holds for all  $x, y, v \in \mathbb{R}^d$  almost surely.

CG: Bar-Rabie, Schuh 2020 "Convergence of uHMC .."



## Theorem (Bar-Rabee and Marsden 2022)

Suppose  $U: \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies

- (A0)  $U(0)=0$  and  $\nabla U(0)=0$
- (A1)  $L$ -gradient Lipschitz
- (A2)  $K$ -strongly convex

Suppose  $T > 0$  and  $h > 0$  satisfy  $LT^2 \leq 1/8$  and  $T/h \in \mathbb{Z}$ .

It holds that

$$\mathcal{W}^2(\tilde{\mu}, \mu) \leq 142 d^{1/2} c' \left(\frac{L}{K}\right)^{1/2} L^{1/4} h^{3/2}$$

↑  $L^2$ -Wasserstein Asymptotic Bias

cf. Dunsmuir, Eberle 2021 "Asymptotic Bias of Inexact MCMC.."

## Idea of Proof

The proof is based on  $L^2$ -accuracy of the underlying randomized midpoint flow

$$\left( \mathbb{E} \left[ |Q_{t_k}(x, v) - q_{t_k}(x, v)|^2 \right] \right)^{1/2} \leq \gamma \cdot (|v| + \sqrt{L}|x|) L^{1/4} h^{3/2}$$

which holds for all  $x, v \in \mathbb{R}^d$ .

The  $3/2$ -order  $L^2$ -accuracy is a consequence of cancellations due to independence of the random midpoints.

See Fundamental Thm. for  $L^2$ -Convergence Milstein and Tret'yakov 2021.

How many uHMC steps till  $\text{Law}(X_n) \approx \mu$ ?

$$W^2(\mu, \pi_n^m) \leq \underbrace{W^2(\mu, \tilde{\mu})}_{\text{asympt. bias}} + \underbrace{W^2(\tilde{\mu}, \pi_n^m)}_{\text{contractivity}}$$

Given accuracy  $\varepsilon > 0$ . Choose  $T \propto L^{-1/2}$ .

Choose  $h \propto \varepsilon^{2/3} d^{-1/3} c^{2/3} (L/K)^{-1/3} L^{-1/6}$ .

Complexity:  $\frac{T}{h} \cdot m \propto \varepsilon^{-2/3} \left(\frac{d}{K}\right)^{1/3} \left(\frac{L}{K}\right)^{5/3}$   
 $\uparrow$   
 # of gradient evaluations  
 per uHMC step.

vs.  $\varepsilon^{-1} \left(\frac{d}{K}\right)^{1/2} \left(\frac{L}{K}\right)^2$

#### 4. Adjustable Randomized Time Integrators

A Metropolis adjustment step can be used to obtain a Markov chain w/out asymptotic bias.

$$\Pi_{\text{ex}} = \Pi \Gamma$$

$\Gamma((q, v), \cdot) = \delta_q \otimes N(0, I_d)$

$\Pi((q, v), \cdot) = \delta_{Q_1(q, v)}$

Since  $\Gamma$  leaves the  $v$ -marginal of  $\mu_{\text{BG}}$  invariant, we focus on  $\Pi$ .

A key property of the exact Hamiltonian flow is reversibility w.r.t.  $S: (q, v) \mapsto (q, -v)$

$$\mathcal{Q}_{-t} = S \circ \mathcal{Q}_t \circ S.$$

Consequently, the corresponding transition Kernel  $\Pi$  satisfies generalized detailed balance w.r.t.  $\mu_{BG}$

$$\mu_{BG}(dq dv) \Pi((q, v), dq' dv') = \mu_{BG}(dq' dv') \Pi(S(q', v'), S(dq dv))$$

Defn. A randomized time integrator is adjustable if the corresponding transition kernel  $\tilde{\Pi}$  satisfies

$$\mu_{BG}(dq dv) \tilde{\Pi}(S(q, v), S(dq' dv')) = \\ \rho((q, v), (q', v')) \mu_{BG}(dq' dv') \tilde{\Pi}((q', v'), dq dv)$$

Adjusted Kernel

$$\tilde{\Pi}_a(q, v, dq' dv') = \alpha((q, v), (q', v')) \tilde{\Pi}(q, v, dq' dv') \\ + r(q, v) \delta_{S(q, v)}(dq' dv')$$

where  $\alpha((q, v), (q', v')) := \min(1, \rho((q', v'), (q, v)))$ .

## Adjustable Randomized Time Integrators

Let  $\{\theta_i\}_{i \in \mathcal{I}}$  be an indexed family of time integrators

such that for each  $i \in \mathcal{I}$

$$(1) \quad \lambda^{\omega} \circ \theta_i = \lambda^{\omega}$$

$\nwarrow$  vol. preserving

$$(2) \quad \theta_i = S \circ \theta_i^{-1} \circ S$$

$\swarrow$  S-reversible

Let  $\rho$  be a probability distribution over  $\mathcal{I}$ .

$$\text{Let } \theta_{i_N \dots i_1} := \theta_{i_N} \circ \dots \circ \theta_{i_2} \circ \theta_{i_1}.$$

$$\text{Then } \tilde{\Pi}((q, v), \cdot) = \int_{\mathcal{I}^N} \delta_{\theta_{h_N \dots h_1}(q, v)} \prod_{i=1}^N \rho(dh_i)$$

is adjustable.

## Example: Randomized 2-Stage Palindromic Integrator

Let  $I = [0, 1]$ ,  $h > 0$ , and  $b \in I$ .

$$\Theta_b := \mathcal{Q}_{bh}^{(A)} \circ \mathcal{Q}_{h/2}^{(B)} \circ \mathcal{Q}_{(1-2b)h}^{(A)} \circ \mathcal{Q}_{h/2}^{(B)} \circ \mathcal{Q}_{bh}^{(A)}$$

where  $\mathcal{Q}_t^{(A)}(q, v) := (q + tv, v)$  and  $\mathcal{Q}_t^{(B)}(q, v) := (q, v + tF(q))$ .

More explicitly,

$\swarrow$   $b = 1/2$  corresponds to p. Verlet

$$(q, v) \mapsto \left( q + hv + (1-b) \frac{h^2}{2} F_+ + \frac{bh^2}{2} F_-, v + \frac{h}{2} [F_+ + F_-] \right)$$

where  $F_+ := F(x + bhv)$  and  $F_- := F(x + (1-b)hv + (1-2b)\frac{h^2}{2} F_+)$



## Adjusted HMC Algorithm

Input: # of integration steps  $N$ , time step size  $h > 0$ , and  $\rho \in \mathcal{P}(\mathcal{I})$ .

Output:  $(X_0, V_0)$ ,  $(X_1, V_1)$ ,  $\dots$

(Step 1) Independently draw  $\xi_k \sim \mathcal{N}(0, I_d)$  and  $\{u_i^k\}_{i=1}^N \stackrel{\text{i.i.d.}}{\sim} \rho$ .

(Step 2)  $(\tilde{X}_{k+1}, \tilde{V}_{k+1}) \leftarrow \Phi_{u_N^k \dots u_1^k}(X_k, \xi_k)$

(Step 3) Draw  $\beta^k \sim \text{Bern}(\alpha((X_k, \xi_k), (\tilde{X}_{k+1}, \tilde{V}_{k+1})))$ .

(Step 4)  $(X_{k+1}, V_{k+1}) \leftarrow \beta^k \cdot \underbrace{(\tilde{X}_{k+1}, \tilde{V}_{k+1})}_{\text{Accept}} + (1 - \beta^k) \cdot \underbrace{S(X_k, \xi_k)}_{\text{Reject}}.$

## - Summary -

- Hamiltonian MCMC provides a flexible framework for constructing implementable exact/in-exact Markov chains.
- Time integrator randomization relaxes regularity requirements on the target measure.
- Stochastic notions of integrator accuracy fit Hamiltonian MCMC and are more lenient.

## - Questions -

- Can one similarly relax the Hessian Lipschitz condition for TV-convergence bounds?

Cf. Bau-Rabee, Eberle 2021 "Mixing Time Guarantees for uHMC"

- Or, the regularity conditions for the discrete version of Moumarché's Modified Entropy Method?

Cf. Camrud, Durmus, Moumarché, Stoltz 2023

- Can these results be extended to uHMC with partial velocity refreshment?

CF: Monmarché 2022

LeinKuhler, Paulin, Whalley 2023

Appendix of Bou-Rabee, Marsden 2022

Complexity  $\in \mathcal{O}\left(\max\left(\left(\frac{d}{K}\right)^{1/4} \cdot \left(\frac{L}{K}\right)^{3/2} \cdot \varepsilon^{-1/2}, \left(\frac{d}{K}\right)^{1/3} \cdot \left(\frac{L}{K}\right)^{4/3} \cdot \varepsilon^{-2/3}\right)\right)$