

# Hamiltonian Monte Carlo<sup>1</sup>

Andrew J. Holbrook

UCLA Biostatistics

# Bayesian inference

We observe data  $y_1, \dots, y_N \stackrel{iid}{\sim} p(y_n|\boldsymbol{\theta})$  and assume  $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$ .  
Here,

- ▶  $p(y|\boldsymbol{\theta}) = \prod_{n=1}^N p(y_n|\boldsymbol{\theta})$  is the *likelihood*,
- ▶  $p(\boldsymbol{\theta})$  is the *prior*,

and the goal of Bayesian inference is to obtain the *posterior*

$$p(\boldsymbol{\theta}|y) = \frac{p(y|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(y)} = \frac{p(y|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int_{\Theta} p(y|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}.$$

# Bayesian inference

We're usually interested in computing another integral

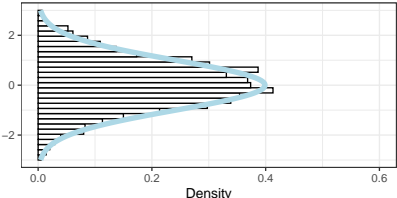
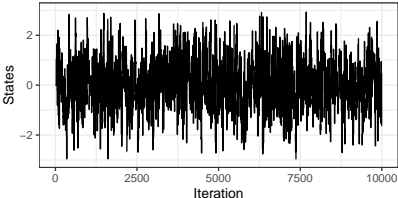
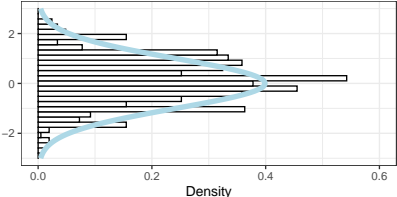
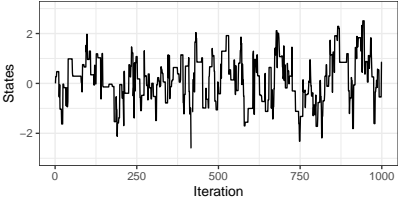
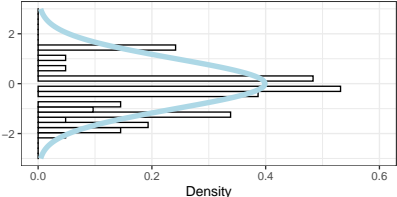
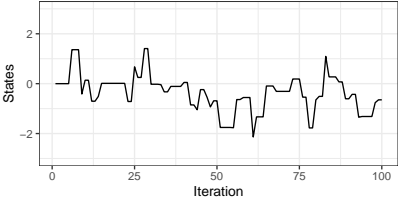
$$\mathbb{E}_{\theta|y} f(\theta) = \int_{\Theta} f(\theta) p(\theta|y) d\theta,$$

so we do what statisticians have been doing forever. We collect samples and rely on the law of large numbers. Suppose  $\theta_1, \dots, \theta_S \stackrel{iid}{\sim} p(\theta|y)$  ( $\mathbb{E}_{\theta|y} |\theta| < \infty$ ) and  $f(\cdot)$  a.s. continuous, then

- ▶ (WLLN)  $\sum_{s=1}^S f(\theta_s)/S \xrightarrow{P} \mathbb{E}_{\theta|y} f(\theta)$
- ▶ (SLLN)  $\sum_{s=1}^S f(\theta_s)/S \xrightarrow{a.s.} \mathbb{E}_{\theta|y} f(\theta)$

But where do we find our samples?

# Markov chain Monte Carlo



# Rejection sampling

We want to sample from generic  $p(\boldsymbol{\theta})$  but only know  $p^*(\boldsymbol{\theta}) \propto p(\boldsymbol{\theta})$ . We can easily sample from  $q(\boldsymbol{\theta})$  and know a number  $M > 0$  s.t.  $p^*(\boldsymbol{\theta}) < Mq(\boldsymbol{\theta})$ .

Algorithm for generating  $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$ :

1. Draw  $\boldsymbol{\theta}^* \sim q(\boldsymbol{\theta})$  and  $U \sim U(0, 1)$
2.  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^*$  if  $U < \frac{p^*(\boldsymbol{\theta}^*)}{Mq(\boldsymbol{\theta}^*)}$

# Funky rejection sampler

We want to sample from generic  $p(\theta)$  and we *can*. But we also feel the compulsion to use the invertible, differentiable function  $T(\cdot)$ .

Algorithm for generating  $\theta \sim p(\theta)$ :

1. Draw  $\theta^\dagger \sim p(\theta)$  and  $U \sim U(0, 1)$
2. Obtain  $\theta^* = T(\theta^\dagger)$
3. Recognize that  $q(\theta^*) = p(T^{-1}(\theta^*)) |\nabla T^{-1}(\theta^*)|$
4.  $\theta \leftarrow \theta^*$  if

$$\begin{aligned} U < \frac{p(\theta^*)}{M q(\theta^*)} &= \frac{p(\theta^*)}{M p(T^{-1}(\theta^*)) |\nabla T^{-1}(\theta^*)|} \\ &= \frac{p(\theta^*)}{M p(\theta^\dagger)} |\nabla T(\theta^\dagger)| \end{aligned}$$

# Metropolis-Hastings-Green

An MCMC algorithm for sampling from  $p(\boldsymbol{\theta})$ . Let  $\boldsymbol{\theta} = \boldsymbol{\theta}^{(s-1)}$  be the former Markov chain state.

1. Draw  $\boldsymbol{\psi} \sim q(\boldsymbol{\theta}, \cdot)$ ,  $U \sim U(0, 1)$ ;
2. Obtain  $(\boldsymbol{\theta}^*, \boldsymbol{\psi}^*) = T(\boldsymbol{\theta}, \boldsymbol{\psi})$ ;
3. Obtain

$$r(\boldsymbol{\theta}, \boldsymbol{\psi}) = \frac{p(\boldsymbol{\theta}^*)q(\boldsymbol{\theta}^*, \boldsymbol{\psi}^*)}{p(\boldsymbol{\theta})q(\boldsymbol{\theta}, \boldsymbol{\psi})} \cdot |\nabla T(\boldsymbol{\theta}, \boldsymbol{\psi})|;$$

4. if  $U < r(\boldsymbol{\theta}, \boldsymbol{\psi})$ , then  $\boldsymbol{\theta}^{(s)} \leftarrow \boldsymbol{\theta}^*$ ,  
else  $\boldsymbol{\theta}^{(s)} \leftarrow \boldsymbol{\theta}$ .

HMC main idea: let  $T(\cdot, \cdot)$  describe the evolution of a Hamiltonian system. But why?

# Change of notation

1.  $\theta \mapsto \mathbf{q}$

2.  $\psi \mapsto \mathbf{p}$

3.  $p(\theta|y) \mapsto \pi(\mathbf{q})$



# Hamiltonian dynamics

- ▶  $q \in \mathbb{R}^D$  is the position of an object
- ▶  $p \in \mathbb{R}^D$  is the momentum of an object
- ▶  $M \in \mathbb{R}^{D \times D}$  is the mass matrix of an object
- ▶  $U(q) = -\log \pi(q)$  is the system's potential energy
- ▶  $K(p) = p^T M^{-1} p / 2$  is the kinetic energy
- ▶  $H(q, p) = U(q) + K(p)$  is the total energy
- ▶ System governed by Hamiltonian equations:

$$\begin{aligned}\frac{dq}{dt} &= \frac{\partial H}{\partial p} = M^{-1} p \\ \frac{dp}{dt} &= -\frac{\partial H}{\partial q} = \nabla \log \pi(q).\end{aligned}$$

- ▶ The unique function  $T_t : (q_0, p_0) \mapsto (q_t, p_t)$  satisfies these equations.

## Hamiltonian dynamics are

1. reversible: replace  $p$  with  $-p$
2. volume preserving ( $|\nabla T| = 1$ ):

$$\begin{aligned}\nabla \cdot \left( \frac{dq}{dt}, \frac{dp}{dt} \right) &= \sum_{d=1}^D \left( \frac{\partial}{\partial q_d} \frac{dq_d}{dt} + \frac{\partial}{\partial p_d} \frac{dp_d}{dt} \right) \\ &= \sum_{d=1}^D \left( \frac{\partial}{\partial q_d} \frac{\partial H}{\partial p_d} - \frac{\partial}{\partial p_d} \frac{\partial H}{\partial q_d} \right) = 0\end{aligned}$$

3. energy conserving:

$$\begin{aligned}\frac{dH}{dt} &= \sum_{d=1}^D \left( \frac{\partial H}{\partial q_d} \frac{dq_d}{dt} + \frac{\partial H}{\partial p_d} \frac{dp_d}{dt} \right) \\ &= \sum_{d=1}^D \left( \frac{\partial H}{\partial q_d} \frac{\partial H}{\partial p_d} - \frac{\partial H}{\partial p_d} \frac{\partial H}{\partial q_d} \right) = 0\end{aligned}$$

# Hamiltonian Monte Carlo (sort of)

Augment parameter space with auxiliary Gaussian variable  $\mathbf{p}$  and construct a Hamiltonian energy function:

$$\begin{aligned} H(\mathbf{q}, \mathbf{p}) &= -\log(\pi(\mathbf{q}) \times \phi_M(\mathbf{p})) \\ &\propto -\log \pi(\mathbf{q}) + \frac{1}{2} \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}. \end{aligned}$$

Given  $\mathbf{q} = \mathbf{q}^{(s-1)}$  a new state of the Markov chain is proposed by forward integrating Hamilton's equations for time  $t$ .

1. Draw  $\mathbf{p}_0 \sim N(0, \mathbf{M})$ ;
2. Obtain  $(\mathbf{q}_t, \mathbf{p}_t) = T_t(\mathbf{q}_0, \mathbf{p}_0)$ ;
3. Accept  $\mathbf{q}_t$  with probability

$$1 \wedge \frac{\pi(\mathbf{q}_t) \phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0) \phi_M(\mathbf{p}_0)} \cdot |\nabla T| = 1 \wedge \frac{\pi(\mathbf{q}_t) \phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0) \phi_M(\mathbf{p}_0)}$$

# Hamiltonian Monte Carlo (sort of)

1. Draw  $\mathbf{p}_0 \sim N(0, M)$ ;
2. Obtain  $(\mathbf{q}_t, \mathbf{p}_t) = T_t(\mathbf{q}_0, \mathbf{p}_0)$ ;
3. Accept  $\mathbf{q}_t$  with probability

$$1 \wedge \frac{\pi(\mathbf{q}_t)\phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0)\phi_M(\mathbf{p}_0)} \cdot |\nabla T| = 1 \wedge \frac{\pi(\mathbf{q}_t)\phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0)\phi_M(\mathbf{p}_0)}$$

But wait!!

$$\begin{aligned} \frac{\pi(\mathbf{q}_t)\phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0)\phi_M(\mathbf{p}_0)} &= \exp\left(\log\left(\frac{\pi(\mathbf{q}_t)\phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0)\phi_M(\mathbf{p}_0)}\right)\right) \\ &= \exp(H(\mathbf{q}_0, \mathbf{p}_0) - H(\mathbf{q}_t, \mathbf{p}_t)) = 1 \end{aligned}$$

# Leapfrog integrator

We need a numerical integrator to solve the Hamiltonian equations. The most popular is the Störmer-Verlet (velocity Verlet) or leapfrog method.

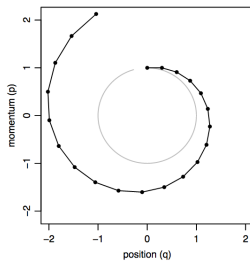
$$p(t + \epsilon/2) = p(t) + \frac{\epsilon}{2} \nabla \log \pi(q(t))$$

$$q(t + \epsilon) = q(t) + \epsilon M^{-1} p(t + \epsilon/2)$$

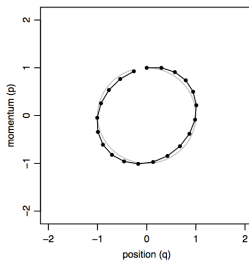
$$p(t + \epsilon) = q(t + \epsilon) + \frac{\epsilon}{2} \nabla \log \pi(q(t + \epsilon))$$

# Leapfrog integrator

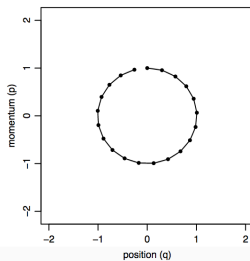
(a) Euler's Method, stepsize 0.3



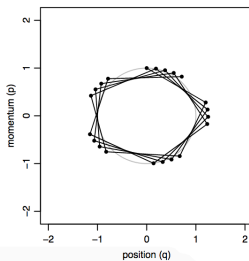
(b) Modified Euler's Method, stepsize 0.3



(c) Leapfrog Method, stepsize 0.3



(d) Leapfrog Method, stepsize 1.2



# Leapfrog integrator

Still reversible (flip  $\mathbf{p}$ ) and volume preserving. To see the latter, the Jacobians are

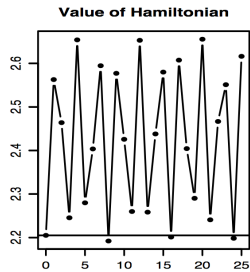
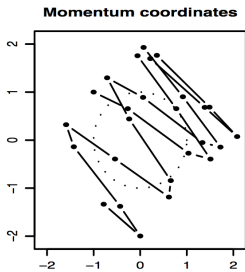
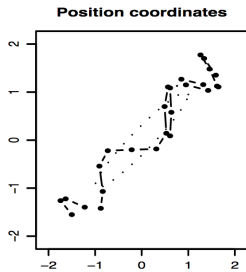
$$\nabla \hat{T}_p(\mathbf{q}, \mathbf{p}) = \begin{pmatrix} 1 & 0 \\ \frac{\epsilon}{2} \nabla^2 \log \pi(\mathbf{q}) & 1 \end{pmatrix}$$

and

$$\nabla \hat{T}_q(\mathbf{q}, \mathbf{p}) = \begin{pmatrix} 1 & \epsilon \mathbf{M}^{-1} \\ 0 & 1 \end{pmatrix}.$$

Unfortunately, we lose energy conservation. The error incurred by a leapfrog trajectory is  $O(\epsilon^2)$ .

# Leapfrog integrator





# Hamiltonian Monte Carlo

For  $t = L \times \epsilon$ :

1. Draw  $\mathbf{p}_0 \sim N(0, M)$ ;
2. Obtain  $(\mathbf{q}_t, \mathbf{p}_t) = \hat{T}_\epsilon^L(\mathbf{q}_0, \mathbf{p}_0)$ ;
3. Accept  $\mathbf{q}_t$  with probability

$$1 \wedge \frac{\pi(\mathbf{q}_t)\phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0)\phi_M(\mathbf{p}_0)} \cdot |\nabla \hat{T}| = 1 \wedge \frac{\pi(\mathbf{q}_t)\phi_M(\mathbf{p}_t)}{\pi(\mathbf{q}_0)\phi_M(\mathbf{p}_0)}$$

# Challenges

1. Ill-conditioned target distributions
2. Multimodal target distributions
3. Big data
4. Fast, flexible and friendly software

# Things we do in practice

1. Jitter  $L$ , the number of leapfrog iterations, at each step.
2. Choose  $\epsilon$  so that majority (55%, say?) of proposals are accepted (for, say,  $L = 100$ ).
3. Adapt  $L$  so that 70% or 80% or 90% of proposals are accepted.
4. Adapt  $M$  using the log posterior Hessian or the empirical covariance of  $p$ .
5. Parallel implementation of computational bottlenecks, i.e., log likelihood and its gradient.

## Neal citation

<sup>1</sup>Many of the images on these slides were taken from

Neal, R. M. (2011). *MCMC using Hamiltonian dynamics*.  
Handbook of Markov chain Monte Carlo, 2(11), 2.