

## Lecture 10:

Topics:

- I) Hamiltonian Monte carlo
- II) Look ahead HMC.

Suppose  $d$  is some data vector drawn from some density  $\rho(d|\theta)$  where  $\theta = (\theta_1, \dots, \theta_n)$  is a parameter vector. Let  $\pi(\theta)$  denote a prior density for  $\theta$ .

Goal: Produce posterior samples

$$\theta^{(1)}, \dots, \theta^{(n)} \sim \rho(\theta|d) = \frac{\rho(d|\theta)\pi(\theta)}{\rho(d)}.$$

Problem: if the dimension  $n$  of  $\theta \in \mathbb{R}^n$  is large then Metropolis-Hastings chains typically have low acceptance rates.

## Hamiltonian Monte Carlo (Hmc)

... an MCMC chain which attempts to solve the low acceptance rate problem when  $\theta$  is high dimensional.

The main requirement for this method is the ability to compute

$\nabla_\theta \log \rho(\theta|d)$  &  $\rho(\theta|d)$  quickly (where  $\rho(\theta|d)$  only needs to be computed up to a unknown normalizing factor).

(1)

To define the HMC sample, start by defining a momentum vector

$$p := (p_1, \dots, p_n) \in \mathbb{R}^n$$

with corresponding mass

$$m := (m_1, \dots, m_n) \in \mathbb{R}^n \text{ s.t. } m_i > 0.$$

Now define the Hamiltonian as follows

$$\begin{aligned} H(\theta, p) &:= -\log \rho(\theta|d) + \sum_{i=1}^n \frac{|p_i|^2}{2m_i} \\ &= -\log \rho(\theta|d) - \log \pi(p) + c \\ &= -\log [\rho(\theta|d)\pi(p)] + c \end{aligned}$$

where  $\pi(p)$  is the density of  $p \sim N(0, \text{diag}(m))$ .

You can think of  $-\log \pi(p)$  as the kinetic energy of the system  $(\theta, p)$  and  $-\log \rho(\theta|d)$  as potential energy. If we are thinking about  $(\theta, p)$  as describing a physical system at time  $t=0$ , then we can evolve the system, producing  $(\theta^t, p^t)$  for  $t \in [0, \infty)$ , with Hamiltonian dynamics:

$$(H1) \quad \frac{d\theta^t}{dt} = \nabla_p H(\theta^t, p^t)$$

$$(H2) \quad \frac{dp^t}{dt} = -\nabla_\theta H(\theta^t, p^t).$$

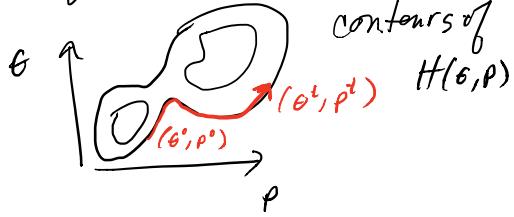
Under this evolution  $H(\theta^t, p^t)$  is preserved

$$H(\theta^t, p^t) = H(\theta^0, p^0) \quad \forall t \in [0, \infty) \quad \text{since}$$

$$\begin{aligned} \frac{d}{dt} H(\theta^t, p^t) &= \nabla_\theta H(\theta^t, p^t) \cdot \frac{d\theta^t}{dt} + \nabla_p H(\theta^t, p^t) \cdot \frac{dp^t}{dt} \\ &= 0. \end{aligned}$$

(2)

$\therefore (\theta^t, p^t)$  moves along a contour of the log likelihood given by  $H(\theta, p) + c$



One can also show that (H1) & (H2) preserve volume in "phase space"  $(\theta, p) \in \mathbb{R}^n \times \mathbb{R}^n$  meaning that for any set  $\mathcal{R} \subset \mathbb{R}^n \times \mathbb{R}^n$  one has

$$\text{vol}(\mathcal{R}) = \text{vol}(\mathcal{H}^t(\mathcal{R}))$$

where

$$\mathcal{H}^t(\mathcal{R}) := \left\{ \mathcal{H}^t(\theta, p) \in \mathbb{R}^n \times \mathbb{R}^n \mid (\theta, p) \in \mathcal{R} \right\}$$

$\mathcal{H}^t(\theta, p) = (\theta^t, p^t)$  evolved by (H1)(H2) with initial condition

$$(\theta^0, p^0) := (\theta, p).$$

To see this simply note that  $\mathcal{H}^t(\theta, p)$  is a div free evolution.

i.e. that

$$\begin{aligned} \text{div}\left((\nabla_p H, \nabla_\theta H)^T\right) &= \sum_{i=1}^n \frac{\partial^2}{\partial \theta_i \partial p_i} H(\theta, p) \\ &\quad - \sum_{i=1}^n \frac{\partial^2}{\partial p_i \partial \theta_i} H(\theta, p) \\ &= 0 \end{aligned}$$

why is this useful? (4)

Let  $X = (\theta^0, p^0)$  be a random vector with some density  $p_X(x)$ . Transform  $X \mapsto Y$  as follows

$$Y = \mathcal{H}^t(X) = (\theta^t, p^t) \text{ under (H1)(H2)}$$

The change of variables formula says

$$p_X(x) dx = p_Y(y) dy$$

$$\therefore p_X(x) = p_Y(y) \left| \frac{\partial y}{\partial x} \right|$$

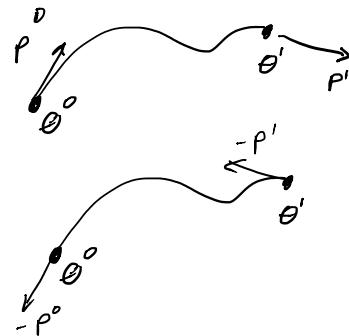
↓ jacobian

$$= p_Y(\mathcal{H}^t(x)) \underbrace{\left| \frac{\partial \mathcal{H}^t(x)}{\partial x} \right|}_{=1 \text{ by vol}} \text{ preserving property}$$

$$\text{i.e. } p_X(\theta^0, p^0) = p_Y(\theta^t, p^t)$$

Finally we note that (H1) & (H2) is reversible in that

flipping the sign on momentum  $p$  runs the chain in reverse.



Let  $s = (\theta, p)$  denote general state vectors. Let  $L$  be the operator on state vectors defined by evolving  $(H_1) \& (H_2)$  to some fixed time  $t$ .

$$\therefore Ls \equiv L(\theta, p)^T := H^t(\theta, p)$$

Also let  $F$  denote the "momentum sign flip operator" so that

$$Fs \equiv F(\theta, p)^T := (\theta, -p).$$

Note that our previous results on  $(H_1) \& (H_2)$  imply

$$L^{-1}s = F L F s \quad (\text{reversability})$$

$$\left| \frac{\partial Ls}{\partial s} \right| = 1 \quad (\text{vol preserving})$$

$$H(Ls) = H(s) \quad (\text{conservation of energy})$$

we also have

$$F^{-1}s = Fs$$

$$\left| \frac{\partial Fs}{\partial s} \right| = 1 \quad (\text{vol preserving})$$

$$H(Fs) = H(s).$$

Finally define the "momentum Randomizing operator"  $R_\beta$  (parameter  $\beta \in [0, 1]$ ) as follows

$$R_\beta s \equiv R_\beta(\theta, p)$$

$$:= (\theta, p\sqrt{1-\beta} + \sqrt{\beta} p^*)$$

$$p^* \sim N(0, \text{diag}(m))$$

Now we can construct a markov chain  $s^{(1)}, s^{(2)}, \dots$  for state vectors  $s^{(i)} \equiv (e^{(i)}, p^{(i)})$  as follows.

### Hamiltonian Markov Chain

for  $i = 1, 2, \dots$

$$\alpha := \min\left(\frac{\exp(-H(Fs^{(i)}))}{\exp(-H(s^{(i)}))}, 1\right)$$

$$\tilde{s} := \begin{cases} Fs^{(i)} & \text{if } U \leq \alpha \\ s^{(i)} & \text{if } U > \alpha \end{cases}$$

$$s^{(i+1)} := R_\beta F \tilde{s}$$

end

The key for showing this works (i.e. has invariant density  $e^{-H(s)} = e^{-H(\theta, p)}$ ) is the Reversability of  $L$  & the vol preserving of  $L$  &  $F$ .

This is basically a proposal.

then an additional random walk step

In this form it is hard to  
see what is going on. We can  
simplify the above algorithm  
when  $\beta = 1$ . Notice the momentum

Components of  $s^{(i)}$  are eventually discarded so we can write the algorithm in terms of  $\theta^{(i)}$  alone.

## Hamiltonian Monte Carlo Chain $\beta=1$

for  $i=1, 2, \dots$   
 simulate  $\rho^* \sim N(0, \text{diag}(m))$

$$\tilde{\zeta} := (\theta^{(i)}, \rho^*)$$

$$\theta^{(i+1)} := \begin{cases} (\tilde{L}\tilde{S})_\theta & \text{if } u \leq \alpha \\ \tilde{S}_\theta & \text{if } u > \alpha. \end{cases}$$

*end*

where  $S_0$  extracts the parameter coordinates from  $S$ .

This is basically RBF's from the previous step

Using the exact  
Hamiltonian evolution operator  
 $L$ , since  $H(L\tilde{s}) = H(\tilde{s})$   
 $\alpha \beta$  always 1 so this  
MCMC always accepts the proposal.

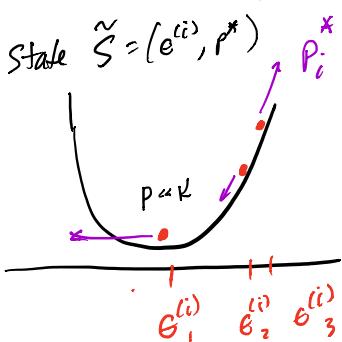
To get an idea of what is going on expand out the Hamiltonian (8)

$$H(\tilde{s}) = H(\theta^{(i)}, p^*)$$

$$= -\log P(G^{(i)} | d) + \sum_{j=1}^n \frac{(P_i^{(k)})^2}{z m_i}$$

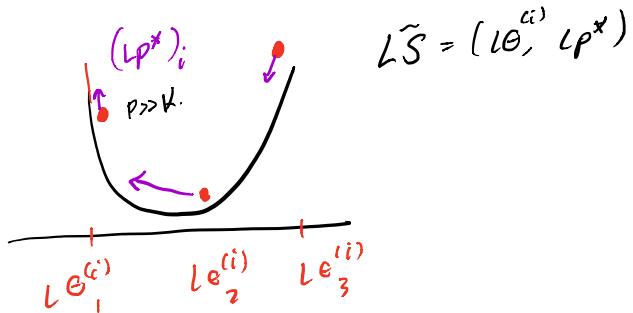
Potential energy ( $T$ )      Kinetic energy ( $K$ )      This will be  
 $\uparrow$                            $\uparrow$                            $\curvearrowright$

is how  $-\log P(\text{child})$   
 show approximately  
 vary when  
 $\text{on } P(\text{child})$



The Hamiltonian is evolved with  $\mathcal{L}$  which where Kinetic & potential energy are exchanged.

$$H(LS) = -\log P(LG^{(i)} | d) + \sum_{i=1}^n \frac{(LP_i^*)^2}{2m_i}$$



This changes the potential energy  $-\log p(\theta|d)$  by an amount expected from samples  $\theta \sim p(\theta|d)$

## Showing detailed balance

(9)

Consider a Markov chain  $X^{(1)}, X^{(2)}, \dots$

Let  $T$  be a random operator s.t. the  $i+1^{\text{st}}$  step can be written

$$X^{(i+1)} = T X^{(i)}$$

$\nwarrow$  the randomness is independently drawn for each  $i$ .

The notion of detailed balance for transition operator  $T$  & a measure  $\pi$  is simply that when  $X \sim \pi$  then  $(X, TX) \stackrel{D}{=} (TX, X)$ .

This clearly implies  $X \stackrel{D}{=} TX$  (when  $X \sim \pi$ ) so detailed balance implies  $\pi$  is an invariant measure for the chain.

Notice that if  $T = T_N \dots T_1$ , i.e. there are  $N$  intermediate steps in each transition  $X^{(i)} \rightarrow X^{(i+1)}$ :

$$X^{(i)} \rightarrow T_1 X^{(i)} \rightarrow T_2 T_1 X^{(i)} \rightarrow \dots \rightarrow T X^{(i)} = X^{(i+1)}$$

Then to prove  $\pi$  is an invariant measure it is sufficient to show each  $T_i$  satisfies detailed balance w.r.t  $\pi$  (follows since each application of  $T_i$  is invariant to  $\pi$ ).

For the HMC each transition

$s^{(i)} \rightarrow s^{(i+1)}$  is composed of two steps

$$s^{(i)} \rightarrow \tilde{s} \rightarrow s^{(i+1)}$$

write this abstractly as

$$s \rightarrow TS \rightarrow R_F T s$$

The goal of this section is to illustrate part of the proof that  $T$  satisfies detailed balance for the measure  $e^{-H(s)}$  and show where the vol preserving nature of  $F$  &  $L$  and reversibility  $L^{-1} = FLF$  are used.

Let  $s \sim e^{-H(s)}$  & consider the two joint R.V.s

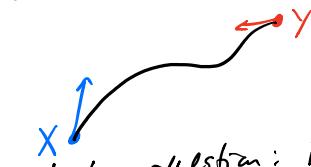
$(s, TS)$  &  $(TS, s)$

we wish to show they have the same distribution where

$$TS := \begin{cases} FLS \text{ with prob } \alpha(FLS|s) \\ s \text{ with prob } 1 - \alpha(FLS|s) \end{cases}$$

$$\alpha(x|y) := \min\left(\frac{e^{-H(x)}}{e^{-H(y)}}, 1\right)$$

Think of proving detail balance as a likelihood inference problem  
Given two states  $X, Y$  in phase space



we ask the question: is it equally likely that

$$\begin{array}{ll} s = X & s = Y \\ TS = Y & TS = X \end{array}$$

assuming  $s \sim e^{-H(s)}$ .

To satisfy detailed balance you need any likelihood discrepancy b/wn

$$e^{-H(X)} \text{ & } e^{-H(Y)}$$

Need to be canceled out by the likelihoods of the transitions

$$T(X) = Y \text{ & } T(Y) = X$$

Let  $P$  be the prob measure for  $(X, Y)$  under the assumption  $(X, Y) := (S, TS)$  when  $S \sim e^{-H(S)}$ . We wish to show

$$P(X \in A, Y \in B) = P(Y \in A, X \in B).$$

We will only do part of the proof, assuming  $A \cap B = \emptyset$ . This implies  $TS$  must have resulted in  $TS = FLS$

$$\begin{aligned} \therefore P(X \in A, Y \in B) &\leftarrow \text{if } A \cap B = \emptyset \text{ then} \\ &T = F^L \text{ &} \\ &Y = F^L X \\ &= P(X \in A, F^L X \in B, T = F^L) \\ &= P(X \in A \cap (F^L)^{-1} B, T = F^L) \\ &\quad \text{since} \\ &\quad (F^L)^{-1} = L^{-1} F \\ &\quad = F L F F = L \\ &\quad \text{by reversability} \\ &= P(X \in A \cap F^L B, T = F^L) \\ &= \int_{A \cap F^L B} \alpha(F^L x | x) e^{-H(x)} dx \\ &= P(T = F^L | X = x) \\ &= \int_{A \cap F^L B} \min(e^{-H(F^L x)}, e^{-H(x)}) dx \\ &= \int_{FLA \cap B} \alpha(F^L y | y) e^{-H(y)} dy \\ &= P(X \in B, Y \in A) \\ &= P(Y \in A, X \in B) \quad \square \end{aligned}$$

*change variables  
 $y = F^L x$   
 $dy = dx$   
by vol preserving of  
 $F^L$*

## Leap frog discretizing $L$

(12)

Recall that  $L$  is defined to be the operator on state vectors  $s$  which evolves  $(H_1) \& (H_2)$  to some fixed time  $t$ .

The key properties of reversibility & vol preserving establish that the HMC algorithm has invariant density given by  $e^{-H(s)}$ .

It turns out that there exists a discrete approximation to  $L$ , denoted by  $L^{\varepsilon, m}$ , that still satisfies reversibility & vol preservation, but for which  $H(s)$  is not preserved (so  $H(s) \neq H(L^{\varepsilon, m} s)$ ) as is the case for  $L$ .

$L^{\varepsilon, m}$  is called leapfrog, or Störmer-Verlet, integration

and is defined as follows

function  $L^{\varepsilon, m}(s)$

set  $s^0 := s$  and  $(\theta^0, p^0) = s^0$ .

set  $\lambda = \text{diag}(m_1, \dots, m_n)$

for  $k = 0, 1, \dots, M$

$$\theta^{k+1} := \theta^k + \varepsilon \lambda^{-1} \left[ p^k - \frac{\varepsilon}{2} \nabla_\theta H(\theta^k, p^k) \right]$$

$$p^{k+1} := p^k - \frac{\varepsilon}{2} \left[ \nabla_\theta H(\theta^k, p^k) + \nabla_\theta H(\theta^{k+1}, p^k) \right]$$

end

$$\text{return } (\theta^{M+1}, p^{M+1})$$

end

Note: the updates for  $\theta^{k+1}, p^{k+1}$  are usually written as

(13)

$$\tilde{p} := p_k - \frac{\varepsilon}{2} \nabla_{\theta} H(\theta^k, p^k)$$

$$\theta^{k+1} := \theta^k + \varepsilon \nabla_p H(\theta^k, \tilde{p})$$

$$p^{k+1} := \tilde{p} - \frac{\varepsilon}{2} \nabla_{\theta} H(\theta^{k+1}, \tilde{p})$$

which more closely resemble discrete Hamiltonian updates (but are exactly equivalent to the algorithm given).

Now the Discrete HMC simply replaces  $L$  with  $L^{\varepsilon, M}$ .

The main difference using  $L^{\varepsilon, M}$  in place of  $L$  is that acceptance rate  $\alpha$  will no longer be exactly 1 since the Hamiltonian is only approximately preserved.

Indeed,  $\varepsilon$  &  $M$  play the role of tuning parameters where larger  $\varepsilon$  and/or larger  $M$  will result in lower acceptance rates ... but give larger moves when accepted.

One of the main drawbacks with HMC is that it can be difficult to tune.

## Tuning Parameters

(14)

\*  $\beta > 0$  can be useful (see Culpepper 2011)

\* The typical suggestion is to set  $\text{diag}(m)$  [or a general  $\Lambda$  where  $p \sim N(0, \Lambda)$ ] to the inverse of the posterior Cov Matrix of  $\theta$ .

\* Then  $\varepsilon$  &  $M$  are chosen so that updating with  $L^{\varepsilon, M}$  gives good acceptance.

## CMB application

Taylor, Ashdown & Hobson used HMC to sample from  $\underbrace{P(T, C_e^{\text{TT}} | d)}_{= P(\theta | d)}$  using the fact that

$$-\nabla_{\theta} \log P(\theta | d) = \begin{pmatrix} -\nabla_{T(\hat{n})} \log P(T, C_e^{\text{TT}} | d) \\ -\nabla_{C_e^{\text{TT}}} \log P(T, C_e^{\text{TT}} | d) \end{pmatrix}$$

where

$$-\nabla_{T(\hat{n})} \log P(T, C_e^{\text{TT}} | d) = -(\bar{\Sigma} + N)^{-1} (d - \bar{T}) + \bar{\Sigma} \bar{T}$$

$$-\nabla_{C_e^{\text{TT}}} \log P(T, C_e^{\text{TT}} | d)$$

$$= \left( \frac{2\ell+1}{2} \right) \frac{1}{C_e^{\text{TT}}} \left( 1 - \frac{\bar{\Sigma} e^T}{C_e^{\text{TT}}} \right)$$

when  $d_{em} = T_{em} + n_{em}$ , one uses a uniform prior on  $C_e^{TT}$  &

$$\sigma_e^T := \frac{1}{2\ell+1} \sum_{m=-\ell}^{\ell} \|T_{em}\|^2.$$

### Look Ahead Hamiltonian Monte Carlo (LAHMC)

The detailed balance property of HMC suggest that even though one gets high acceptance rates there is still random walk behavior in HMC leading to slow mixing.

Sohl-Dickstein et al. developed a similar chain to HMC with the same stationary distribution but which doesn't satisfy detailed balance & can have a faster mixing rate.

(13)

LAHMC Produces  $s^{(1)}, s^{(2)}, \dots$  in phase space (14)

Fix  $K \in \{1, 2, \dots\}$ ,  $0 \leq \beta \leq 1$  and let  $L$  denote the leapfrog operator  $L^{e,m}$  where  $L^j := \underbrace{L L \dots L}_{j \text{ times}}$ . Define

$$\alpha_1(s) := \min \left[ \frac{e^{-H(FLS)}}{e^{-H(s)}}, 1 \right]$$

$$\alpha_k(s) := \min \left[ \frac{e^{-H(FL^k s)}}{e^{-H(s)}}, \left( 1 - \sum_{i=1}^{k-1} \alpha_i(FL^k s) \right), 1 - \sum_{i=1}^{k-1} \alpha_i(s) \right]$$

for  $k = 2, \dots, K$  and

$$\alpha_{K+1}(s) := 1 - \sum_{k=1}^K \alpha_k(s).$$

for  $i = 1, 2, \dots$

$$\tilde{s} := \begin{cases} FL's^{(i)} & \text{with prob } \alpha_1(s^{(i)}) \\ FL^2 s^{(i)} & \text{with prob } \alpha_2(s^{(i)}) \\ \vdots & \vdots \\ FL^K s^{(i)} & \text{with prob } \alpha_K(s^{(i)}) \\ s^{(i)} & \text{with prob } \alpha_{K+1}(s^{(i)}) \end{cases}$$

$$s^{(i+1)} = R_\beta F \tilde{s}$$

and

To see why LAHMC works lets again take a look at the transitions

$$s^{(i)} \xrightarrow{T} \tilde{s}$$

and show  $e^{-H(s)}$  is an invariant measure.

So suppose  $s \sim e^{-H(s)}$

$$\therefore TS = \begin{cases} FL's & \text{w/prob } \alpha_1(s) \\ \vdots & \vdots \\ FL^K s & \text{w/prob } \alpha_K(s) \\ s & \text{w/prob } \alpha_{K+1}(s) \end{cases}$$

(16)

$$\begin{aligned}
 \therefore P(Ts \in A) &= \int \underbrace{P(Ts \in A | s)}_{\sum_{k=1}^K P(F_L^k s \in A | s) \alpha_k(s)} e^{-H(s)} ds \\
 &= \sum_{k=1}^K P(F_L^k s \in A | s) \alpha_k(s) \\
 &\quad + P(F_S \in A | s) \alpha_{K+1}(s) \\
 &= \sum_{k=1}^K \int_{(F_L^k)^{-1}A} \alpha_k(s) e^{-H(s)} ds + \int_A \alpha_{K+1}(s) e^{-H(s)} ds \\
 &= \sum_{k=1}^K \int_A \alpha_k(F_L^k s) e^{-H(F_L^k s)} ds + \int_A \alpha_{K+1}(s) e^{-H(s)} ds
 \end{aligned}$$

Since  $(F_L^k)^{-1} = F_L^{-k}$  is vol preserving

Now

$$\begin{aligned}
 \alpha_k(F_L^k s) e^{-H(F_L^k s)} &= s \\
 &= \min \left[ e^{-H(s)} \left( 1 - \sum_{i=1}^{k-1} \alpha_i(F_L^k F_L^i s) \right), \right. \\
 &\quad \left. e^{-H(F_L^k s)} \left( 1 - \sum_{i=1}^{k-1} \alpha_i(F_L^k s) \right) \right] \\
 &= \alpha_k(s) e^{-H(s)}
 \end{aligned}$$

and

$$\begin{aligned}
 \therefore P(Ts \in A) &= \sum_{k=1}^K \int_A \alpha_k(s) e^{-H(s)} ds \\
 &\quad + \int_A \alpha_{K+1}(s) e^{-H(s)} ds \\
 &= \int_A \underbrace{(\alpha_1(s) + \dots + \alpha_{K+1}(s))}_{=1} e^{-H(s)} ds
 \end{aligned}$$

$\therefore Ts \sim e^{-H(s)}$  as was to be shown.