

Examples of Metropolis-Hastings algorithms

① Metropolis algorithm

$$q(x'|x) = N(x', \mathbf{C})$$

$$q(x'|x) \propto \exp(-\frac{1}{2}(x-x')^T \mathbf{C}^{-1}(x-x'))$$

$$q(x|x') \propto \exp(-\frac{1}{2}(x-x')^T \mathbf{C}^{-1}(x-x')) \propto q(x'|x)$$

$$\Rightarrow \text{acceptance ratio: } \alpha = \min \left\{ 1, \frac{p(x')}{p(x)} \right\}$$

↳ Moves are likely towards regions of higher probability ($p(x') > p(x)$)

↳ If $p(x') < p(x)$, chain will likely stay put.

How to choose \mathbf{C} ?

A popular choice: set $\mathbf{C} = \mathbf{I}$

↳ Random Walk Metropolis.

This is great and very popular because it is

(i) easy to implement

(ii) problem independent (target only shows up in α)

(iii) Means also that it is slow for most problems.

An "optimal choice": set $\mathbf{C} = \sigma \mathbf{I}$

choose σ so that acceptance probability ~ 0.24
(Roberts et al 1997)

For high-dimensional problems, some paper suggests
that $\sigma \sim \frac{1}{n}$

You can see why this becomes problematic in large n problems:

Proposed Sample: $x' = x_k + \sqrt{\sigma} I$
 $= x_k + \underbrace{\frac{1}{\sqrt{n}} I}$

the proposed move is tiny if n is large

The small moves cause the samples to be correlated, which we will see will be a problem.

Other MH algorithms use different, often problem specific, proposals but this does not help with such high-dimension issues.

Situation is similar as in importance sampling: "good" proposals are useful, but it does not "break" this issue. Localization does.

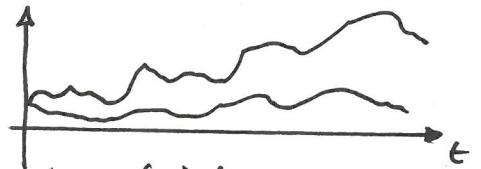
(2) Langevin MCMC

This is a stochastic differential eqn:

$$dx = f(x) dt + \sqrt{2} dW$$



Brownian motion.



two trajectories of an SDE
starting at the same initial
condition

Def of BM:

1. $W(0) = 0$

2. $W(t+T) - W(t) \sim N(0, T)$

3. $W(t)$ is almost surely continuous
(not differentiable)

4. $W(t)$ has independent increments.

$$0 \leq t_1 \leq t_2 \leq t_3 \leq t_4$$

then $W(t_1) - W(t_2)$ and $W(t_3) - W(t_2)$
are independent random variables.

Note: $x(t)$ is a random variable with distribution $p_t(x)$

Question: What is the distribution of x as $t \rightarrow \infty$

Answer: $P_\infty \propto \exp(-U(x))$ $\Rightarrow \nabla U(x) = f(x)$
? "potential"

Idea: Generate samples from a given target distribution $p(x)$ by solving the SDE

$$dx = \nabla \log p(x) dt + \sqrt{2\sigma} dw$$

Th. 3 is good because it is problem dependent

Problems: solving SDE is hard.

Solving SDE numerically is also difficult to do accurately.

An "easy" solver is a variation of the Euler scheme

$$dw \approx \Delta w \approx N(0, \Delta t)$$

$$dt \approx \Delta t$$

Their: $\textcircled{*} x_k = x_{k-1} + \nabla \log p(x_{k-1}) \Delta t + \sqrt{2\sigma} g_k$

$$g_k \sim N(0, I) \text{ iid}$$

How to use it?

"Simulate" SDE using $\textcircled{*}$ (or another numerical scheme) until you reach k^* which is large.

Repeat N times.

Samples are approximately distributed as $p(x)$.

The approximation is due to the discretization of the SDE and errors go to zero as $\Delta t \rightarrow 0$.

But a small Δt slows down the algorithm \rightarrow you need to simulate to a large $T = N \Delta t$ for each sample.

A great thing is that the algorithm is very easy to parallelize.

3) Metropolis-adjusted Langevin algorithm "MALA"

Idea: use Langevin algorithm in discrete time as a proposal,
get rid of discretization error by an accept/reject step.

$$\text{Proposal: } x' = x_k + \nabla \log p(x_k) \Delta t + \sqrt{2\Delta t} \xi, \quad \xi \sim N(0, I)$$

$$\hookrightarrow q(x'|x_k) = N(x_k + \nabla \log p(x_k) \Delta t, 2\Delta t)$$

$$q(x'|x_k) \propto \exp\left(-\frac{1}{4\Delta t} \|x' - x_k - \nabla \log p(x_k) \Delta t\|^2\right)$$

$$\alpha = \min\left\{1, \frac{p(x') q(x_k|x')}{p(x_k) q(x'(x_k))}\right\}$$

How to choose Δt , especially in high-dimensions?

Roberts & Rosenthal 1998: Optimal acceptance ratio is
0.574

$$\text{and } \Delta t \sim n^{-1/3}$$

\nearrow
much better than RWM
which was n^{-1} .

④ Hamiltonian MCMC (hybrid Mark Carlo)

Hamiltonian dynamics.

$H(q, p)$
position ↑ momentum.

Eqs of motion:

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{dp_i}{dt} = -\frac{\partial H}{\partial q_i}$$

Hamiltonian is conserved:

$$\frac{d}{dt}H = \underbrace{\sum_i \frac{\partial H}{\partial q_i} \frac{dq_i}{dt}}_{\frac{\partial H}{\partial q_i} \frac{\partial H}{\partial p_i}} + \underbrace{\sum_i \frac{\partial H}{\partial p_i} \frac{dp_i}{dt}}_{= -\frac{\partial H}{\partial p_i} \frac{\partial H}{\partial q_i}} = 0$$

Distribution of q and p ? $p(q, p) \propto \exp(-H(q, p))$

We often have: $H(q, p) = U(q) + K(p)$

↑
potential energy ↑
kinetic energy

Example: $U(q) = \frac{1}{2}q^2$ $\frac{\partial H}{\partial p} = \frac{\partial K}{\partial p} = \frac{p}{m} \quad (= v, \text{ velocity})$
 $K(p) = \frac{p^2}{2m}$ $\frac{\partial H}{\partial q} = \frac{\partial U}{\partial q} = q$

EDM: $\frac{dq}{dt} = \frac{p}{m}; \frac{dp}{dt} = -q$

In this case: $p(q, p) \propto \exp(-H(q, p)) \propto \exp(-U(q)) \exp(-K(p))$
 $\Rightarrow q, p$ are independent!

Idea of HMC Pick $U(x) = -\log p(x)$, make up a kinetic energy $K(y)$

Solve Ham.2tonian eqns of motion

$$\frac{dx}{dt} = \frac{\partial K}{\partial y}$$

$$\frac{dy}{dt} = -\frac{\partial U}{\partial x}$$

$$\exp(-\log p(x)) = p(x)$$

Marginal distribution of x : uniform ✓

Similar to Langevin, same problems: discretizing Ham.2tonian dynamics is hard because you need a solver that preserves H. The typical ones won't be good enough, you need symplectic integrators
+ you have errors due to discretization

HMC algorithm: use Leapfrog to discretize the dynamics.

$$H(x, y) = \frac{1}{2} y^T y + \cancel{F(x)} \quad \text{where } F(x) = \log p(x)$$

$$p(x, y) \propto \exp(-\frac{1}{2} y^T y) \exp(-F(x)) \Rightarrow y \sim N(0, I)$$

Leap-frog: given x_{k-1} sample $y_k \sim \underline{N(0, I)}$

$$y_{k-1} = y - \frac{h}{2} \nabla F(x)$$

$$x_k = x + h y_{k-1}$$

$$y_k = y_{k-1} - \frac{h}{2} \nabla F(x_k)$$

x_k, y_k are proposed states.

accept with prob: $\alpha = \min \left\{ 1, \frac{\exp(-H(x_k, y_k))}{\exp(-H(x, y))} \right\}$.

How to pick step-size?

acceptance prob. ≈ 0.651

$$\underline{h \approx n^{-\frac{1}{4}}} \quad \text{even better than MALA!}$$

Recap:

	Acceptance prob.	step size
RWM	0.24	n^{-1}

MALA	0.574	$n^{-\frac{1}{3}}$
------	-------	--------------------

HMC	0.651	$n^{-\frac{1}{4}}$
-----	-------	--------------------

Beskos et al. 2011

⑤ pCN MCMC.

pCN = preconditioned BGCRANK N. Richardson
pde solve

Variation of RWM:

$$x' = \sqrt{1-\beta^2} x + \beta \xi, \quad \xi \sim N(0, I)$$

$$q(x|x') \propto \exp\left(-\frac{1}{2\beta^2} (x - \sqrt{1-\beta^2} x')^T C^{-1} (x - \sqrt{1-\beta^2} x')\right)$$

$$\propto \exp\left(-\frac{1}{2\beta^2} (x^T C^{-1} x - 2x^T C^{-1} \sqrt{1-\beta^2} x' + (1-\beta^2) x'^T C^{-1} x')\right)$$

$$q(x'|x) \propto \exp\left(-\frac{1}{2\beta^2} (x'^T C^{-1} x' - 2x'^T C^{-1} \sqrt{1-\beta^2} x + (1-\beta^2) x^T C^{-1} x)\right)$$

$$\begin{aligned} \frac{q(x|x')}{q(x'|x)} &= \exp\left(-\frac{1}{2\beta^2} (x^T C^{-1} x - x'^T C^{-1} x' + (1-\beta^2)(x^T C^{-1} x - x'^T C^{-1} x))\right) \\ &= \exp\left(-\frac{1}{2\beta^2} (-\beta^2(x'^T C^{-1} x' - x^T C^{-1} x))\right) \\ &= \exp\left(-\frac{1}{2} (x^T C^{-1} x - x'^T C^{-1} x')\right) \end{aligned}$$

Assume that: $p(x) \propto \exp(-\frac{1}{2} x^T C^{-1} x - \phi(x))$

Then: $\frac{p(x')}{p(x)} \propto \exp\left(-\frac{1}{2}(x'^T C^{-1} x' - x^T C^{-1} x) - \phi(x') + \phi(x)\right)$

Acceptance ratio:

$$\alpha = \min \left\{ 1, \underbrace{\frac{p(x') q(x|x')}{p(x) q(x'|x)}} \right\}$$

$$= \exp \left(-\frac{1}{2} (x'^T C^{-1} x' - x^T C^{-1} x) - \phi(x') + \phi(x) \right)$$

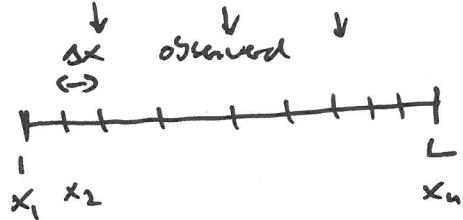
$$\cdot \exp \left(-\frac{1}{2} (x^T C^{-1} x - x'^T C^{-1} x') \right)$$

$$= \exp(\phi(x) - \phi(x'))$$

$$\hookrightarrow \alpha = \min \left\{ 1, \exp(\phi(x) - \phi(x')) \right\}$$

It's a small modification of RWM, but it can perform better.

If random variable to be estimated is as follows:



$$n = \frac{L}{\Delta x}$$

Fixed observations at δ locations ($\delta < n$)

As $\Delta x \rightarrow 0$ we approach the continuous function / infinite dimensional random variable.

Then we keep δ fixed, then pCN performs well as $\Delta x \rightarrow 0$ and $n \rightarrow \infty$.

This does not mean that pCN is a good algorithm. When L is large and δ is large pCN might not perform so well.

x is Continuous fcn
on $[0, L]$ (arbitrary D
version of R.W.)

⑥ Heat bath / partial resampling / Gibbs Sampling

Single Component Metropolis-Hastings

↳ Don't update all of x "at once", update step by step

We discuss updating one component of x at a time, but this can also be done in "blocks" (chunks of x).

Note: Single component proposes local moves, so there are probably connections w.r.t. localization.

Notation: $x_K = (x_K^1, x_K^2, \dots, x_K^n)$ state of Markov Chain
 x_K^i {^{location}
_{t time}} at time K

x_K^i is a scalar

$x_K^{(i)} = \{x_K^1, x_K^2, \dots, x_K^{i-1}, x_K^{i+1}, \dots, x_K^n\}$

everyting except its component.

At time K : iterate over blocks.

1st block: $\hat{x}^1 \sim q_1(x^1 | x_K^{(1)})$

accept with prob:

$$\alpha = \min \left\{ 1, \frac{P(\hat{x}^1 | x_K^{(1)}) q_1(x_K^1 | \hat{x}^1, x_K^{(1)})}{P(x^1 | x_K^{(1)}) q_1(x_K^1 | x_K^{(1)}, \hat{x}^1)} \right\}$$

↳ update x_K^1 to x_{K+1}^1 , if accepted, by \hat{x}^1 , otherwise set $x_{K+1}^1 = x_K^1$.

2nd block: $\hat{x}^2 \sim q_2(x^2 | x_{K+1}^1, x_K^{(2)})$
 ↳ updated state!



We need more notation.

Define: $\hat{x}_k^{(i)} = \{ \underbrace{x_{k+1}^1, x_{k+1}^2, \dots, x_{k+1}^{i-1}}_{\text{updated}}, \underbrace{x_k^{i+1}, \dots, x_k^n}_{\text{not updated}} \}$

At i th iteration:

$$\hat{x}^i \sim q_i(\cdot | \hat{x}_k^{(i)})$$

$$\alpha = \min \left\{ 1, \frac{p(\hat{x}^i | \hat{x}_k^{(i)})}{p(x^i | \hat{x}_k^{(i)})} \frac{q_i(x^i | \hat{x}_k^{(i)})}{q_i(\hat{x}^i | \hat{x}_k^{(i)})} \right\}.$$

In pseudo code:

for samples : u

for blocks/ coordinates : i

$$q_i(\hat{x}^i | \hat{x}_k^{(i)})$$

$$\alpha = \min \left\{ 1, \frac{p(\hat{x}^i | \hat{x}_k^{(i)})}{p(x^i | \hat{x}_k^{(i)})} \frac{q_i(x^i | \hat{x}_k^{(i)})}{q_i(\hat{x}^i | \hat{x}_k^{(i)})} \right\}$$

end.
end.

This can work well, but requires that we know the conditional distributions $p(x^i | \hat{x}_k^{(i)})$.

This is sometimes the case, but more often than not it is not the case.

The heat bath / partial resampling / Gibbs sampler is a special case of a Single-Component MH sampler:

$$\text{We chose: } q_i(\cdot | \mathbf{x}_u^{(i)}) = p(\cdot | \mathbf{x}_u^{(i)})$$

↳ proposal is full conditional of the target distribution.

↳ We accept every move (closed interval!)

This has even stronger assumptions:

(i) we need to know how to sample the conditional distributions.

This is easy for Gaussians, but often it is quite difficult.

Note: For the "Birch" problem of sampling a high-dimensional isotropic Gaussian this sampler is great!

This is another hint that there may be deep connections between localization and Gibbs sampling / Single Component MH.