

# Hamiltonian Monte Carlo

Dr. Jarad Niemi

Iowa State University

September 12, 2017

Adapted from Radford Neal's MCMC Using Hamiltonian Dynamics in Handbook of Markov Chain Monte Carlo (2011).

# Hamiltonian system

Considering a body in a frictionless 1-dimensional environment, let

- $m$  be its mass,
- $q$  be its position, and
- $p$  be its momentum.

The mass has

- potential energy  $U(q)$  (which is proportional to its height) and
- kinetic energy  $K(p) = p^2/(2m)$ .

# Hamilton's equations

Extending this to  $d$  dimensions, we have

- position vector  $q$  and
- momentum vector  $p$ .

The Hamiltonian  $H(q, p)$  describes the time evolution of the system through

$$\begin{aligned}\frac{dq_i}{dt} &= \frac{\partial H}{\partial p_i} \\ \frac{dp_i}{dt} &= -\frac{\partial H}{\partial q_i}\end{aligned}$$

for  $i = 1, \dots, d$ .

# Potential and kinetic energy

For Hamiltonian Monte Carlo, we usually use Hamiltonian functions that can be written as follows:

$$H(q, p) = U(q) + K(p)$$

where

- $U(q)$  is called the potential energy and will be defined to be minus the log probability density of the distribution for  $q$  (plus any constant that is convenient) and
- $K(p)$  is called the kinetic energy and is usually defined as

$$K(p) = p^\top M^{-1} p / 2$$

where  $M$  is a symmetric, positive-definite “mass matrix”, which is typically diagonal, and is often a scalar multiple of the identity matrix. This form for  $K(p)$  corresponds to minus the log probability density (plus a constant) of the zero-mean Gaussian distribution with covariance matrix  $M$ .

The resulting Hamilton's equations are

$$\frac{dq_i}{dt} = [M^{-1}p]_i, \quad \frac{dp_i}{dt} = -\frac{\partial U}{\partial q_i}.$$

# One-dimensional example

Suppose

$$H(q, p) = U(q) + K(p), \quad U(q) = q^2/2, \quad K(p) = p^2/2$$

The dynamics resulting from this Hamiltonian are

$$\frac{dq}{dt} = p, \quad \frac{dp}{dt} = -q.$$

Solutions of the form

$$q(t) = r \cos(a + t), \quad p(t) = -r \sin(a + t)$$

for some constants  $r$  and  $a$ .

# One-dimensional example simulation

Hamiltonian dynamics is reversible, i.e. the mapping  $T_s$  from the state at time  $t$ ,  $(q(t), p(t))$ , to the state at time  $t + s$ ,  $(q(t + s), p(t + s))$ , is one-to-one, and hence as an inverse,  $T_{-s}$ . Under our usual assumptions for HMC, the inverse mapping can be obtained by negating  $p$ , applying  $T_s$ , and then negating  $p$  again. The reversibility of Hamiltonian dynamics is important for showing convergence of HMC.

# Conservation of the Hamiltonian

The dynamics conserve the Hamiltonian since

$$\begin{aligned}\frac{dH}{dt} &= \sum_{i=1}^d \left[ \frac{dq_i}{dt} \frac{\partial H}{\partial q_i} + \frac{dp_i}{dt} \frac{\partial H}{\partial p_i} \right] \\ &= \sum_{i=1}^d \left[ \frac{\partial H}{\partial p_i} \frac{\partial H}{\partial q_i} - \frac{\partial H}{\partial q_i} \frac{\partial H}{\partial p_i} \right]\end{aligned}$$

If  $h$  is conserved, then the acceptance probability based on Hamiltonian dynamics is 1. In practice, we can only make  $H$  approximately invariant.



# Conservation of the Hamiltonian

# Volume preservation

If we apply the mapping  $T_s$  to point in some region  $R$  of  $(q, p)$  space with volume  $V$ , the image of  $R$  under  $T_s$  will also have volume  $V$ . This feature simplifies calculation of the acceptance probability for Metropolis updates.

# Euler's method

For simplicity, assume

$$H(q, p) = U(q) + K(p), \quad K(p) = \sum_{i=1}^d \frac{p_i^2}{2m_i}.$$

One way to simulate Hamiltonian dynamics is to discretize time into increments of  $e$ , i.e.

$$\begin{aligned} p_i(t+e) &= p_i(t) + e \frac{dp_i}{dt}(t) = p_i(t) - e \frac{\partial U}{\partial q_i}(q(t)) \\ q_i(t+e) &= q_i(t) + e \frac{dq_i}{dt}(t) = q_i(t) + e \frac{p_i(t)}{m_i} \end{aligned}$$

# Leapfrog method

An improved approach is the **leapfrog** method which has the following updates:

$$\begin{aligned}p_i(t + e/2) &= p_i(t) - (e/2) \frac{\partial U}{\partial q_i}(q(t)) \\q_i(t + e) &= q_i(t) + e \frac{p_i(t+e/2)}{m_i} \\p_i(t + e) &= p_i(t + e/2) - (e/2) \frac{\partial U}{\partial q_i}(q(t + e))\end{aligned}$$

The leapfrog method is reversible and preserves volume exactly.

# Leap-frog simulator

```
leap_frog = function(U, grad_U, e, L, theta, omega) {  
  omega = omega - e/2 * grad_U(theta)  
  
  for (l in 1:L) {  
    theta = theta + e * omega  
    if (l<L) omega = omega - e * grad_U(theta)  
  }  
  omega = omega - e/2 * grad_U(theta)  
  return(list(theta=theta, omega=omega))  
}
```

# Leap-frog simulator

# Conservation of the Hamiltonian

# Probability distributions

The Hamiltonian is an energy function for the joint state of “position”,  $q$ , and “momentum”,  $p$ , and so defines a joint distribution for them, via

$$P(q, p) = \frac{1}{Z} \exp(-H(q, p))$$

where  $Z$  is the normalizing constant.

If  $H(q, p) = U(q) + K(p)$ , the joint density is

$$P(q, p) = \frac{1}{Z} \exp(-U(q)) \exp(-K(p)).$$

If we are interested in a posterior distribution, we set  $q = \theta$  and

$$U(\theta) = -\log[p(y|\theta)p(\theta)].$$



# Hamiltonian Monte Carlo algorithm

Set tuning parameters

- $L$ : the number of steps
- $e$ : stepsize
- $D = \{d_i\}$ : covariance matrix for  $\omega$

Let  $\theta^{(i)}$  be the current value of the parameter  $\theta$ . The leap-frog Hamiltonian Monte Carlo algorithm is

1. Sample  $\omega \sim N_d(0, D)$ .
2. Simulate Hamiltonian dynamics on location  $\theta^{(i)}$  and momentum  $\omega$  via the leapfrog method (or any reversible method that preserves volume) for  $L$  steps with stepsize  $e$ . Call these updated values  $\theta^*$  and  $-\omega^*$ .
3. Set  $\theta^{(i+1)} = \theta^*$  with probability  $\min\{1, \rho(\theta^{(i)}, \theta^*)\}$  where

$$\rho(\theta^{(i)}, \theta^*) = \frac{p(\theta^*|y)}{p(\theta^{(i)}|y)} \frac{p(\omega^*)}{p(\omega^{(i)})} = \frac{p(y|\theta^*)p(\theta^*)}{p(y|\theta^{(i)})p(\theta^{(i)})} \frac{N_d(\omega^*; 0, D)}{N_d(\omega^{(i)}; 0, D)}$$

otherwise set  $\theta^{(i+1)} = \theta^{(i)}$ .

# Reversibility

A reversible simulation means that

- if you simulate from  $(\theta, \omega)$  to  $(\theta', \omega')$  for some step size  $e$  and number of steps  $L$  then
- if you simulate from  $(\theta', \omega')$  for the same  $e$  and  $L$ , you will end up at  $(\theta, \omega)$ .

If we use  $q$  to denote our simulation “density”, then reversibility means

$$q(\theta', \omega' | \theta, \omega) = q(\theta, \omega | \theta', \omega')$$

and thus in the Metropolis-Hastings calculation, the proposal is symmetric. In order to ensure reversibility of our proposal, we need to negate momentum after we complete the leap-frog simulation, but so long as  $p(\omega) = p(-\omega)$  this will not affect our acceptance probability.

# Volume preserving results in perfect acceptance

Recall that we accept with probability  $\min\{1, \rho(\theta^{(i)}, \theta^*)\}$  where

$$\rho(\theta^{(i)}, \theta^*) = \frac{p(\theta^*|y)}{p(\theta^{(i)}|y)} \frac{p(\omega^*)}{p(\omega^{(i)})}$$

Volume is preserved if

$$p(\theta^{(i)}|y)p(\omega^{(i)}) = p(\theta^*|y)p(\omega^*) \implies \frac{p(\theta^*|y)}{p(\theta^{(i)}|y)} \frac{p(\omega^*)}{p(\omega^{(i)})} = 1$$

This will only be the case if the simulation is perfect! But we have discretization error. The acceptance probability accounts for this error.

```

HMC_neal = function(U, grad_U, epsilon, L, current_q) {
  q = current_q
  p = rnorm(length(q),0,1)
  current_p = p

  p = p-epsilon*grad_U(q)/2

  for (i in 1:L) {
    q = q+epsilon*p
    if (i!=L) p = p -epsilon * grad_U(q)
  }
  p = p-epsilon * grad_U(q)/2

  p = -p

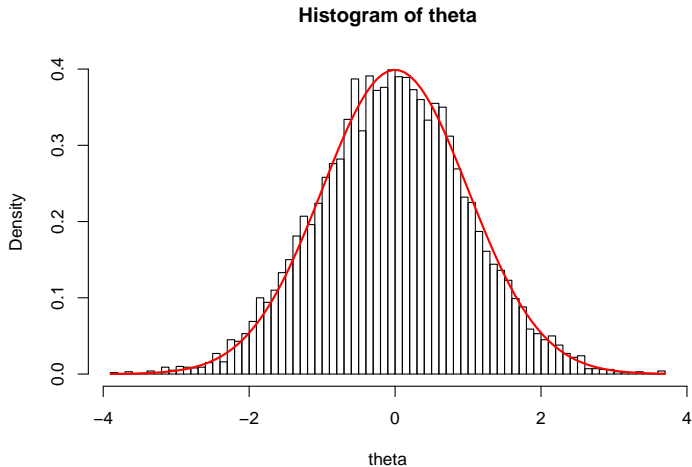
  current_U = U(current_q)
  current_K = sum(current_p^2)/2
  proposed_U = U(q)
  proposed_K = sum(p^2)/2

  if (runif(1) < exp(current_U-proposed_U+current_K-proposed_K))
  {
    return(q)
  }
  else {
    return(current_q)
  }
}

```

```
HMC = function(n_reps, log_density, grad_log_density, tuning, initial) {  
  theta = rep(0, n_reps)  
  theta[1] = initial$theta  
  
  for (i in 2:n_reps) theta[i] = HMC_neal(U = function(x) -log_density(x),  
                                           grad_U = function(x) -grad_log_density(x),  
                                           e = tuning$e,  
                                           L = tuning$L,  
                                           theta[i-1])  
  
  theta  
}
```

```
theta = HMC(1e4, function(x) -x^2/2, function(x) -x, list(e=1,L=1), list(theta=0))  
hist(theta, freq=F, 100)  
curve(dnorm, add=TRUE, col='red', lwd=2)
```



# Tuning parameters

There are three tuning parameters:

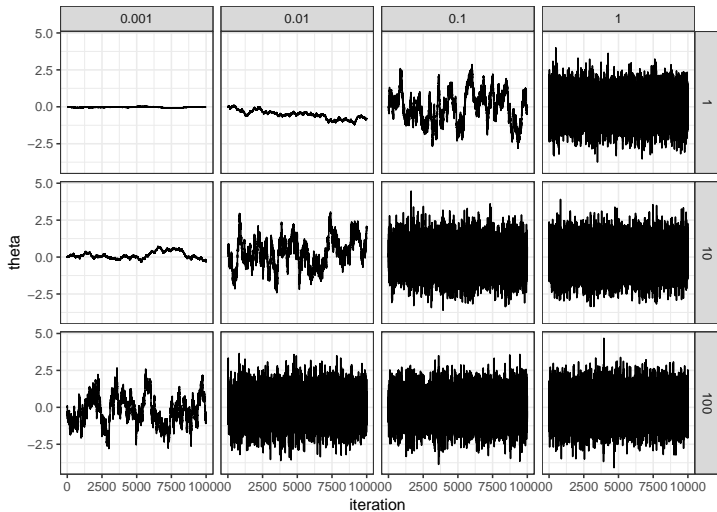
- $\epsilon$ : step size
- $L$ : number of steps
- $D$ : covariance matrix for momentum

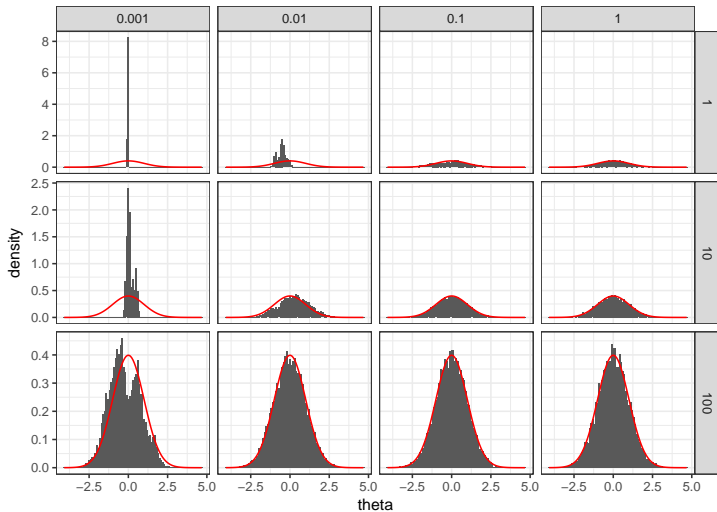
Let  $\Sigma = V(\theta|y)$ , then an optimal normal distribution for  $\omega$  is  $N(0, \Sigma^{-1})$ . Typically, we do not know  $\Sigma$ , but we can estimate it using posterior samples. We can update this estimate throughout burn-in (or warm-up).

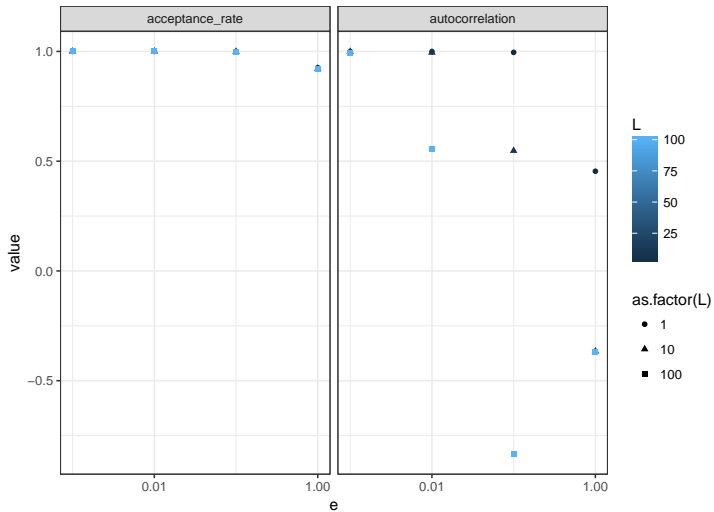
# Effect of $e$ and $L$

```
n_reps = 1e4
d = expand.grid(e=10^seq(-3,0,by=1), L=10^seq(0,2))
r = ddply(d, .(e,L), function(xx) {
  data.frame(
    iteration = 1:n_reps,
    theta = HMC(n_reps, function(x) -x^2/2, function(x) -x, list(e=xx$e,L=xx$L), list(theta=0)))
  })
```









# Random-walk vs HMC

<https://www.youtube.com/watch?v=Vv3f0QNWvWQ>

# Summary

Hamiltonian Monte Carlo (HMC) is a Metropolis-Hastings method using parameter augmentation and a sophisticated proposal distribution based on Hamiltonian dynamics such that

- the acceptance probability can be kept near 1
- while still efficiently exploring the posterior.

HMC still requires us to set tuning parameters

- $\epsilon$ : step size
- $L$ : number of steps
- $D$ : covariance matrix for momentum

and can only be run in models with continuous parameters in  $\mathbb{R}^d$  (or transformed to  $\mathbb{R}^d$ ).