

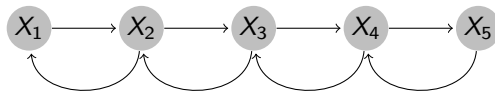
Lecture 8: Markov chain Monte Carlo

- **Goal:** Generation of samples or approximation of an expected value for a (possibly high-dimensional) density $\pi(x)$.
- Application of ordinary Monte Carlo methods is difficult.
- However, **Markov chain Monte Carlo (MCMC) methods** will then be a useful alternative.



Andrey Markov (1856 – 1922),
Russian mathematician.

Markov chain:



en.wikipedia.org/wiki/Markov_chain

Given the previous observation X_{i-1} , X_i is independent of the sequence of events that preceded it.

Review: Markov chains

A Markov chain is a stochastic process $\{X_i\}_{i=0}^{\infty}$, $X_i \in S$, where given the present state, past and future states are independent (**Markov assumption**):

$$P(X_{i+1} = x_{i+1} \mid X_0 = x_0, X_1 = x_1, \dots, X_i = x_i) = P(X_{i+1} = x_{i+1} \mid X_i = x_i)$$

Idea of Markov chain Monte Carlo

Idea

Simulate a **Markov chain** X_1, \dots, X_i, \dots , which is designed in a way such that $P(X_i = x)$ **converges to the target distribution $\pi(x)$** , e.g. the **posterior distribution**.

Properties:

- After convergence, one obtains random samples from the target distribution, which can be used to estimate posterior characteristics.
- Samples will typically be **dependent**.

Central algorithms:

- **Metropolis-Hasting algorithm**
- **Gibbs sampling**

Review: Markov chains

A Markov chain with stationary transition probabilities can be specified by:

- **the initial distribution** $P(X_0 = x_0) = g(x_0)$
- **the transition matrix**

$$P(x^* \mid x) = P(X_{i+1} = x^* \mid X_i = x) \quad [= P_{xx^*}]$$

Review: Markov chains

A Markov chain has a **unique limiting distribution** $\pi(x)$ if the chain is **irreducible**, **aperiodic**, and **positive recurrent**. If so, the limiting distribution $\pi(x) = \lim_{i \rightarrow \infty} P(X_i = x)$ is given by

$$\begin{aligned} \pi(x^*) &= \sum_{x \in S} \pi(x) P(x^* | x) \quad \text{for all } x^* \in S \\ \sum_{x \in S} \pi(x) &= 1 \end{aligned} \tag{1}$$

A sufficient condition for (1) is the **detailed balance condition**:

$$\pi(x) P(x^* | x) = \pi(x^*) P(x | x^*) \quad \text{for all } x, x^* \in S \tag{2}$$

which gives a **time-reversible Markov chain**.

Idea

Focus on (2) instead. **We want to find $P(x^* | x)$ that solves**

$$\pi(x) P(x^* | x) = \pi(x^*) P(x | x^*) \quad \text{for all } x, x^* \in S$$

Here, we still have many solutions. However, we do not need a general solution, one (good) solution is enough.

Problem statement

In stochastic processes course: The Markov chain is given, i.e. $P(x^* | x)$ is given, find $\pi(x)$.

Now: $\pi(x)$, $x \in S$ is given, **want to find $P(x^* | x)$** , $x, x^* \in S$ so that

$$\begin{aligned} \pi(x^*) &= \sum_{x \in S} \pi(x) P(x^* | x) \quad \text{for all } x^* \in S \\ \sum_{x \in S} \pi(x) &= 1 \end{aligned}$$

However, # unknowns: $|S| \cdot (|S| - 1)$; # equations: $|S|$.

\Rightarrow **many solutions exist – we want one!**

(Note: $|S|$ can be huge, so solving this as a matrix equation is not possible.)

Idea II

Try:

$$P(x^* | x) = Q(x^* | x) \alpha(x^* | x), \quad x^* \neq x$$

$$P(x | x) = 1 - \sum_{x^* \neq x} Q(x^* | x) \alpha(x^* | x)$$

so that $\sum_{x^* \in S} P(x^* | x) = 1$.

Here:

- $Q(x^* | x)$ is an almost arbitrary transition matrix (**proposal kernel**) for some other irreducible Markov chain.
- $\alpha(x^* | x) \in [0, 1]$ is an **acceptance probability**.

Metropolis-Hastings algorithm

```
1: Init  $x_0 \sim g(x_0)$ 
2: for  $i = 1, 2, \dots$  do
3:   Generate a proposal  $x^* \sim Q(x^* | x_{i-1})$ 
4:   Compute the acceptance probability  $\alpha(x^* | x_{i-1})$ 
5:    $u \sim U(0, 1)$ 
6:   if  $u < \alpha(x^* | x_{i-1})$  then
7:      $x_i \leftarrow x^*$ 
8:   else
9:      $x_i \leftarrow x_{i-1}$ 
10:  end if
11: end for
```

See blackboard

Metropolis-Hastings algorithm

```
1: Init  $x_0 \sim g(x_0)$ 
2: for  $i = 1, 2, \dots$  do
3:   Generate a proposal  $x^* \sim Q(x^* | x_{i-1})$ 
4:    $u \sim U(0, 1)$ 
5:   if  $u < \underbrace{\min \left( 1, \frac{\pi(x^*)}{\pi(x_{i-1})} \times \underbrace{\frac{Q(x_{i-1} | x^*)}{Q(x^* | x_{i-1})}}_{\text{Proposal ratio}} \right)}_{\text{Acceptance probability } \alpha}$  then
6:      $x_i \leftarrow x^*$ 
7:   else
8:      $x_i \leftarrow x_{i-1}$ 
9:   end if
10: end for
```

What is Q and α ?

Acceptance step

- In the acceptance step the proposal x^* is accepted with probability α as new value of the Markov chain.
- This is similar to rejection sampling. However, here no constant c needs to be determined.
- Further, if we reject, then we retain the sample.

History of Metropolis-Hastings

- The algorithm was presented 1953 by Metropolis, Rosenbluth, Rosenbluth, Teller and Teller from the Los Alamos group. It is named after the first author **Nicholas Metropolis**.
- **W. Keith Hastings** extended it to the more general case in 1970.
- It was then ignored for a long time.
- Since 1990 it has been used more intensively.

Toy example

We considered the Poisson distribution

$$\pi(x) = \frac{10^x}{x!} e^{-10}, \quad x = 0, 1, 2, \dots$$

Choose proposal kernel

- If $x = 0$

$$Q(x^*|0) = \begin{cases} \frac{1}{2} & \text{for } x^* \in \{0, 1\} \\ 0 & \text{otherwise} \end{cases}$$

- For $x > 0$

$$Q(x^*|x) = \begin{cases} \frac{1}{2} & \text{for } x^* \in \{x-1, x+1\} \\ 0 & \text{otherwise} \end{cases}$$

Toy example

Setting see blackboard and computer

Toy example

- If $x = 0$

$$\alpha(0|0) = \min \{1, 1\} = 1$$

$$\alpha(1|0) = \min \{1, 10\} = 1$$

- If $x > 0$

$$\alpha(x-1|x) = \min \left\{ 1, \frac{\frac{10^{x-1}}{(x-1)!} e^{-10}}{\frac{10^x}{x!} e^{-10}} \cdot \frac{\frac{1}{2}}{\frac{1}{2}} \right\} = \min \left\{ 1, \frac{x}{10} \right\} \quad (3)$$

$$\alpha(x+1|x) = \min \left\{ 1, \frac{\frac{10^{x+1}}{(x+1)!} e^{-10}}{\frac{10^x}{x!} e^{-10}} \cdot \frac{\frac{1}{2}}{\frac{1}{2}} \right\} = \min \left\{ 1, \frac{10}{x+1} \right\} \quad (4)$$

From (3) we see that $\alpha = 1$ if $x > 9$ and $x/10$ else.

From (4) we see that $\alpha = 1$ if $x \leq 9$ and $10/(x+1)$ else.

Toy example

Note this gives for $x > 0$:

$$P(x-1|x) = \frac{1}{2} \min \left\{ 1, \frac{x}{10} \right\} = \begin{cases} \frac{x}{20} & \text{for } x \leq 9 \\ \frac{1}{2} & \text{for } x > 9 \end{cases}$$

$$P(x+1|x) = \frac{1}{2} \min \left\{ 1, \frac{10}{x+1} \right\} = \begin{cases} \frac{1}{2} & \text{for } x \leq 9 \\ \frac{5}{x+1} & \text{for } x > 9 \end{cases}$$

$P(x|x)$ follows directly.

(For $x = 0$ we have $P(0|0) = 1/2$ and $P(1|0) = 1/2$).

However, we do not have to compute these values! (Show R-code `demo_toyMCMC2.R`)

What about

- **Irreducible**: Must be checked in each case. Must choose $Q(x^* | x)$ so that this is ok.
- **Aperiodic**: Sufficient that $P(x | x) > 0$ for one $x \in S$, so sufficient that $\alpha(x^* | x) < 1$ for one pair $x^*, x \in S$.
- **Positive recurrent**: for finite S , irreducibility is sufficient. More difficult in general, but if Markov chain is not recurrent we will see this as drift in the simulations. (In practice usually no problem).

Remarks on the Metropolis-Hastings algorithm

- Under some regularity conditions it can be shown that the **Metropolis-Hasting algorithm converges to the target distribution** regardless of the specific choice of $Q(x|x_{i-1})$.
- However, the **speed of convergence** and the **dependence between the successive samples** depends strongly on the proposal distribution.
- Since we only need to compute the ratio $\pi(x^*)/\pi(x_{i-1})$, the **proportionality constant is irrelevant**.
- Similarly, we only care about $Q(\cdot)$ up to a constant.
- Often it is advantageous to calculate the acceptance probability on **log-scale**, which makes the computations more stable.

Special cases of the Metropolis-Hastings algorithm

Depending on the choice of $Q(x^*|x)$ different special cases result. In particular, two classes are important

- **The independence proposal**
- **The Metropolis algorithm**

Independence proposal

- The proposal distribution does not depend on the current value x_{i-1}

$$Q(x|x_{i-1}) = Q(x).$$

- $Q(x)$ is an approximation to $\pi(x)$.
- The sampler is closer to rejection sampler. However, here if we reject, then we retain the sample.

Experience:

- Performance is either very good or very bad, usually very bad.
- The tails of the proposal distribution should be at least as heavy as the tails of the target distribution.

The Metropolis algorithm

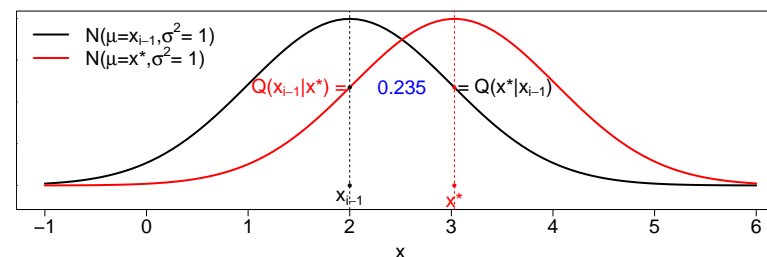
The proposal density is symmetric around the current value, that means

$$Q(x_{i-1}|x^*) = Q(x^*|x_{i-1}).$$

Hence,

$$\alpha = \min \left(1, \frac{\pi(x^*)}{\pi(x_{i-1})} \times \frac{Q(x_{i-1}|x^*)}{Q(x^*|x_{i-1})} \right) = \min \left(1, \frac{\pi(x^*)}{\pi(x_{i-1})} \right)$$

A particular case is the **random walk proposal**, defined as the current value x_{i-1} plus a random variate of a 0-centred symmetric distribution.



Examples for random walks proposal

Assume x is scalar.

Then all proposal kernels, which **add a random variable generated from a zero-symmetrical distribution to the current value** x_{i-1} , are random walk proposals. For example:

$$x^* \sim \mathcal{N}(x_{i-1}, \sigma^2)$$

$$x^* \sim t_\nu(x_{i-1}, \sigma^2)$$

$$x^* \sim \mathcal{U}(x_{i-1} - d, x_{i-1} + d)$$

Efficiency of the Metropolis-Hastings algorithm

The efficiency and performance of the Metropolis-Hastings algorithm depends crucially on the **relative frequency of acceptance**.

An acceptance rate of one is not always good. Consider the random walk proposal:

- Too large acceptance rate \Rightarrow Slow exploration of the target density.
- Too small acceptance rate \Rightarrow Large moves are proposed, but rarely accepted.

Tuning the acceptance rate:

- For **random walk proposals**, acceptance rates between **20% and 50%** are typically recommended. They can be achieved by changing the variance of the proposal distribution.
- For **independence proposals** a **high acceptance rate** is desired, which means that the proposal density is close to the target density.

Example: Random walk proposal

Exploration of a standard Gaussian distribution ($\mathcal{N}(0, 1)$) using a random walk Metropolis algorithm. As proposal assume a Gaussian distribution with variance σ^2 , where.

- $\sigma^2 = 0.24$
- $\sigma^2 = 2.4$
- $\sigma^2 = 24$

See R-code `demo_mcmcRW.R`.

Rao: Independence proposal

$$\theta^* \sim \mathcal{N}(\text{Mod}(\theta|\mathbf{y}), F^2 \times I_p^{-1}), \quad (5)$$

where $\text{Mod}(\theta|\text{data})$ denotes the posterior mode, I_p the negative curvature of the log posterior at the mode, and F a factor to blow up the standard deviation.

Of note, asymptotically the posterior distribution follows (5) for $F = 1$.

Example of Rao (1973)

The vector $\mathbf{y} = (y_1, y_2, y_3, y_4) = (125, 18, 20, 34)$ is multinomial distributed with probabilities

$$\left\{ \frac{1}{2} + \frac{\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4}, \frac{\theta}{4} \right\}$$

We would like to simulate from the posterior distribution (assuming a uniform prior)

$$f(\theta|\mathbf{y}) \propto (2 + \theta)^{y_1} (1 - \theta)^{y_2 + y_3} \theta^{y_4}.$$

using MCMC and compare two proposal kernels:

1. independence proposal
2. random walk proposal

See R-code `demo_mcmcRao.R`.

Rao: Random walk proposal

$$\theta^* \sim \text{U}(\theta^{(k)} - d, \theta^{(k)} + d),$$

where $\theta^{(k)}$ denotes the current state of the Markov chain and $d = \sqrt{12}/2 \cdot 0.1$.