

# **CPSC 661: Sampling Algorithms in ML**

---

Andre Wibisono

April 28, 2021

Yale University

## Last time

- Langevin dynamics
- Unadjusted Langevin Algorithm
- Variants of Langevin dynamics
  1. {Weighted, Mirror, Newton, Interacting} Langevin Dynamics
  2. Underdamped Langevin Dynamics

**Today:** Hamiltonian Monte Carlo

# Hamiltonian Flow

Given  $H: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  differentiable (we call  $H$  the Hamiltonian)  
 $(x, v) \mapsto H(x, v)$

The Hamiltonian flow generated by  $H$  is

$$(HF) \quad \begin{cases} \dot{x}_t = \frac{\partial H}{\partial v}(x_t, v_t) \\ \dot{v}_t = -\frac{\partial H}{\partial x}(x_t, v_t) \end{cases}$$

Fact: Hamiltonian flow preserves Hamiltonian  $H$

$$H(x_t, v_t) = H(x_0, v_0)$$

Proof:  $\frac{d}{dt} H(x_t, v_t) = \left\langle \frac{\partial H}{\partial x}(x_t, v_t), \dot{x}_t \right\rangle + \left\langle \frac{\partial H}{\partial v}(x_t, v_t), \dot{v}_t \right\rangle$

$$= \left\langle \frac{\partial H}{\partial x}, \frac{\partial H}{\partial v} \right\rangle + \left\langle \frac{\partial H}{\partial v}, -\frac{\partial H}{\partial x} \right\rangle$$
$$= 0.$$

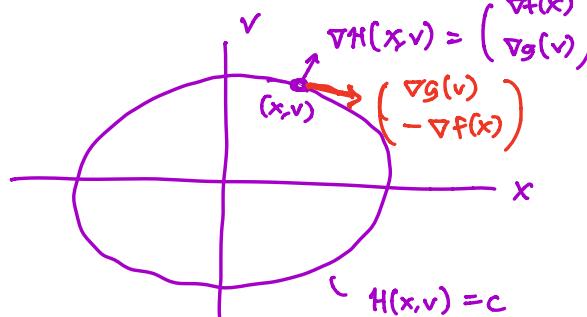
- Eg. Newton's Law:  $\ddot{x}_t = -\nabla f(x_t)$

$$\Leftrightarrow \begin{cases} \dot{x}_t = v_t & = \frac{\partial}{\partial v} \left( \frac{1}{2} \|v_t\|^2 \right) = \frac{\partial H}{\partial v}(x_t, v_t) \\ \dot{v}_t = -\nabla f(x_t) & = -\frac{\partial}{\partial x} (f(x_t)) = -\frac{\partial H}{\partial x}(x_t, v_t) \end{cases}$$

this is the Hamiltonian flow of

$$H(x, v) = \underbrace{f(x)}_{\text{potential energy}} + \underbrace{\frac{1}{2} \|v\|^2}_{\text{kinetic energy}}$$

- Hamiltonian flow of  $H(x, v) = f(x) + g(v)$   
is  $\begin{cases} \dot{x} = \nabla g(v) \\ \dot{v} = -\nabla f(x) \end{cases}$

$$\Leftrightarrow \begin{pmatrix} \dot{x} \\ \dot{v} \end{pmatrix} = \underbrace{\begin{pmatrix} 0 & I \\ -I & 0 \end{pmatrix}}_{\text{Rotation by } 90^\circ} \begin{pmatrix} \nabla f(x) \\ \nabla g(v) \end{pmatrix}$$


- e.g.  $f(x) = \frac{\alpha}{2} x^2$  on  $\mathbb{R}$

Hamiltonian flow  $\begin{cases} \dot{x} = v \\ \dot{v} = -\alpha x \end{cases} \Rightarrow \ddot{x} = -\alpha x$

Solution is 
$$x_t = \cos(\sqrt{\alpha}t) x_0 + \frac{1}{\sqrt{\alpha}} \sin(\sqrt{\alpha}t) \dot{x}_0$$

$$\dot{x}_t = -\sqrt{\alpha} \sin(\sqrt{\alpha}t) x_0 + \cos(\sqrt{\alpha}t) \dot{x}_0$$

$$\ddot{x}_t = -\alpha \cos(\sqrt{\alpha}t) x_0 - \sqrt{\alpha} \sin(\sqrt{\alpha}t) \dot{x}_0$$

$$= -\alpha x_t$$

# Hamiltonian Monte Carlo

To sample from  $\nu(x) \propto e^{-f(x)}$  on  $\mathbb{R}^n$

HMC: ( Hybrid Monte Carlo [Duane et al. 1987]  
Hamiltonian Monte Carlo )

\* From  $X_k \sim g_k$ :

- Run Hamiltonian flow for some integration time  $\eta$ :

$$\begin{aligned}\dot{x}_t &= v_t \\ \dot{v}_t &= -\nabla f(x_t)\end{aligned}\quad \left. \begin{array}{l} \text{this is Hamiltonian flow of} \\ H(x, v) = f(x) + \frac{1}{2} \|v\|^2 \end{array} \right\}$$

from  $(X_0, v_0) = (x_k, z_k)$ , where  $z_k \sim \mathcal{N}(0, I)$  is independent

- set  $x_{\text{int}} = X_\eta$

## Example: Gaussian

$$v = \mathcal{N}(0, \frac{1}{\alpha}) \quad \text{on } \mathbb{R}^1, \quad f(x) = \frac{\alpha}{2} x^2 + \text{constant}$$

Hamiltonian flow  $\dot{x}_t = v_t$

$$\dot{v}_t = -\alpha x_t$$

$$\Rightarrow x_t = \cos(\sqrt{\alpha} t) x_0 + \frac{1}{\sqrt{\alpha}} \sin(\sqrt{\alpha} t) v_0$$

HMC: from  $x_n \sim S_n$ :

\* draw  $z_n \sim \mathcal{N}(0, I)$

\* set  $x_{n+1} = \cos(\sqrt{\alpha} \eta) x_n + \frac{1}{\sqrt{\alpha}} \sin(\sqrt{\alpha} \eta) z_n$

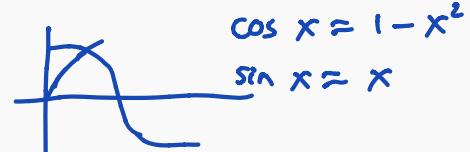
$$e^{i\theta} = \cos \theta + i \sin \theta$$

Consider free integration time  $\eta$ :

i)  $0 < \eta \ll 1$ :  $\cos(\sqrt{\alpha} \eta) \approx 1 - \alpha \eta^2$

$$\sin(\sqrt{\alpha} \eta) \approx \sqrt{\alpha} \eta$$

then  $x_{n+1} \approx (1 - \alpha \eta^2) x_n + \eta z_n$



c.f. OU (Langevin dynamics for  $\nu = \mathcal{N}(0, \frac{1}{\alpha})$ ) is:

$$x_{n+1} = e^{-\alpha\eta} x_n + \sqrt{\frac{1 - e^{-2\alpha\eta}}{\alpha}} z_n, \quad z_n \sim \mathcal{N}(0, I)$$

2) Large  $\eta$ ?

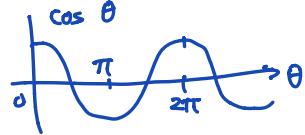
\* if  $\sqrt{\alpha}\eta = 2m\pi$  for some  $m \in \mathbb{N}$

$$\cos(\sqrt{\alpha}\eta) = 1$$

$$\sin(\sqrt{\alpha}\eta) = 0$$

then  $x_{n+1} = x_n$  in HMC

$\rightarrow$  this is bad, no mixing to  $\nu = \mathcal{N}(0, \frac{1}{\alpha})$



\* If  $\sqrt{\alpha}\eta = (2m-1)\pi$

$$\cos(\sqrt{\alpha}\eta) = -1$$

$$\sin(\sqrt{\alpha}\eta) = 0$$

then  $x_{n+1} = -x_n$

$\rightarrow$  this is also bad, no mixing to  $\nu = \mathcal{N}(0, \frac{1}{\alpha})$

\* if  $\sqrt{\alpha}\eta = \frac{\pi}{2} + m\pi$

$$\cos(\sqrt{\alpha}\eta) = 0$$

$$\sin(\sqrt{\alpha}\eta) = \pm 1$$

then HMC becomes  $x_{n+1} = 0 \cancel{x_n} + \frac{1}{\sqrt{\alpha}} z_n = \frac{1}{\sqrt{\alpha}} z_n \sim \mathcal{N}(0, \frac{1}{\alpha}) \Rightarrow$

$\therefore$  optimal integration time is

$$\boxed{\eta = \frac{\pi}{2\sqrt{\alpha}}}$$

then HMC converges in 1 step to  $\nu = \mathcal{N}(0, \frac{1}{\alpha})$

# Why does HMC work?

**Lemma:** In one step of HMC, for any integration time  $\eta > 0$ ,  $x_k \sim \rho_k$  satisfies

$$\text{KL}(\rho_{k+1} \| \nu) \leq \text{KL}(\rho_k \| \nu)$$

$$\text{KL}(s \| \nu) = H_\nu(s) = \int s(x) \log \frac{s(x)}{\nu(x)} dx$$

Proof: From  $x_n \sim s_n$  } set  $(X_0, V_0) = (x_n, v_n) \sim \tilde{\nu}_n(x, v)$   
let  $v_n \sim \mathcal{N}(0, I)$  independent }

In the joint space  $(x, v)$ , the target distribution is

$$\tilde{\nu}(x, v) = \nu(x) \cdot \mathcal{N}(v; 0, I)$$

$$\propto e^{-f(x)} \cdot e^{-\frac{1}{2}\|v\|^2} = e^{-H(x, v)}$$

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



$$\text{In the initialization, } \text{KL}(\tilde{s}_0 \parallel \tilde{\nu}) = \text{KL}(s_0 \parallel \nu) \quad (*)$$

then run Hamiltonian flow

$$\begin{aligned} \dot{x}_t &= v_t \\ \dot{v}_t &= -\nabla f(x_t) \end{aligned} \quad \left. \begin{array}{l} (x_t, v_t) \sim \tilde{s}_t(x, v) \\ \text{satisfies continuity equation} \end{array} \right\}$$

Can show: Hamiltonian flow preserves KL divergence to  $\tilde{\nu}$

$$(*) \quad \text{KL}(\tilde{s}_t \parallel \tilde{\nu}) = \text{KL}(\tilde{s}_0 \parallel \tilde{\nu}) \stackrel{(*)}{=} \text{KL}(s_0 \parallel \nu)$$

then at time  $t=\eta$ , want to set  $x_{\text{init}} = X_\eta$

$$(X_\eta, v_\eta) \sim \tilde{s}_\eta(x, v)$$

$$\text{let } X_\eta \text{ have marginal } \tilde{s}_{\eta, x}(x) = \int_{\mathbb{R}^n} \tilde{s}_\eta(x, v) dv$$

$$v_\eta \text{ have marginal } \tilde{s}_{\eta, v}(v) = \int_{\mathbb{R}^n} \tilde{s}_\eta(x, v) dx$$

then can show: since  $\tilde{\nu}(x, v) = \nu(x) \cdot \mathcal{N}(v; 0, I)$

$$\begin{aligned} (*) \Rightarrow \text{KL}(\tilde{s}_\eta \parallel \tilde{\nu}) &= \text{KL}(\tilde{s}_{\eta, x} \parallel \nu) + \text{KL}(\tilde{s}_{\eta, v} \parallel \mathcal{N}(0, I)) \\ &\quad + I(\tilde{s}_\eta) \end{aligned}$$

$$I(\tilde{s}_\eta) = \int \tilde{s}_\eta(x, v) \log \frac{\tilde{s}_\eta(x, v)}{\tilde{s}_{\eta, x}(x) \cdot \tilde{s}_{\eta, v}(v)}$$

↑ mutual information

Since each term is  $\geq 0$ , when we set

$$x_{\text{init}} = X_\eta$$

$$s_{\text{init}} = \tilde{s}_{\eta, x}$$

$$\Rightarrow \text{KL}(s_{\text{init}} \parallel \nu) \stackrel{(*)}{=} \text{KL}(\tilde{s}_\eta \parallel \tilde{\nu}) = \text{KL}(\tilde{s}_0 \parallel \tilde{\nu}) = \text{KL}(s_0 \parallel \nu).$$

□

# Choosing integration time

How to choose integration time  $\eta$ ?

- Too small  $\rightarrow$  not much progress
- Too large  $\rightarrow$  come back to beginning

# Choosing integration time

How to choose integration time  $\eta$ ?

- Too small  $\rightarrow$  not much progress
- Too large  $\rightarrow$  come back to beginning

Will see:

1. Random integration time [Bou-Rabee & Sanz-Serna 2017]
2. Ideal HMC with short integration time [Chen & Vempala 2019]
3. Metropolized HMC with leapfrog integrator [Chen et al. 2020]

Other ways, e.g. No-U-Turn Sampler [Hoffman & Gelman 2014]

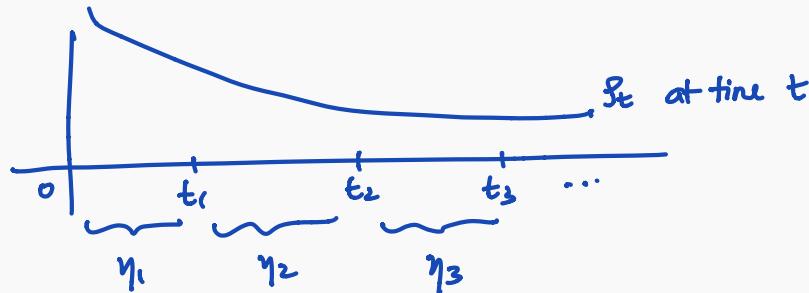
# 1. Randomized Hamiltonian Monte Carlo

HMC with integration time at step  $k$ :

[RHMC]  $\eta_k \sim \text{Exp}(\gamma)$  exp. r.v. with mean  $\frac{1}{\gamma}$   
indep.

$$P(\eta_k \leq t) = 1 - e^{-\frac{t}{\gamma}}$$

Let  $t_h = \eta_1 + \dots + \eta_h$



[Bou-Rabee &  
Sanz-Serna 2018]

Thm: [Lu & Wang 2021]:

Assume  $v \propto e^{-f}$  satisfies  $\alpha$ -Poincaré ineq.

Assume  $\|\nabla^2 f(x)\| \leq M(1 + \|\nabla f(x)\|)$ ,  $f$  is superlinear.

then RHMC with any  $\gamma > 0$  satisfies

$$x_v^2(s_t) \leq c_0 \cdot e^{-2\lambda t} x_v^2(s_0)$$

where  $\lambda = \frac{\alpha\gamma}{(\sqrt{\alpha} + \gamma + R)^2}$

where:

1) if  $f$  is convex, then  $R=0$

then with  $\gamma = \sqrt{\alpha} \Rightarrow \boxed{\lambda = \frac{\sqrt{\alpha}}{4}}$  accelerated rate  
(similar to undamped,  
c.f. Langevin dynamics)

2) if  $\nabla^2 f(x) \geq -K \cdot I$ , then  $R = \sqrt{K}$

then with  $\gamma = \sqrt{\alpha} + \sqrt{K} : \lambda = \frac{\alpha}{4(\sqrt{\alpha} + \sqrt{K})}$

## 2. Ideal HMC with short integration time

Consider HMC with constant  $\eta = \frac{1}{2\sqrt{L}}$

Assume  $f$  is  $\alpha$ -strongly convex,  $L$ -smooth ( $\alpha I \leq \nabla^2 f(x) \leq L I$ )

$$\kappa = \frac{L}{\alpha} \text{ is condition number}$$

Lemma: Hamiltonian flow is a contraction

$$\begin{aligned} \dot{x}_0 &= \dot{y}_0 = v \\ x_0 &\xrightarrow{\quad v \quad \dots} \quad \ddot{x}_t = -\nabla f(x_t) \\ y_0 &\xrightarrow{\quad v \quad \dots} \quad \ddot{y}_t = -\nabla f(y_t) \end{aligned} \quad \left. \begin{array}{l} \text{for } 0 \leq t \leq \frac{1}{2\sqrt{L}} \\ \|x_t - y_t\|^2 \leq \left(1 - \frac{\alpha t^2}{4}\right) \|x_0 - y_0\|^2 \end{array} \right\}$$

Theorem: [Chan & Vempala 2020]

$$\text{Along HMC, } W_2^2(s_n, v) \leq \left(1 - \frac{1}{16\kappa}\right)^k W_2^2(s_0, v)$$

$\Rightarrow$  mixing time is  $\tilde{O}(\kappa)$

\* Matching lower bound:

Consider  $\nu = \mathcal{N}(0, \begin{pmatrix} \frac{1}{\alpha} & 0 \\ 0 & \frac{1}{L} \end{pmatrix})$ , so  $\nabla^2 f(x) = \begin{pmatrix} \alpha & 0 \\ 0 & L \end{pmatrix}$

Theorem: HMC with  $\eta = \frac{1}{c\sqrt{L}}$

has spectral gap  $\gamma \leq \frac{1}{2c^2 K}$

$\Rightarrow$  so mixing time is  $\Omega(K)$

### 3. Metropolis HMC with Leapfrog Integrator

To sample from  $\nu(x) \propto e^{-f(x)}$

\* From  $x_0 \sim \mathcal{S}_0$ :

- Run Leapfrog integrator (to approx. Hamiltonian flow)

from  $(\tilde{x}_m, \tilde{v}_m)$ :

$$\begin{cases} \tilde{v}_{m+\frac{1}{2}} = \tilde{v}_m - \frac{\eta}{2} \nabla f(\tilde{x}_m) \\ \tilde{x}_{m+1} = \tilde{x}_m + \eta \tilde{v}_{m+\frac{1}{2}} \\ \tilde{v}_{m+1} = \tilde{v}_{m+\frac{1}{2}} - \frac{\eta}{2} \nabla f(\tilde{x}_{m+1}) \end{cases}$$

for M times, from  $(\tilde{x}_0, \tilde{v}_0) = (x_0, z_0)$ ,  $z_0 \sim \mathcal{N}(0, I)$

- proposed point  $\tilde{x}_M$

- (Metropolis-Hastings) : accept/reject

$$\begin{aligned} \dot{x}_t &= v_t \\ \dot{v}_t &= -\nabla f(x_t) \end{aligned}$$

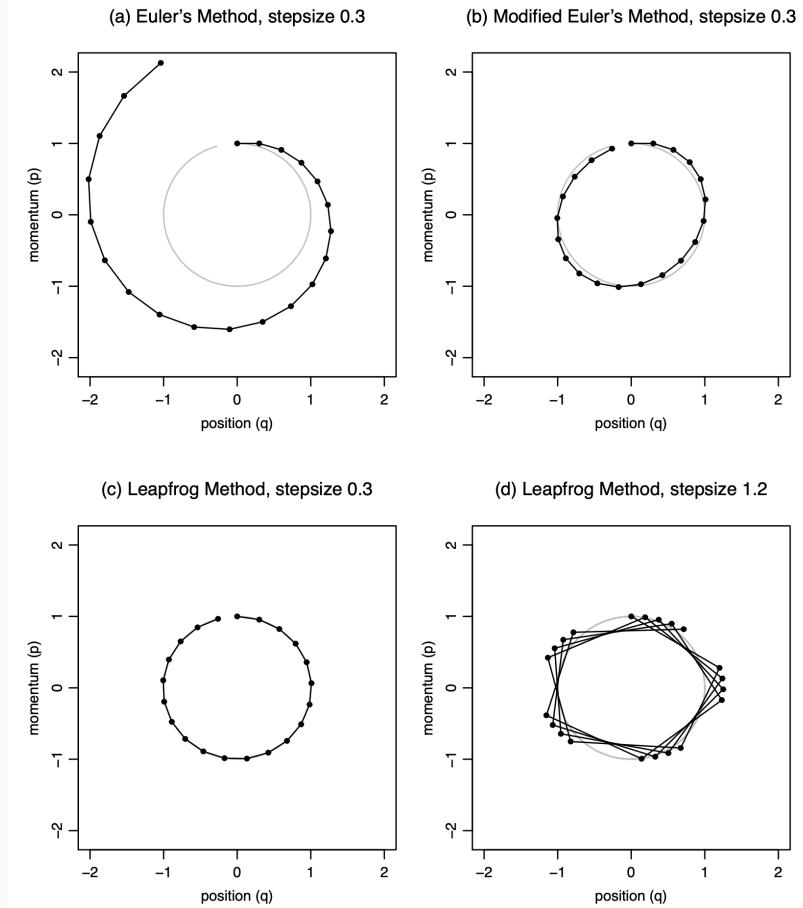
$$\text{set } x_{n+1} = \begin{cases} \tilde{x}_M & \text{with prob. } \alpha_n = \min \left\{ 1, \frac{\tilde{\nu}(\tilde{x}_M, \tilde{v}_M)}{\tilde{\nu}(\tilde{x}_n, \tilde{v}_n)} \right\} \\ x_n & \text{with prob. } 1 - \alpha_n. \end{cases}$$

Notes:

- 1) Acceptance prob.  $\alpha_n$  is simple, because leapfrog is invertible
- 2)  $M = \# \text{ of steps, each with step size } \eta$   
 $\Rightarrow \sim \text{approximating Hamiltonian flow for time } \eta \cdot M = T$
- 3)  $M = 1$ , then this recovers MALA

because  $\tilde{x}_{n+1} = x_n - \frac{\eta^2}{2} \nabla f(x_n) + \eta z_n, \quad z_n \sim N(0, I)$

# Leapfrog Integrator



# Metropolis HMC with Leapfrog

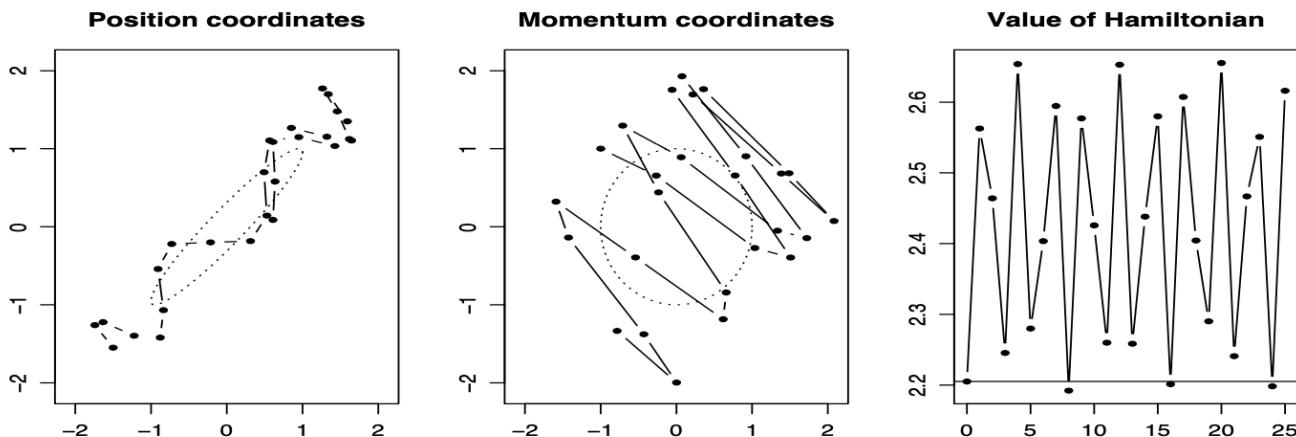


Figure 3: A trajectory for a 2D Gaussian distribution, simulated using 25 leapfrog steps with a stepsize of 0.25. The ellipses plotted are one standard deviation from the means. The initial state had  $q = [-1.50, -1.55]^T$  and  $p = [-1, 1]^T$ .

[Neal, *MCMC using Hamiltonian dynamics*, Handbook of MCMC, 2010]

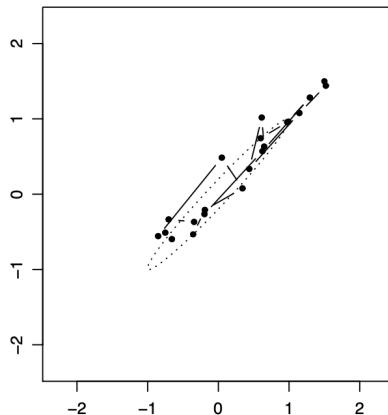
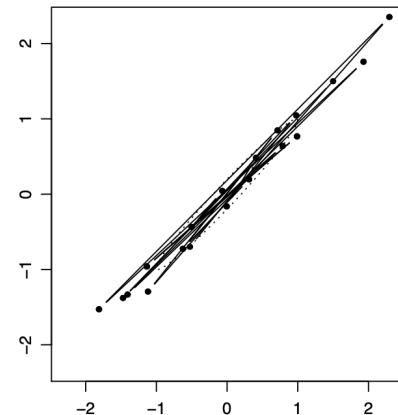
**Random-walk Metropolis****Hamiltonian Monte Carlo**

Figure 4: Twenty iterations of the random-walk Metropolis method (with 20 updates per iteration) and of the Hamiltonian Monte Carlo method (with 20 leapfrog steps per trajectory) for a 2D Gaussian distribution with marginal standard deviations of one and correlation 0.98. Only the two position coordinates are plotted, with ellipses drawn one standard deviation away from the mean.

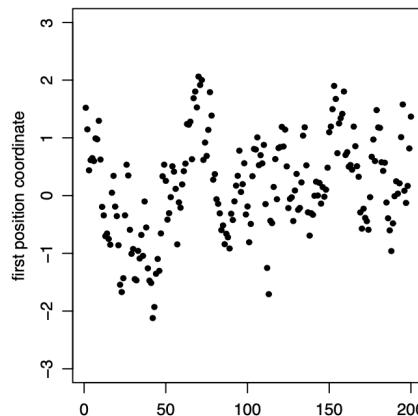
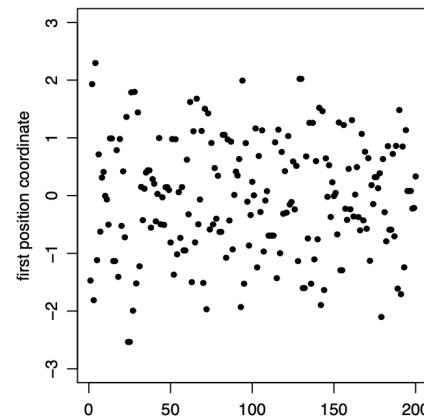
**Random-walk Metropolis****Hamiltonian Monte Carlo**

Figure 5: Two hundred iterations, starting with the twenty iterations shown above, with only the first position coordinate plotted.

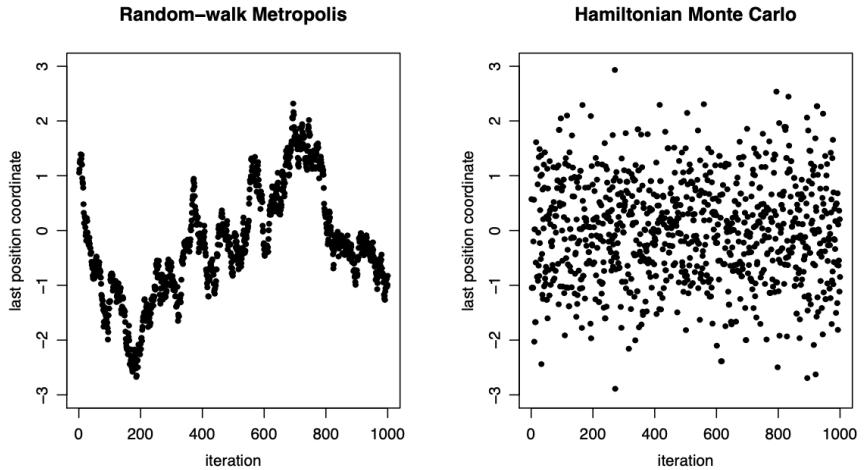


Figure 6: Values for the variable with largest standard deviation for the 100-dimensional example, from a random-walk Metropolis run and an HMC run with  $L = 150$ . To match computation time, 150 updates were counted as one iteration for random-walk Metropolis.

[Neal, *MCMC using Hamiltonian dynamics*, Handbook of MCMC, 2010]

# HMC Mixing Time

$$r(s) := 1 + \max \left\{ \left( \frac{\log(1/s)}{d} \right)^{1/4}, \left( \frac{\log(1/s)}{d} \right)^{1/2} \right\}, \quad (11a)$$

for  $s > 0$ , and involves the step-size choices

$$\eta_{\text{warm}} = \sqrt{\frac{1}{cL \cdot r\left(\frac{\epsilon^2}{2\beta}\right) d^{\frac{7}{6}}}}, \quad \text{and} \quad \eta_{\text{feas}} = \sqrt{\frac{1}{cL \cdot r\left(\frac{\epsilon^2}{2\kappa^d}\right)} \min \left\{ \frac{1}{d\kappa^{\frac{1}{2}}}, \frac{1}{d^{\frac{2}{3}}\kappa^{\frac{5}{6}}}, \frac{1}{d^{\frac{1}{2}}\kappa^{\frac{3}{2}}} \right\}}. \quad (11b)$$

With these definitions, we have the following:

**Corollary 1.** Consider an  $(L, L_H, m)$ -strongly log-concave target distribution  $\Pi^*$  (cf. Assumption (B)) such that  $L_H^{2/3} = O(L)$ , and any error tolerance  $\epsilon \in (0, 1)$ .

(a) Suppose that  $\kappa = O(d^{\frac{2}{3}})$  and  $\beta = O\left(\exp\left(d^{\frac{2}{3}}\right)\right)$ . Then with any  $\beta$ -warm initial distribution  $\mu_0$ , hyper-parameters  $K = d^{\frac{1}{4}}$  and  $\eta = \eta_{\text{warm}}$ , the HMC- $(K, \eta)$  chain satisfies

$$\tau_2^{\text{HMC}}(\epsilon; \mu_0) \leq c d^{\frac{2}{3}} \kappa r\left(\frac{\epsilon^2}{2\beta}\right) \log\left(\frac{\log \beta}{\epsilon}\right). \quad (12a)$$

(b) With the initial distribution  $\mu_\dagger = \mathcal{N}(x^*, \frac{1}{L}\mathbb{I}_d)$ , hyper-parameters  $K = \kappa^{\frac{3}{4}}$  and  $\eta = \eta_{\text{feas}}$ , the HMC- $(K, \eta)$  chain satisfies

$$\tau_2^{\text{HMC}}(\epsilon; \mu_\dagger) \leq c r\left(\frac{\epsilon^2}{2\kappa^d}\right) \max \left\{ d \log \kappa, \max \left[ d, d^{\frac{2}{3}}\kappa^{\frac{1}{3}}, d^{\frac{1}{2}}\kappa \right] \log\left(\frac{d \log \kappa}{\epsilon}\right) \right\}. \quad (12b)$$

# HMC Mixing Time: Warm Start

Assume  $\rho_0$  is  $\beta$ -warm with respect to  $\nu$

Sampling algorithm	Mixing time	#Gradient evaluations
MRW [24, Theorem 2]	$d\kappa^2 \cdot \log \frac{1}{\epsilon}$	NA
MALA [24, Theorem 1]	$d\kappa \cdot \log \frac{1}{\epsilon}$	$d\kappa \cdot \log \frac{1}{\epsilon}$
HMC-( $K, \eta$ ) [ours, Corollary 1]	$d^{\frac{2}{3}} \kappa \cdot \log \frac{1}{\epsilon}$	$d^{\frac{11}{12}} \kappa \cdot \log \frac{1}{\epsilon}$

**Table 2.** Summary of the  $\epsilon$ -mixing time and the corresponding number of gradient evaluations for MRW, MALA and HMC from a *warm start* with an  $(L, L_H, m)$ -strongly-log-concave target. These statements hold under the assumption  $L_H^{2/3} = O(L)$ ,  $\kappa = \frac{L}{m} \ll d$ , and omit logarithmic terms in dimension.

[Chen, Dwivedi, Wainwright, & Yu, *Fast mixing of Metropolized Hamiltonian Monte Carlo: Benefits of multi-step gradients*, JMLR 2020]

# HMC Mixing Time: Feasible Start

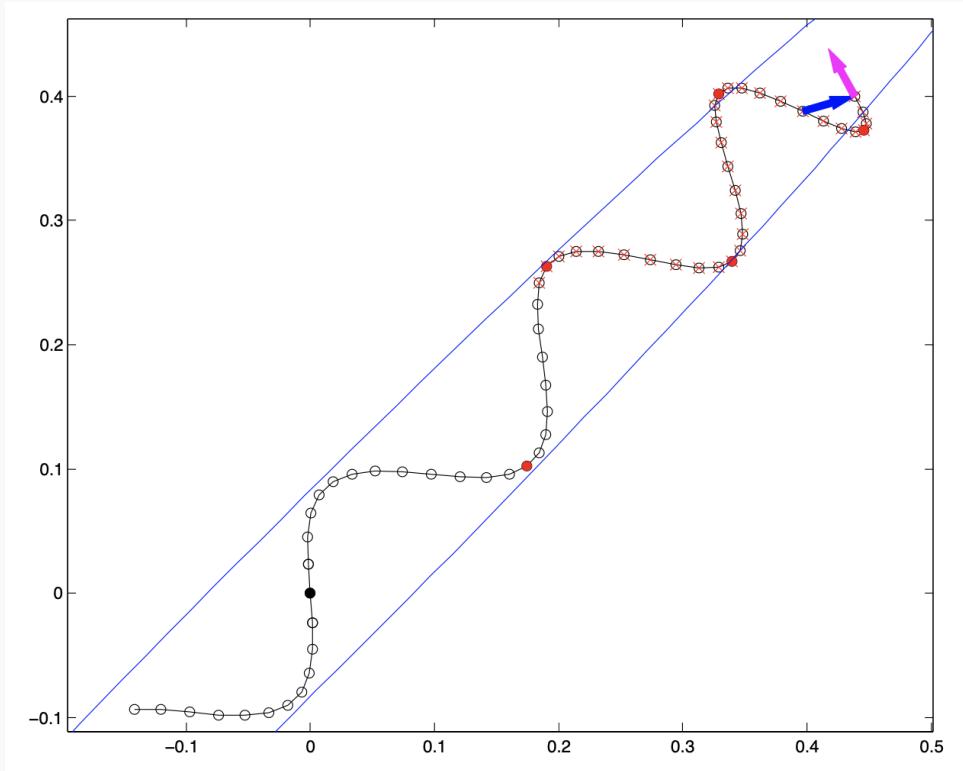
Assume  $\rho_0 = \mathcal{N}(x^*, \frac{1}{L}I)$ , which is warm with  $\beta = \kappa^{d/2}$

Sampling algorithm	Mixing time	# Gradient Evaluations general $\kappa$	# Gradient Evaluations $\kappa \ll d^{\frac{1}{2}}$
MRW [ours, Theorem 2]	$d\kappa^2$	NA	NA
MALA [ours, Theorem 2]	$\max \left\{ d\kappa, d^{\frac{1}{2}}\kappa^{\frac{3}{2}} \right\}$	$\max \left\{ d\kappa, d^{\frac{1}{2}}\kappa^{\frac{3}{2}} \right\}$	$d\kappa$
HMC-( $K, \eta$ ) [ours, Corollary 1]	$\max \left\{ d, d^{\frac{2}{3}}\kappa^{\frac{1}{3}}, d^{\frac{1}{2}}\kappa \right\}$	$\max \left\{ d\kappa^{\frac{3}{4}}, d^{\frac{2}{3}}\kappa^{\frac{13}{12}}, d^{\frac{1}{2}}\kappa^{\frac{7}{4}} \right\}$	$d\kappa^{\frac{3}{4}}$

**Table 3.** Summary of the  $\epsilon$ -mixing time and the corresponding number of gradient evaluations for MRW, MALA and HMC from the *feasible start*  $\mu_{\dagger} = \mathcal{N}(x^*, \frac{1}{L}\mathbb{I}_d)$  for an  $(L, L_H, m)$ -strongly-log-concave target. Here  $x^*$  denotes the unique mode of the target distribution. These statements hold under the assumption  $L_H = O(L^{\frac{3}{2}})$ , and hide the logarithmic factors in  $\epsilon, d$  and  $\kappa = L/m$ .

[Chen, Dwivedi, Wainwright, & Yu, *Fast mixing of Metropolized Hamiltonian Monte Carlo: Benefits of multi-step gradients*, JMLR 2020]

# No-U-Turn Sampler



[Hoffman & Gelman, *The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo*, JMLR 2014]

# References

- Duane, Kennedy, Pendleton, & Roweth, *Hybrid Monte Carlo*, Physics Letters B, 1987
- Neal, *MCMC using Hamiltonian dynamics*, Handbook of MCMC, 2010
- Girolami & Calderhead, *Riemann Manifold Langevin and Hamiltonian Monte Carlo Methods*, Journal of the Royal Statistical Society: Series B, 2011
- Hoffman & Gelman, *The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo*, JMLR 2014
- Bou-Rabee & Sanz-Serna, *Randomized Hamiltonian Monte Carlo*, The Annals of Applied Probability, 2017
- Lee & Vempala, *Convergence rate of Riemannian Hamiltonian Monte Carlo and faster polytope computation*, STOC 2018
- Chen & Vempala, *Optimal Convergence Rate of Hamiltonian Monte Carlo for Strongly Logconcave Distributions*, APPROX/RANDOM 2019
- Chen, Dwivedi, Wainwright, & Yu, *Fast mixing of Metropolized Hamiltonian Monte Carlo: Benefits of multi-step gradients*, JMLR 2020
- Lu & Wang, *On Explicit  $L^2$ -Convergence Rate Estimate for Piecewise Deterministic Markov Processes in MCMC Algorithms*, arXiv 2021

# Recap

---

## I. Classical Sampling

1. Introduction
2. Markov Chain
3. Spectral Theory
4. Conductance
5.  $s$ -Conductance
6. Metropolis-Hastings
7. Isoperimetry
8. From Isoperimetry to Conductance
9. MRW
10. MALA

## II. Sampling as Optimization

11. Optimization and Dynamics
12. Optimization Algorithms
13. Optimization on Manifold
14. Wasserstein Metric
15. Otto Calculus
16. Potential Energy
17. Entropy
18. Langevin Dynamics
19. Langevin Dynamics under Isoperimetry
20. ULA
21. ULA Analysis

### III. Variations

- 22. Variations of Langevin Dynamics
- 23. Underdamped Langevin Dynamics
- 24. Hamiltonian Monte Carlo
- 25. Student presentation I
- 26. Student presentation II

Thank you!