

# MCMC & HMC

Assume BNN model:  $p(w) = N(0, I)$   
 $p(y|x, w) = \prod_n N(y_n | \theta(x_n, w), \Sigma)$

Goal: Draw samples  $w^s \sim p(w|x, y)$  for our BNN

Using these samples, can compute  $f(x_*, w^s)$  for each sample,

goal  $\mathbb{E}_{p(w|x, y)} [f(x_*, w)]$   
 $\text{Var}[f(x_*, w)]$  } estimate via empirical methods on  $S$  samples

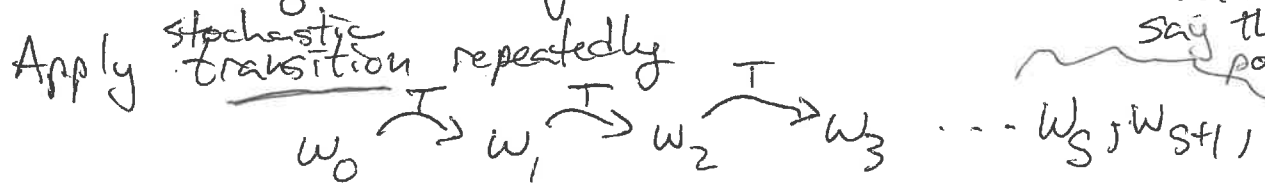
Common problem in Bayesian inference:

sample from posterior distribution given model and data  
 ↳ likelihood density  
 ↳ prior density

we only know posterior up to constant density

$$p(w|x, y) \propto \underbrace{p(y|x)}_{\text{unknown}} \underbrace{p(y|x, w)}_{\text{lik}} \underbrace{p(w)}_{\text{prior}}$$

Idea:  $w_0$  starting value



if we choose transition well, can we say these match posterior samples

# What is MCMC?

02

This is a Markov Chain b/c  $w_t$  only depends on  $w_{t-1}$

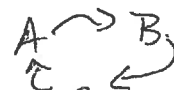
knowing  $w_{t-2}, \dots, w_1$  doesn't tell us any more about  $p(w_t)$

This is Monte Carlo b/c transitions are random/stochastic not deterministic.

Markov chain theory says that if a transition operator is

(1) irreducible: good at exploring pos. proba. of visiting all states

(2) aperiodic: no cycles possible



then after many applications of  $T$ , we will reach a stationary distrib<sup>n</sup>.

Goal of MCMC: construct  $T$  so stationary distr. is the posterior!

How? One way: make  $T$  satisfy detailed balance

$$p(w_t) T(w_{t+1} | w_t) = p(w_{t+1}) T(w_t | w_{t+1})$$

can show this joint distr. leaves  $p(w_{next})$  equal to  $p(w_{current})$ . transitions produce draw for stationary mean cover  $\frac{1}{2}$

Metropolis algorithm

Given  $w_0$

for  $t$  in  $1, 2, \dots$   
 $w^{prop} \sim N(\text{mean}, \sigma^2)$

accept-prob  $\leftarrow \min(1,$

if  $\text{rand}() < \text{accept-prob}$ :

$w^t \leftarrow w^{prop}$

else

$w^t \leftarrow w^{t-1}$

$T = \text{Normal}(w_t | w_{t-1}, \Sigma)$

Symmetric ~~proposals~~  
transitions  
proposals

$$\frac{p(w^{prop})}{p(w_{t-1})}$$

Exercise: write out accept ratio for our BNN model

# HMC overview

Want: Transition proposal that explores "typical set" efficiently

Idea: Use gradients of target distr. in proposal

One idea:

$$T(w^{t+1}/w^t) = \text{Normal}(w^t + \epsilon \nabla_w \log P(w) \mid \sigma^2 I)$$

but this would be mode seeking, not a good sampler

We need better ideas.

## Motivation from Physics

add "momentum"  $p \sim N(0, I)$   
same size as  $w$

$$\text{Hamiltonian}(p, w) = \underbrace{-\log P(w)}_{\text{potential}} + \underbrace{-\log P(p)}_{\text{kinetic}}$$

Sample from joint  $(p, w)$   
and drop  $p$   
to get a sample of  $w$

can (approximately) evolve position & momentum while keeping the Hamiltonian (energy) conserved

$$\frac{\partial w}{\partial t} = \frac{\partial}{\partial p} [-\log P(p)]$$

$$\frac{\partial p}{\partial t} = \frac{\partial}{\partial w} [-\log P(w)]$$

gradient of our target posterior

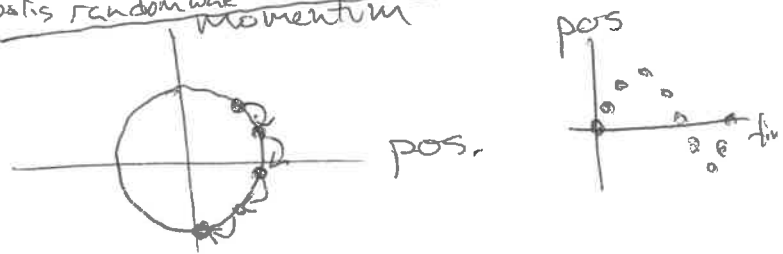
If we can run a simulation of these eqs:  $w_0, p_0, \dots, w_T, p_T$   
Cool way to explore level sets of joint distr over  $p, w$   
Accept just like Metropolis random walk

Betancourt Fig. 22 & Neal Fig 1:

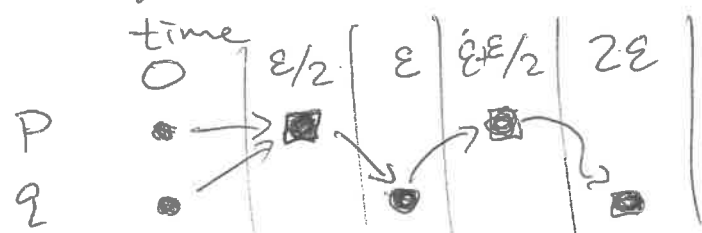
$$w(t + \frac{\epsilon}{2}) = w(t) - \frac{\epsilon}{2} \nabla_w \log P(w)$$

$$w_{t+\frac{\epsilon}{2}} \leftarrow w_t + \epsilon P(t + \frac{\epsilon}{2})$$

$$p_{t+\epsilon+\frac{\epsilon}{2}} \leftarrow p_{t+\frac{\epsilon}{2}} - \epsilon \nabla_w \log P(w)$$



Challenge: numerical approx of the differential eqns.  
(Not a goal of this course)



each step only one var changes using delta from other var

# Demos + For Next Time

04

Challenge: Write calc-potential-energy function for HMC  
for BNN model  
with partner

Demos: ((while working, open projector))  
interactive random walk vs. HMC

"Harlem Shake" with MCMC

HWZ:

Work with peers! Get help early (at office)  
use code from class on Tues (autograd + NNs)  
break down into parts

simplify! could I make a random walk sampler?  
could I do this sampler for linear regression?

Reading: black box VI + blog post  
Can I give high level punctive for each figure?  
Could I use Alg 1 pseudocode to train BNN?  
- Skim: lots of extensions (ctrl variates, etc)