

This algorithm manages to sample from "almost any" probability distribution (up to the burn in period) in a very simple way. Moreover, it only requires to know the target π up to a multiplicative constant (usually the normalisation constant). In addition, you will see that the possibilities that the MH offers make it possible to think of the currently employed MCMC methods as special cases of this algorithm.

Let now π be our target distribution which we assume to have a pdf w.r.t the Lebesgue measure. We only know π up to a multiplicative constant. We also assume that we are able to sample from what is called the proposal density $q(x, \cdot)$ (w.r.t. the Lebesgue measure).

1) The MH algorithm:

Start with an initial value x_0 .

Given x_n , a "candidate" y_{n+1} is generated according to $q(x_n, \cdot)$. It is then accepted with probability $\alpha(x_n, y_{n+1})$ given by:

$$\alpha(x_n, y_{n+1}) = \begin{cases} \min \left(1, \frac{\pi(y_{n+1})}{\pi(x_n)} \frac{q(x_n, y_{n+1})}{q(y_{n+1}, x_n)} \right) & \text{if } \text{denom} > 0 \\ 1 & \text{otherwise} \end{cases}$$

If y_{n+1} is accepted then $x_{n+1} = y_{n+1}$ otherwise $x_{n+1} = x_n$.

This produces a Markov chain (x_0, \dots, x_n, \dots)

this is a "reject-type" algorithm. However using all simulations and not only those accepted.

Pseudo-code: $\{x_0 \text{ given then Repeat}$

$\left. \begin{array}{l} n=0 \\ \end{array} \right\}$	$\left\{ \begin{array}{l} U_{n+1} \sim \mathcal{U}(0,1) \\ Y_{n+1} \sim q(x_n, \cdot) \\ \text{if } U_{n+1} < \alpha(x_n, Y_{n+1}) \\ \quad \text{then } x_{n+1} = Y_{n+1} \\ \quad \text{otherwise } x_{n+1} = x_n \\ n = n+1 \end{array} \right\}$	<p>← this enables to manage the acceptance or reject.</p>
---	--	---

2) Transition kernel:

Intuitively: let P be the transition kernel. 2 cases appear:

→ if we accept y which has been proposed: $p(x, y) = q(x, y) \alpha(x, y)$ with $y \neq x$. Indeed, the "Probability" to go from x to y is equalled to the probability of simulating y from x : $q(x, y)$ times its proba to be accepted $\alpha(x, y)$.

→ If we have rejected, it means that we stay at x and this probability is $P(x, \{x\}) = \int q(x, y) (1 - \alpha(x, y)) dy$ i.e. we sum over all y which have been proposed and rejected.

this writes: $\forall A \in \mathcal{X}$

$$P(x, A) = \int_A q(x, y) \alpha(x, y) dy + \delta_x(A) \int_{\mathcal{X}} (1 - \alpha(x, y)) q(x, y) dy$$

Rq: this kernel depends on π through α which only depends on $\frac{\pi(y)}{\pi(x)}$. This explains why we only need π up to a multiplicative constant.

Formally: Now we can get this expression by a calculation:

Let f be a Borel function continuous on X .

$$\begin{aligned} \mathbb{E}[f(X_{n+1}) | \mathcal{F}_n] &= \mathbb{E}_{(x, y)} [f(y) \mathbf{1}_{\{0 \leq \alpha(x, y)\}} + f(x) \mathbf{1}_{\{0 > \alpha(x, y)\}}] \\ &\stackrel{Y \perp \mathcal{F}_n}{=} \mathbb{E}_y [f(y) \alpha(x, y) + f(x) (1 - \alpha(x, y))] \\ &= \int f(y) \alpha(x, y) q(x, y) dy + \delta_x(f) \int (1 - \alpha(x, y)) q(x, y) dy \\ &= \int f(y) P(x, y) dy \text{ where } P \text{ is given above.} \end{aligned}$$

3) Invariant measure:

Prop: $\pi P = \pi$

Proof: We will use the detailed balanced relation. We need to prove that:

" $\pi(x) P(x, y) = \pi(y) P(y, x)$ " or in a more formal way:

$$\forall A, B \in \mathcal{X}, \text{ show that } \int_{A \times B} \pi(dx) P(x, dy) = \int_{A \times B} \pi(dy) P(y, dx)$$

we first write $P(x, A)$ as: $P(x, A) = P_1(x, A) + P_2(x, A)$ where

$$P_2(x, A) = c(x) \delta_x(A) \quad \text{and} \quad P_1(x, A) = \alpha(x, y) q(x, y) \mathbb{1}_{\{y \in A\}}$$

We will show the previous inequality for both P_j ($j=1, 2$) which by linearity will conclude the proof. For clarity we will work with densities w.r. to Lebesgue.

• let us start with P_1 :

By definition of the acceptance ratio α , we have

$$\begin{aligned} \pi(x) \alpha(x, y) q(x, y) &= \pi(x) \left[\frac{\pi(y) q(y, x)}{q(x, y) \pi(x)} \wedge 1 \right] q(x, y) \\ &= [\pi(y) q(y, x) \wedge \pi(x) q(x, y)] \\ &= \pi(y) \alpha(y, x) q(y, x) \end{aligned}$$

$$\begin{aligned} \text{therefore: } \int_{A \times B} \pi(dx) P_1(x, dy) &= \int_{A \times B} \pi(x) \alpha(x, y) q(x, y) dx dy \\ &= \int_{A \times B} \pi(y) \alpha(y, x) q(y, x) dx dy \\ &= \int_{A \times B} \pi(dy) P_1(y, dx) \end{aligned}$$

• For P_2 we will use the convention that $\int_{\emptyset} = 0$

$$\begin{aligned} \int_{A \times B} \pi(dx) P_2(x, dy) &= \int_{A \times B} \pi(dx) c(x) \delta_x(y) dy \\ &= \int_{A \cap B} \pi(dx) c(x) = \int_{A \cap B} \pi(dy) c(y) \\ &= \int_{A \times B} \pi(dy) c(y) \delta_y(dx) = \int_{A \times B} \pi(dy) P_2(y, dx) \end{aligned}$$

4) Examples: what choice for q ?

a) Independent sampler:

This is the most simple case since the proposal does not depend on the current state of the chain: $q(x, y) = q(y)$. Of course, the dependency will appear

through the acceptance ratio: $\alpha(x, y) = \left(\frac{\pi(y) q(x)}{q(y) \pi(x)} \wedge 1 \right)$

(This looks like a generalisation of the acceptance-reject method.)

The convergence properties of the generated chain are obviously depending on q

Prop: $(X_n)_n$ is irreducible and aperiodic iff q is positive (>0) (4)
on the support of π

this also enable to get the stronger convergence property + the geometric ergodicity under a weak condition. Here I give you the theorem which you will be able to prove with the tools introduced in the following classes.

let us first define the Total variation norm:

$$\|P(x, \cdot) - \pi\|_{TV} = \sup_{A \in \mathcal{X}} \left| \int_A (P(x, y) - \pi(y)) dy \right| \quad \text{where } P \text{ is used for both the kernel and its corresponding probability measure.}$$

then:

thm: the mcmc algorithm given by the Independent sampler (IS) M-H algorithm produces a uniformly ergodic Markov chain if there exist a constant $M > 0$ s.t. : $\pi(x) \leq M q(x) \quad \forall x \in \text{supp}(\pi)$ (11)

$$\text{In this case : } \|P^n(x, \cdot) - \pi(\cdot)\|_{TV} \leq \left(1 - \frac{1}{M}\right)^n$$

Moreover, if $\forall M > 0 \exists A \in \mathcal{X}$ s.t. $\pi(A) \neq 0$ and (11) is not satisfied on A , then (X_n) is not geometrically ergodic.

This provides a simple condition to prove the convergence property.

One can also wonder: how efficient is this sampler? Its efficiency may be measured by its mean acceptance ratio: $\mathbb{E}[\alpha(X, Y)]$

lemma: If the chain is stationary, then if q satisfies (11) then

$$\mathbb{E}[\alpha(X, Y)] \geq \frac{1}{M}$$

$$\text{Proof: } \mathbb{E}[\alpha(X, Y)] = \int \left(1, \frac{\pi(y) q(x)}{q(y) \pi(x)}\right) \pi(x) q(y) dx dy$$

$$= \int (\pi(x) q(y) \wedge \pi(y) q(x)) dx dy$$

$$\stackrel{q(y) \geq \frac{1}{M} \pi(y)}{=} \frac{1}{M} \int [\pi(x) \pi(y) \wedge \pi(y) \pi(x)] dx dy = \frac{1}{M}$$

Ex: Simulation of $\pi = \mathcal{D}(0,1)$ using $q = \mathcal{D}(0, \sigma^2)$

(5)

$$\alpha(x,y) = \min\left[1, \exp\left(-\frac{1}{2}y^2 + \frac{1}{2}x^2 + \frac{1}{2\sigma^2}y^2 - \frac{1}{2\sigma^2}x^2\right)\right]$$
$$= \min\left[1, \exp\left(-\frac{1}{2}(y^2 - x^2)(1 - \sigma^{-2})\right)\right]$$

→ if $\sigma^2 > 1$: we have a proposal with a heavier tail than the target distribution

We will propose many elements far from 0: i.e. $|Y_{n+1}| > |X_n|$. But

$$1 - \frac{1}{\sigma^2} > 0 \Rightarrow \exp\left(-\frac{1}{2}(Y_{n+1}^2 - X_n^2)(1 - \sigma^{-2})\right) < 1 \text{ and even more}$$

if $Y_{n+1}^2 \gg X_n^2$. These elements are therefore very likely to be rejected

The elements Y_{n+1} s.t. $Y_{n+1}^2 \leq X_n^2$ will be accepted with probability 1

Therefore, to compensate the heavy tail proposal, many elements will be rejected

→ On the other hand, if $\sigma^2 < 1$, all the candidates s.t. $|Y_{n+1}| > |X_n|$ will

be accepted with probability 1 but not all the elements in the complementary set.

This compensates the fact that many elements will be proposed close to 0.

Demo: (able: $\pi = \mathcal{D}(0, \mu) |_{(0,1)}$ avec $\mu = 0,01$)

General Independence sampler → draw a distribution with
Small variance than large one.

b) Symmetric Random Walk M-H (SRWM)

Assume that \mathcal{X} is a vector space (typically \mathbb{R}^d) on which we can define a random walk with transition density $q(x,y) = q(|x-y|)$ with q even function.

In this case, the acceptance ratio reduces to $\alpha(x,y) = \min\left(1, \frac{\pi(y)}{\pi(x)}\right)$

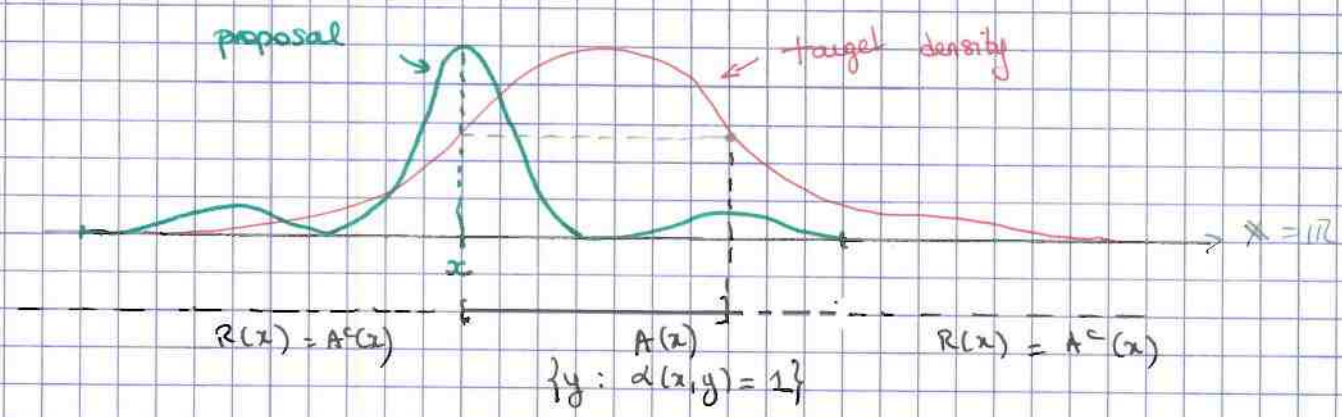
This has the following interpretation:

With the usual M-H, we accept systematically Y_{n+1} such that the "likelihood ratio" $\frac{\pi(y)}{q(x,y)}$ is greater than the opposite move ($Y_{n+1} \rightarrow X_n$)

Here, the condition becomes: π increases \Rightarrow accept with probability 1.

The algorithm biases the Random Walk by promoting samples which move toward the modes of π

Schematically this can be viewed here:



Although quite simple, this algorithm is very versatile as one has a large range of possible q . In particular, it is appealing as π only shows up in α .

However, it has some drawbacks also; in particular:

thm: If π has non-compact support and if q is symmetric, the Markov chain $(X_n)_n$ is not uniformly ergodic.

However with additional assumptions on π and q , one can achieve the geometric ergodicity.

Ex: $\pi(x) \propto \exp(-\gamma|x|^\beta)$, $\gamma > 0$, $\beta > 0$

$$q(x, y) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-y)^2}{2\sigma^2}\right)$$

Given X_n , the SRW samples $Y_{n+1} \sim \mathcal{N}(X_n, \sigma^2)$ and gives

$$\alpha(X_n, Y_{n+1}) = \min\left(1, \exp\left(\gamma(|X_n|^\beta - |Y_{n+1}|^\beta)\right)\right)$$

this implies as previously 2 cases: if $|Y_{n+1}| \leq |X_n|$ then it is automatically accepted; on the other hand if $|Y_{n+1}| > |X_n|$ it may be rejected

It seems that σ^2 does not have a real impact on the results. To see how it actually does, we compute the auto-correlation of the generated chain.

$$C_i = \frac{\sum_{j=1}^{N-i} (X_j - \bar{X})(X_{j+i} - \bar{X})}{N-i} \times \frac{N}{\sum_{j=1}^N (X_j - \bar{X})^2} \quad \text{with } \bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$$

demo: π Exponential density on $(0, 1)$

$\gamma = 1$, $\beta = 1$

q : Gaussian RW

$\sigma^2 = 2 \rightarrow$ good C
 $\sigma^2 = 0,1 \rightarrow$ bad C.

A second question is the dependency on X_0 .

slides

$$\pi = \exp$$

$$q = \mathcal{P}(0, \sigma^2)$$

$$X_0 = 1 \text{ or } X_0 = 100$$

↳ burn in 200 steps

The auto-correlation measure how "independent" are the first samples w.r.t. the last ones. You expect that it will decrease quickly so that your chain does not depend too long on its initial value.

this is what is illustrated here.

You can also introduce a burn-in and the autocorrelation will help you have a first guess of its value - this plot reflects the stationarity behaviour.

Ex2: A heavy tailed Cauchy distribution:

$$\pi(x) \propto (1 + |x|)^{-2}$$

$$q \sim \mathcal{P}(0, \sigma^2) \text{ where } \sigma^2 \text{ so that } \mathbb{E}[\alpha] \approx 35\%$$

slide 5 : quite good at first sight but! the maximum value of our sampler is less than 30 even though we have done 10.000 iterations. The Cauchy distribution has a heavy tail and one can expect a maximum value around 200.

Since the cumulative distribution function is available here, one can plot a q-q plot to visually assess the goodness of fit of this output with our target (q-q plot: plots the quantiles against each other → goal: to be as close as possible to the diagonal)

slide 6: 3 independent runs of the algo together with their qq plot.

→ good agreement for small values

→ fail completely for large ones

→ ⊕ erratic behaviour 3 runs → 3 very different shapes.

demo: can multimodal

c.) MALA:

One may wonder how one can better use the target into the proposal. Many solutions were proposed although one appeared to be well adapted to high dimensional problems. It starts from a diffusion equation in continuous time:

$$dX(t) = a(t, X(t)) dt + \sigma^2(t, X(t)) dB(t) \quad \text{where } B \text{ is a Brownian motion.}$$

This equation is given as elements of $L^2(dP)$

In terms of diffusion process π is the stationary distribution under conditions on a and σ^2 .

The goal of the MALA algorithm is to go from continuous to discrete schemes. This enables (thanks to an Euler scheme) to propose the following algorithm:

for a C^2 target π ; we define $q(x, y)$ by:

$$q(x, y) = \mathcal{N}(x + \delta \nabla_x \log \pi; \delta \text{Id})$$

This introduces a complex dependency of q w.r.t. π

→ Expliquer sur ex gaussienne

$$\pi \propto e^{-\frac{x^2}{2}}$$
$$\log \pi = -\frac{x^2}{2}$$
$$\nabla \log \pi = -x$$

Rappel vers les valeurs du mode

→ favorable au multimodal. → demo