# Hamiltonian Monte Carlo

Dr. Jarad Niemi

Iowa State University

September 12, 2017

Adapted from Radford Neal's MCMC Using Hamltonian 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,
- $\theta$ be its position, and
- $\omega$ be its momentum.

The mass has

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

# Hamilton's equations

Extending this to $d$ dimensions, we have

- position vector $\theta$ and
- momentum vector $\omega$.

The Hamiltonian $H(\theta, \omega)$ describes the time evolution of the system through

$$\begin{array}{rcl} \frac{d\theta_i}{dt} & = & \frac{\partial H}{\partial \omega_i} \\ \frac{d\omega_i}{dt} & = & -\frac{\partial H}{\partial \theta_i} \end{array}$$

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

# Potential and kinetic energy

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

$$H(\theta, \omega) = U(\theta) + K(\omega)$$

where

- $U(\theta)$ is called the potential energy and will be defined to be minus the log probability density of the distribution for $\theta$ (plus any constant that is convenient) and

- $K(\omega)$ is called the kinetic energy and is usually defined as

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

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 $K(\omega)$ 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

$$\frac{d\theta_i}{dt} = [M^{-1}\omega]_i, \qquad \frac{d\omega_i}{dt} = -\frac{\partial U}{\partial \theta_i}.$$

## One-dimensional example

Suppose

$$H(\theta, \omega) = U(\theta) + K(\omega), \quad U(\theta) = \theta^2/2, \quad K(\omega) = \omega^2/2$$

The dynamics resulting from this Hamiltonian are

$$\frac{d\theta}{dt} = \omega, \quad \frac{d\omega}{dt} = -\theta.$$

Solutions of the form

$$\theta(t) = r\cos(a + t), \quad \omega(t) = -r\sin(a + t)$$

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

# Conservation of the Hamiltonian

The dynamics conserve the Hamiltonian since

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

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

# Conservation of the Hamiltonian

# Volume preservation

If we apply the mapping $T_s$ to point in some region $R$ of $(\theta, \omega)$ 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

$$H(\theta, \omega) = U(\theta) + K(\omega), \qquad K(\omega) = \sum_{i=1}^{d} \frac{\omega_i^2}{2m_i}.$$

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

$$
\begin{aligned}
\omega_i(t + e) &= \omega_i(t) + e\frac{d\omega_i}{dt}(t) &= \omega_i(t) - e\frac{\partial U}{\partial \theta_i}(\theta(t)) \\
\theta_i(t + e) &= \theta_i(t) + e\frac{d\theta_i}{dt}(t) &= \theta_i(t) + e\frac{\omega_i(t)}{m_i}
\end{aligned}
$$

# Leapfrog method

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

$$\begin{aligned}
\omega_i(t + e/2) &= \omega_i(t) - (e/2)\frac{\partial U}{\partial \theta_i}(\theta(t)) \\
\theta_i(t + e) &= \theta_i(t) + e\frac{\omega_i(t+e/2)}{m_i} \\
\omega_i(t + e) &= \omega_i(t + e/2) - (e/2)\frac{\partial U}{\partial \theta_i}(\theta(t + e))
\end{aligned}$$

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", $\theta$, and "momentum", $\omega$, and so defines a joint distribution for them, via

$$P(\theta, \omega) = \frac{1}{Z} \exp\left(-H(\theta, \omega)\right)$$

where $Z$ is the normalizing constant.
If $H(\theta, \omega) = U(\theta) + K(\omega)$, the joint density is

$$P(\theta, \omega) = \frac{1}{Z} \exp\left(-U(\theta)\right) \exp\left(-K(\omega)\right).$$

If we are interested in a posterior distribution, we set

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

# 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

Reversibility for the leapfrog 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 $\theta$ to denote our simulation "density", then reversibility means

$$\theta(\theta^*, \omega^* | \theta, \omega) = \theta(\theta, \omega | \theta^*, \omega^*)$$

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. So long as $p(\omega) = p(-\omega)$, which is true for a multivariate normal centered at 0, this will not affect our acceptance probability.

# Conservation of Hamiltonian results in perfect acceptance

The Hamiltonian is conserved if $H(\theta, \omega) = H(\theta^*, \omega^*)$ which implies

$$\begin{aligned} p(\theta^*|y)p(\omega^*) &= \exp\left(-H(\theta^*, \omega^*)\right) \\ &= \exp\left(-H(\theta, \omega)\right) \\ &= p(\theta|y)p(\omega) \end{aligned}$$

and thus the Metropolis-Hastings acceptance probability is

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

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

```r
HMC_neal = function(U, grad_U, e, L, current_theta) {
  theta = current_theta
  omega = rnorm(length(theta),0,1)
  current_omega = omega

  omega = omega - e * grad_U(theta) / 2

  for (i in 1:L) {
    theta = theta + e * omega
    if (i!=L) omega = omega - e * grad_U(theta)
  }
  omega = omega - e * grad_U(theta) / 2

  omega = -omega

  current_U  = U(current_theta)
  current_K  = sum(current_omega^2)/2
  proposed_U = U(theta)
  proposed_K = sum(omega^2)/2

  if (runif(1) < exp(current_U-proposed_U+current_K-proposed_K))
  {
    return(theta)
  }
  else {
    return(current_theta)
  }
}
```
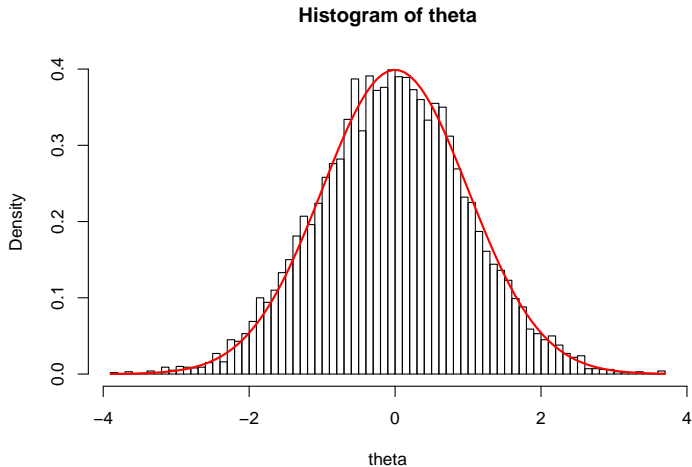
```
HMC = function(n_reps, log_density, grad_log_density, tuning, initial) {
  theta = rep(0, n_reps)
  theta[1] = initial$theta

  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
}
```

```r
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)
```



**Histogram of theta**

# Tuning parameters

There are three tuning parameters:

- $e$: 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*

```r
n_reps = 1e4
d = expand.grid(e=10^c(-3:3), L=10^seq(0,3))
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)))
})

## Error in if (runif(1) < exp(current_U - proposed_U + current_K - proposed_K)) {: missing value where
TRUE/FALSE needed
```

## Error: ggplot2 doesn't know how to deal with data of class numeric

```
## Error:  ggplot2 doesn't know how to deal with data of class numeric
```

```
## Error in if (empty(.data)) return(.data):  missing value where TRUE/FALSE needed
## Error in melt(s, id.var = c("e", "L"), variable.name = "statistic"):  object 's' not found
## Error in ggplot(m, aes(e, value, color = L, shape = as.factor(L))): object 'm' not found
```

# 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

- $e$: 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 $\mathbb{R}^d$).