

# 12: Computationally Efficient Markov chain Simulation

Taylor

University of Virginia

We mention:

- ① an example where adding auxiliary variables increases computational efficiency
- ② a few tuning tips for Random-Walk Metropolis-Hastings
- ③ Metropolis-adjusted Langevin Algorithm (MALA)
- ④ Hamiltonian Monte Carlo (HMC)
- ⑤ Pseudo-Marginal Metropolis-Hastings (PMMH).

# Example: Data Augmentation

- $y_1, \dots, y_n \mid \mu, \sigma^2 \stackrel{\text{iid}}{\sim} t_\nu(\mu, \sigma^2)$
- $\nu$  is assumed known
- $p(y_i \mid \mu, \sigma^2) \propto \left(1 + \frac{1}{\nu} \left(\frac{y_i - \mu}{\sigma}\right)^2\right)^{-(\nu+1)/2}$
- $p(\mu) \propto 1$
- $p(\sigma^2) \propto (\sigma^2)^{-1}$  (uniform for  $\log \sigma$ )

# Example: Data Augmentation

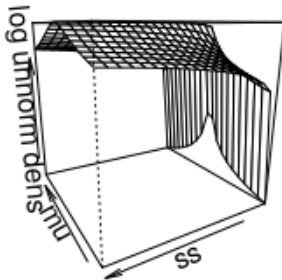
- $y_1, \dots, y_n \mid \mu, \sigma^2 \stackrel{\text{iid}}{\sim} t_\nu(\mu, \sigma^2)$
- $\nu$  is assumed known
- $p(y_i \mid \mu, \sigma^2) \propto \left(1 + \frac{1}{\nu} \left(\frac{y_i - \mu}{\sigma}\right)^2\right)^{-(\nu+1)/2}$
- $p(\mu) \propto 1$
- $p(\sigma^2) \propto (\sigma^2)^{-1}$  (uniform for  $\log \sigma$ )

Normally we would do

$$p(\mu, \sigma \mid y) \propto p(y \mid \mu, \sigma^2)p(\mu)p(\sigma^2)$$

# Example: Data Augmentation

Gibbs sampler not available :(



# Data Augmentation: Option 1

Instead, we introduce  $V_i$  (hidden/latent/unobserved data):

- $\nu$  is assumed known still
- $p(\mu) \propto 1$  still
- $p(\sigma^2) \propto (\sigma^2)^{-1}$  (uniform for  $\log \sigma$ ) still

$$p(y_i | V_i, \mu, \sigma^2) = (2\pi V_i)^{-1/2} \exp \left[ -\frac{1}{2V_i} (y_i - \mu)^2 \right] \quad (1)$$

$$p(V_i | \sigma^2) = \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \sigma^\nu V_i^{-(\nu/2+1)} \exp \left[ -\nu\sigma^2/(2V_i) \right] \quad (2)$$

Homework: show  $p(y_i | \mu, \sigma^2)$  is the same as before.

# Data Augmentation: Option 1

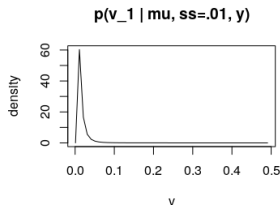
Homework question: verify all of these:

- ①  $V_i \mid \mu, \sigma^2, y \sim \text{Inv-}\chi^2 \left( \nu + 1, \frac{\nu\sigma^2 + (y_i - \mu)^2}{\nu + 1} \right)$
- ②  $\mu \mid \sigma^2, V_{1:n}, y \sim \text{Normal} \left( \frac{\sum_i \frac{1}{V_i} y_i}{\sum_i \frac{1}{V_i}}, \frac{1}{\sum_i \frac{1}{V_i}} \right)$
- ③  $\sigma^2 \mid \mu, V_{1:n}, y \sim \text{Gamma} \left( \frac{n\nu}{2}, \frac{\nu}{2} \sum_i \frac{1}{V_i} \right)$

# Data Augmentation: Option 1

Note:

$$\begin{aligned} V_i \mid \mu, \sigma^2, y &\sim \text{Inv-}\chi^2 \left( \nu + 1, \frac{\nu\sigma^2 + (y_i - \mu)^2}{\nu + 1} \right) \\ &= \text{Inv-Gamma} \left( \frac{\nu + 1}{2}, \frac{\nu\sigma^2 + (y_i - \mu)^2}{2} \right) \end{aligned}$$



(see `t_visualization.r`)

Near-zero values of  $V_i$ s lead to  $\sigma^2$  being near zero, too.



# Data Augmentation: Option 2

Add another parameter:  $\alpha > 0$

Rename a few things:

$$\tau^2 = \sigma^2 / \alpha^2 \quad (3)$$

$$U_i = V_i / \alpha^2 \quad (4)$$

Assume a noninformative prior for  $\alpha$ :

$$p(\alpha^2) \propto (\alpha^2)^{-1}$$

Page 295 has the full algorithm. Homework question: prove that the model is not identifiable in the full parameter space!

# Random Walk M-H: Some Tricks

Last section, when we were using the Metropolis-Hastings algorithm, we struggled with tuning our proposal's covariance matrix:

$$q(\theta^* \mid \theta^{t-1}) = \text{Normal}(\theta^{t-1}, \Sigma).$$

The book recommends setting

$$\Sigma \approx \frac{2.4^2}{d} \text{Var}(\theta \mid y)$$

after transforming  $\theta$  be roughly normal. Here  $d$  is the dimension of  $\theta$ . A rough approximation of the posterior covariance matrix is required.

The book also recommends aiming for an acceptance rate of about 22% for problems where  $d > 5$ .

The **Metropolis-adjusted Langevin Algorithm** is a Metropolis-Hastings algorithm with a special proposal distribution:

$$q(\theta^* \mid \theta^{t-1}) = \text{Normal}(\theta^{t-1} + \frac{\sigma_1^2}{2} \nabla \log p(\theta \mid y), \sigma_1^2 D).$$

It jumps more often in the direction of the mode, and is usually more efficient than regular Random Walk Metropolis-Hastings.

Tune  $\sigma_1^2 > 0$ ,  $D$ , hopefully have an acceptance rate around (.574) (source)

# Pseudo-Marginal Metropolis-Hastings

Sometimes you are not able to evaluate the likelihood at all!

Example 1: consider a model with a lot of hidden/missing data. Here the likelihood is a high-dimensional integral or sum that may be either impossible or computationally difficult to evaluate. Examples include state space models, factors models, “regular models” with missing data.

# Pseudo-Marginal Metropolis-Hastings

Sometimes you are not able to evaluate the likelihood at all!

Example 1: consider a model with a lot of hidden/missing data. Here the likelihood is a high-dimensional integral or sum that may be either impossible or computationally difficult to evaluate. Examples include state space models, factors models, “regular models” with missing data.

Example 2: you have a simple model, but too much data. Evaluating the density either takes too long to iterate through all of your observations, or the observations can't be read into memory. So, you can't do it several thousand times to run an MCMC sampler.

The marginal Metropolis-Hastings' acceptance ratio is

$$\frac{p(y_{1:T} \mid \theta') p(\theta')}{p(y_{1:T} \mid \theta) p(\theta)} \frac{q(\theta \mid \theta')}{q(\theta' \mid \theta)},$$

and the **pseudo-marginal Metropolis-Hastings** acceptance ratio is

$$\frac{\hat{p}(y_{1:T} \mid \theta') p(\theta')}{\hat{p}(y_{1:T} \mid \theta) p(\theta)} \frac{q(\theta \mid \theta')}{q(\theta' \mid \theta)}.$$

# Pseudo-Marginal Metropolis-Hastings

Rewrite it a little differently:

$$\frac{\overbrace{\hat{p}(y_{1:T} | u', \theta') p(\theta') \psi(u' | y_{1:T}, \theta')}^{\text{unnormalized target}}}{\underbrace{\hat{p}(y_{1:T} | u, \theta) p(\theta) \psi(u | y_{1:T}, \theta)}_{\text{proposal distribution}}} \frac{\psi(u | y_{1:T}, \theta) q(\theta | \theta')}{\psi(u' | y_{1:T}, \theta') q(\theta' | \theta)}$$

- The full normalized target is  $p(\theta, u | y_{1:T}) = \frac{\hat{p}(y_{1:T} | u, \theta) \psi(u | y_{1:T}, \theta) p(\theta | y_{1:T})}{p(y_{1:T} | \theta)}$
- When likelihood is unbiased, the marginal is  $p(\theta | y_{1:T})$
- $u$  is the collection of all auxiliary random variables (e.g. importance sampling output, a particle filter's samples and ancestor indices, etc.)

Hamiltonian Monte Carlo can be quite effective at sampling from a high-dimensional posterior. It makes use of the derivative of the log-likelihood as well.

We will describe it in three steps:

- 1 Describing Hamiltonian dynamics in continuous time
- 2 Describing how to discretize Hamiltonian dynamics
- 3 Describing how to use these in a proposal distribution in the Metropolis-Hastings algorithm.



Hamiltonian Monte Carlo can be quite effective at sampling from a high-dimensional posterior. It makes use of the derivative of the log-likelihood as well.

We will describe it in three steps:

- 1 Describing Hamiltonian dynamics in continuous time
- 2 Describing how to discretize Hamiltonian dynamics
- 3 Describing how to use these in a proposal distribution in the Metropolis-Hastings algorithm.

Video time!

- 1 <https://www.youtube.com/watch?v=mzjErXqBXw4>
- 2 <https://www.youtube.com/watch?v=87E0DKs5bok>

Say you have  $\theta_1, \dots, \theta_d$ . You add  $d$  auxiliary variables:  $\phi_1, \dots, \phi_d$ .

It's customary to use the notation  $q_1, \dots, q_d$  (the positions), and  $p_1, \dots, p_d$  (the momenta).

Say you have  $\theta_1, \dots, \theta_d$ . You add  $d$  auxiliary variables:  $\phi_1, \dots, \phi_d$ .

It's customary to use the notation  $q_1, \dots, q_d$  (the positions), and  $p_1, \dots, p_d$  (the momenta).

A HMC proposal follows a two-step procedure:

- 1 sample a random momentum vector
- 2 transform the momentum and position **nonrandomly** using Hamilton's equations

Both steps are transition kernels that preserve the stationary distribution.

# HMC: a 1-d example

Start with one-dimensional  $q$  (position) and one-dimensional  $p$  (momentum). Also,  $m$  is mass (a tuning parameter).

① potential energy:  $U(q)$  is negative the logarithm of the unnormalized posterior.

② kinetic energy:  $K(p) = \frac{p^2}{2m}$

$p = m \times \text{velocity}$

Two good resources:

① <https://arxiv.org/pdf/1206.1901.pdf> (primary reference),

② <https://arxiv.org/pdf/1701.02434.pdf>

When the particle goes up the hill, it loses kinetic energy, and gains potential energy.

Define the **Hamiltonian** as

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

and define **Hamilton's Equations** as

$$\frac{dq}{dt} = \frac{\partial H(q, p)}{\partial p} \quad (5)$$

$$\frac{dp}{dt} = -\frac{\partial H(q, p)}{\partial q} \quad (6)$$

# HMC: a first example

Assume the posterior is a Gaussian with mean  $y = 0$  and variance 1.  
Negative log of the posterior is proportional to

$$U(q) = \frac{q^2}{2}.$$

Also, assume kinetic energy is of the form

$$K(p) = \frac{p^2}{2m}.$$

# HMC: a first example

Assume the posterior is a Gaussian with mean  $y = 0$  and variance 1.  
Negative log of the posterior is proportional to

$$U(q) = \frac{q^2}{2}.$$

Also, assume kinetic energy is of the form

$$K(p) = \frac{p^2}{2m}.$$

so

$$\frac{dq}{dt} = \frac{\partial H(q, p)}{\partial p} = \frac{dK(p)}{dp} = p(t)/m \quad (7)$$

$$\frac{dp}{dt} = -\frac{\partial H(q, p)}{\partial q} = -\frac{dU(q)}{dq} = -q(t) \quad (8)$$

# HMC: a first example

If  $m = 1$ , a solution to

$$\frac{dq}{dt} = \frac{\partial H(q, p)}{\partial p} = \frac{dK(p)}{dp} = p(t) \quad (9)$$

$$\frac{dp}{dt} = -\frac{\partial H(q, p)}{\partial q} = -\frac{dU(q)}{dq} = -q(t) \quad (10)$$

is

$$q(t) = r \cos(a + t) \quad (11)$$

$$p(t) = -r \sin(a + t) \quad (12)$$



# HMC: a first example

For this particular model,  $(q, p)'$  rotates clockwise in phase-space.

$$q(t) = r \cos(a + t) \quad (13)$$

$$p(t) = -r \sin(a + t) \quad (14)$$

Can be written as

$$\begin{bmatrix} r \cos(a + t) \\ -r \sin(a + t) \end{bmatrix} = \underbrace{\begin{bmatrix} \cos([t - s]) & \sin([t - s]) \\ \sin([t - s]) & \cos([t - s]) \end{bmatrix}}_{T_{t-s}} \begin{bmatrix} r \cos(a + s) \\ -r \sin(a + s) \end{bmatrix}$$

(use Angle-Sum trig identity)

# HMC: a first example

When  $m \neq 1$ ,  $T_s$  might look like this.

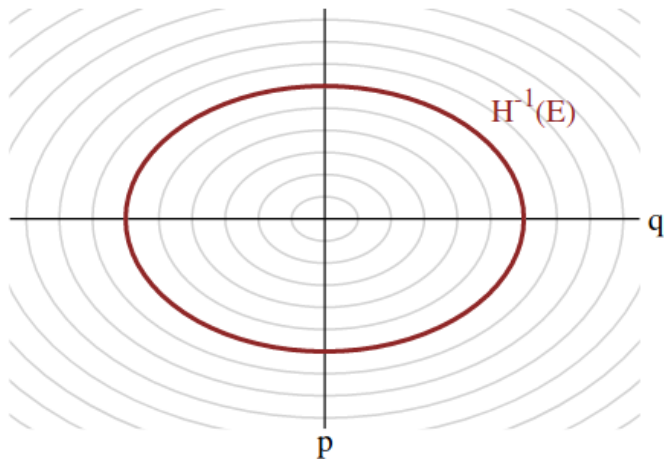


image source: <https://arxiv.org/pdf/1701.02434.pdf>

Now assume  $d$ -dimensional posterior  $\mathbf{q} = (q_1, \dots, q_d)$  and  $\mathbf{p} = (p_1, \dots, p_d)$ .

When

$$K(\mathbf{p}) = \frac{\mathbf{p}' M^{-1} \mathbf{p}}{2}$$

Hamilton's equations become

$$\frac{dq_i}{dt} = \frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial p_i} = \frac{dK(\mathbf{p})}{dp_i} = [M^{-1} \mathbf{p}]_i \quad (15)$$

$$\frac{dp_i}{dt} = -\frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial q_i} = -\frac{dU(\mathbf{q})}{dq_i} \quad (16)$$

Recall that  $U(\mathbf{q})$  is the negative log of the posterior you're interested in.

# Property 1: Reversibility of $T_s$

$T_s : [\mathbf{q}(0), \mathbf{p}(0)] \mapsto [\mathbf{q}(s), \mathbf{p}(s)]$  always has an easy-to-find inverse.

Proof: just take the negative of the derivatives.

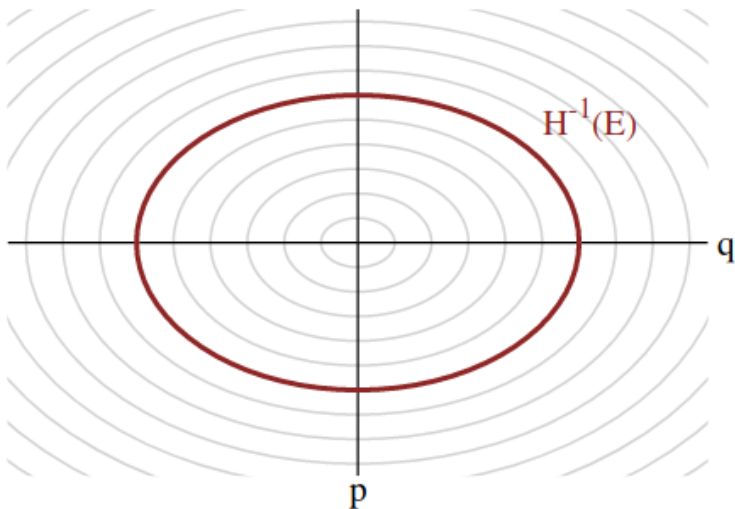
## Property 2: Conservation of the Hamiltonian

Using the chain rule:

$$\begin{aligned}\frac{dH}{dt} &= \sum_{i=1}^d \frac{\partial H}{\partial p_i} \frac{dp_i}{dt} + \sum_{i=1}^d \frac{\partial H}{\partial q_i} \frac{dq_i}{dt} \\ &= \sum_{i=1}^d \frac{dK}{dp_i} \frac{dp_i}{dt} + \sum_{i=1}^d \frac{dU}{dq_i} \frac{dq_i}{dt} \\ &= \sum_{i=1}^d \frac{dK}{dp_i} \left( -\frac{dU}{dq_i} \right) + \sum_{i=1}^d \frac{dU}{dq_i} \frac{dK}{dp_i} \\ &= 0\end{aligned}$$

Moving through time keeps you on the same contour or level-set in the phase space.

$T_S$  keeps you on a level-set/contour:



# HMC Property 3 and 4: Volume Preservation and Symplecticness

Volume preservation:  $\{(q, p) : (q, p) \in A\}$  and  $\{T_s(q, p) : (q, p) \in A\}$  have the same volume.

Symplecticness: a nice property of the Jacobian (matrix of time derivatives) of  $T_s$ .

Another thing: when we approximate these dynamics in our proposal distribution, these properties are preserved!

# HMC: looking back at the big picture

HMC will work as follows: given that we are currently at position  $\mathbf{q}(t)$ , we are going to sample a momentum vector (which puts us on one of the level-sets), and then we are going to follow  $T_s$  for a deterministic amount of time (how much time is a tuning parameter we decide on).

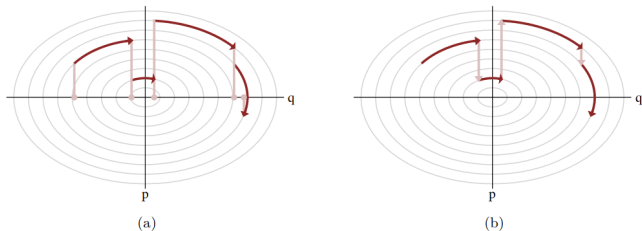


FIG 22. (a) Each Hamiltonian Markov transition lifts the initial state onto a random level set of the Hamiltonian, which can then be explored with a Hamiltonian trajectory before projecting back down to the target parameter space. (b) If we consider the projection and random lift steps as a single momentum resampling step, then the Hamiltonian Markov chain alternates between deterministic trajectories along these level sets (dark red) and a random walk across the level sets (light red).



# HMC: looking back at the big picture

HMC will work as follows: given that we are currently at position  $\mathbf{q}(t)$ , we are going to sample a momentum vector (which puts us on one of the level-sets), and then we are going to follow  $T_s$  for a deterministic amount of time (how much time is a tuning parameter we decide on).

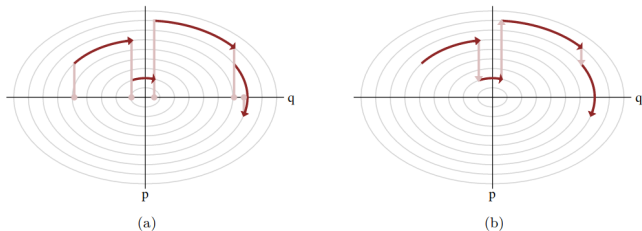


FIG 22. (a) Each Hamiltonian Markov transition lifts the initial state onto a random level set of the Hamiltonian, which can then be explored with a Hamiltonian trajectory before projecting back down to the target parameter space. (b) If we consider the projection and random lift steps as a single momentum resampling step, then the Hamiltonian Markov chain alternates between deterministic trajectories along these level sets (dark red) and a random walk across the level sets (light red).

Following a contour line is impossible in continuous time though...

# Discretizing Hamilton's Equations: Version 1.0

We need to be able to approximate  $T_S$  using the derivatives. To do that, we pick a small change in time called  $\epsilon$ . Then we take  $L$  steps of size  $\epsilon$ .

Two procedures are described. The last one is the one that is most commonly used.

For simplicity, assume the mass matrix is diagonal, making

$$K(\mathbf{p}) = \mathbf{p}' M^{-1} \mathbf{p} = \sum_{i=1}^d \frac{p_i^2}{2m_i}.$$

# Discretizing Hamilton's Equations: Version 1.0

When  $K(\mathbf{p}) = \sum_{i=1}^d \frac{p_i^2}{2m_i}$ , **Euler's method** approximates the solution of

$$\frac{dq_i}{dt} = \frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial p_i} = \frac{dK(\mathbf{p})}{dp_i} = \frac{p_i}{m_i} \quad (17)$$

$$\frac{dp_i}{dt} = -\frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial q_i} = -\frac{dU(\mathbf{q})}{dq_i} \quad (18)$$

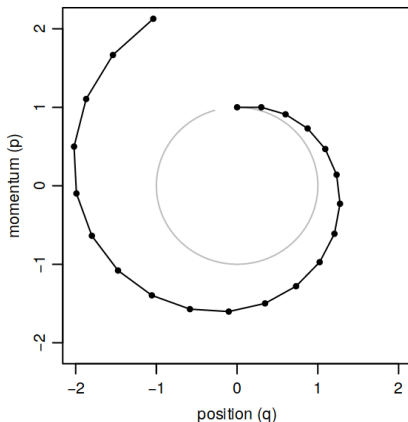
as

$$q_i(t + \epsilon) = q_i(t) + \epsilon \frac{p_i(t)}{m_i} \quad (19)$$

$$p_i(t + \epsilon) = p_i(t) - \epsilon \frac{dU(\mathbf{q}(t))}{dq_i} \quad (20)$$

# Discretizing Hamilton's Equations: Version 1.0

(a) Euler's Method, stepsize 0.3



Twenty steps when  $H(q, p) = p^2/2 + q^2/2$ , the initial state is  $(q, p) = (0, 1)$ .

# Discretizing Hamilton's Equations: Version 2.0

When  $K(\mathbf{p}) = \sum_{i=1}^d \frac{p_i^2}{2m_i}$ , **the leap-frog method** approximates

$$\frac{dq_i}{dt} = \frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial p_i} = \frac{dK(\mathbf{p})}{dp_i} = p_i/m_i \quad (21)$$

$$\frac{dp_i}{dt} = -\frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial q_i} = -\frac{dU(\mathbf{q})}{dq_i} \quad (22)$$

with

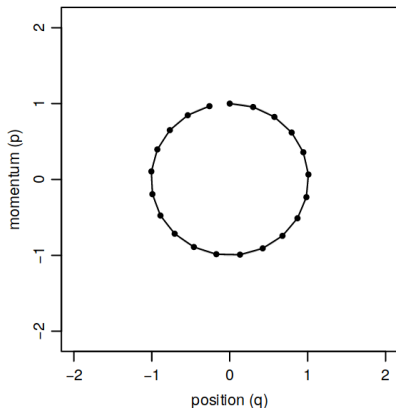
$$p_i(t + \epsilon/2) = p_i(t) - (\epsilon/2) \frac{dU(\mathbf{q}(t))}{dq_i} \quad (23)$$

$$q_i(t + \epsilon) = q_i(t) + \epsilon \frac{p_i(t + \epsilon/2)}{m_i} \quad (24)$$

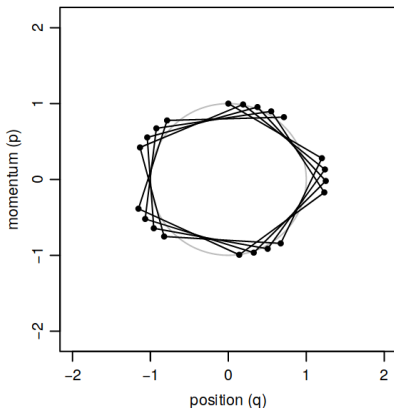
$$p_i(t + \epsilon) = p_i(t + \epsilon/2) - (\epsilon/2) \frac{dU(\mathbf{q}(t + \epsilon))}{dq_i} \quad (25)$$

# Discretizing Hamilton's Equations: Version 2.0

(c) Leapfrog Method, stepsize 0.3



(d) Leapfrog Method, stepsize 1.2



Twenty steps when  $H(q, p) = p^2/2 + q^2/2$ , the initial state is  $(q, p) = (0, 1)$ .

# Describing the HMC algorithm

The algorithm targets the distribution for  $(\mathbf{q}, \mathbf{p})$ :

$$\begin{aligned}\frac{1}{Z} \exp \left[ -\frac{H(\mathbf{q}, \mathbf{p})}{T} \right] &= \frac{1}{Z} \exp \left[ -\frac{K(\mathbf{p}) + U(\mathbf{q})}{T} \right] \\ &= \frac{1}{Z} \exp \left[ -\frac{K(\mathbf{p})}{T} \right] \exp \left[ -\frac{U(\mathbf{q})}{T} \right] \\ &= \frac{1}{Z} \exp \left[ -\frac{K(\mathbf{p})}{T} \right] \times \\ &\quad \exp \left[ -\frac{-\log \{ \text{prior}(\mathbf{q}) \times \text{likelihood}(\mathbf{q}) \}}{T} \right]\end{aligned}$$

# Describing the HMC algorithm

Step 1:

Sample  $p$  from the conditional target distribution

$$\frac{1}{Z} \exp \left[ -\frac{K(\mathbf{p})}{T} \right].$$

In our case, this is the same as the marginal, due to independence.

Notice how this is a Gibbs-like step! It preserves the stationary distribution, and it has 100% chance of being accepted.



# Describing the HMC algorithm

## Step 2:

If we could integrate Hamilton's equations exactly, then our proposal would be deterministic, and we would accept with probability 1, because the Hamiltonian is preserved (property 2).

However, because we are using numerical leap-frog integration, there will be some change in the Hamiltonian. We think of the  $L$  leap-frog steps as a proposal distribution. This is a deterministic proposal, and it's symmetrical (we don't prove this). So what we end up with is a Metropolis-like acceptance probability:

$$\min \left[ 1, \frac{\exp[-H(\mathbf{q}^*, \mathbf{p}^*)]}{\exp[-H(\mathbf{q}, \mathbf{p})]} \right]$$

# Describing the HMC algorithm

Code from <https://arxiv.org/pdf/1206.1901.pdf> that performs one iteration of HMC can be found in the file `hmc.r`.

Here's a visualization:

<https://chi-feng.github.io/mcmc-demo/>