

Lecture 3: Markov Chain Monte Carlo

Tianpei Xie

Sep. 28th., 2022

Contents

1	Basic Concept of Markov Chain	2
2	Markov Chain Monte Carlo	5
2.1	From Vanilla Monte Carlo to MCMC	5
2.2	Metropolis-Hastings Algorithm	6
3	Some Special Metropolis-Hastings Algorithms	10
3.1	Random-walk Metropolis-Hastings	10
3.2	Independent Metropolis-Hastings (IMH) Algorithm	11
3.3	Configurational bias Monte Carlo	12

1 Basic Concept of Markov Chain

- **Markov Chain** $(X_t)_t$ is a **probabilistic graphical model** over a chain graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}_C)$, where each random variable X_t only has exactly one children X_{t+1} and one parent X_{t-1} . Denote the index of variable t as the time. Markov chain $(X_t)_t$ is also a **stochastic process**.
- By Markov property,

$$P(X_{t+1}|X_t, X_{t-1}, \dots, X_1) = P(X_{t+1}|X_t).$$

It is seen that the transition probability does not depend on the time t , i.e. Markov chain is **time-invariant**.

- Define the **kernel** of Markov Chain as the **time-invariant transition probability**

$$K(x, y) = p(x, y) := P(X_{t+1} = y | X_t = x). \quad (1)$$

Then the **m -step transition probability** is defined as

$$K^m(x, y) = P(X_{t+m} = y | X_t = x). \quad (2)$$

- For $X_t \in \mathcal{X} := \{1, \dots, n\}$ as discrete random variable with $|\mathcal{X}| = n$, we can define the **transition matrix**

$$\mathbf{K} = [K(i, j)]_{n \times n} \quad (3)$$

- We have the **Chapman-Kolmogorov equation** [Ross, 2014]:

$$\begin{aligned} K^{m+n}(x, y) &= \sum_z K^m(x, z) K^n(z, y), \quad \forall x, y \in \mathcal{X} \\ \Rightarrow \mathbf{K}^{m+n} &= \mathbf{K}^m \mathbf{K}^n \end{aligned} \quad (4)$$

That is, we split the $(m+n)$ -step path from $x \rightarrow y$ into all possible combination of a m -step path from $x \rightarrow z$ and a n -step path from $z \rightarrow y$ for some intermediate state z .

- Define $T_j = \min \{t \geq 1 : X_t = j\}$ as the time steps for Markov Chain $(X_t)_t$ to **hit** state j for the **first time**. T_j is called the state j 's **first hitting time**.

Denote $f_{i,j}$ be the **probability of ever hitting state j (within finite time) starting from state i** . That is

$$f_{i,j} := P(T_j < \infty | X_0 = i) \quad (5)$$

Denote $f_{i,j}^{(m)}$ be the **probability of hitting at state j at time m starting from state i**

$$f_{i,j}^{(m)} := P(T_j = m | X_0 = i) \quad (6)$$

- Define $N(y) := \sum_{t=0}^{\infty} \mathbb{1} \{X_t = y\}$ is the **total number of times hitting the state y** .

$$P(N(y) \geq 1 | X_0 = x) = P(T_y < \infty | X_0 = x) = f_{x,y} \quad (7)$$

$$P(N(y) \geq m | X_0 = x) = P(N(y) \geq 1 | X_0 = x) P^{m-1}(N(y) \geq 1 | X_0 = y) \quad (8)$$

$$= f_{x,y} f_{y,y}^{m-1}$$

Note that in order to visit y at least m times, we need to visit y first time and stating from y recurrently visit y ($m - 1$) times.

The random variable $N(x)|X_0 = x$ follows a **geometric distribution** with mean $1/(1 - f_{x,x})$.

$$P(N(x) = m | X_0 = x) = (1 - f_{x,x}) f_{x,x}^{m-1} \quad (9)$$

- For any pair $x, y \in \mathcal{X}$, if there exists $n \in \mathbb{N}_+$ so that $K^n(x, y) > 0$, then the state y is **accessible** from state x . This is equivalent to say that the probability of hitting time of y being finite starting from x is above zero, i.e. $f_{x,y} > 0$.
- If x is accessible from y , and y is accessible from x , then we say that x and y **coummunicate**, $x \leftrightarrow y$. It is easy to check that this is an **equivalence relation**.
- A Markov Chain is **irreducible** if it has **only one equivalence class**, i.e. all states in \mathcal{X} communicate to each other.
- Based on hitting time, we can categorize states into two groups:
 - A state i is **recurrent** if and only if $f_{i,i} = P(T_i < \infty | X_0 = i) = 1$, i.e. the Markov Chain will definitely revisit the state i after stating from i .
 - A recurrent state i is **positive recurrent** if the *expected returning time* $\mathbb{E}[T_i | X_0 = i] < \infty$; otherwise we say it is **null recurrent**.
 - A state i is called **transient** if $f_{i,i} < 1$.

- **Proposition 1.1** *If i is recurrent, and $i \rightarrow j$, then j is also recurrent. Therefore, in any equivalent class, either all states are recurrent or all are transient. In particular, if the chain is irreducible, then either all states are recurrent or all are transient.*

Proposition 1.2 *An irreducible finite state Markov chain must be positive recurrent.*

- **Definition** The probability of states $\{\pi(x), \forall x \in \mathcal{X}\}$ is a **stationary distribution** if and only if

$$\pi(y) = \sum_{x \in \mathcal{X}} K(x, y) \pi(x), \forall y \in \mathcal{X} \quad (10)$$

$$\pi^T = \pi^T \mathbf{K} \quad (11)$$

That is, π is the eigenvector of stochastic matrix \mathbf{K} corresponding to eigenvalue $\lambda_0 = 1$.

The **stationary distribution does not change over time**.

- **Proposition 1.3** *Suppose that the **limit distribution** $\lim_{t \rightarrow \infty} P(X_t = y)$ exists, and*

$$\lim_{t \rightarrow \infty} K^t(x, y) = \pi(y), \quad \forall x, y \in \mathcal{X}$$

*which is independent of where it starts from, then the Markov Chain has a **unique stationary distribution** and*

$$\lim_{t \rightarrow \infty} P(X_t = y) = \pi(y), \quad \forall y \in \mathcal{X} \quad (12)$$

i.e. the limit distribution is stationary distribution.

- **Proposition 1.4** (*Global Balance Equation*)

The stationary distribution $\{\pi(x), \forall x \in \mathcal{X}\}$ satisfies the following **global balance equation**:

$$\sum_{j \in \mathcal{X}} \pi(i)K(i, j) = \sum_{j \in \mathcal{X}} \pi(j)K(j, i). \quad (13)$$

This means the total flow out of i (LHS) is equal to the total flow into i (RHS) in steady state.

- **Proposition 1.5** (*Detailed Balance Equation*)

For distribution $\{\pi(x), \forall x \in \mathcal{X}\}$, if the following **detailed balance equation** is satisfied

$$\pi(i)K(i, j) = \pi(j)K(j, i), \quad \forall i, j \in \mathcal{X} \quad (14)$$

then $\{\pi(x), \forall x \in \mathcal{X}\}$ is a stationary distribution.

- **Theorem 1.6** (*Stationary distribution for transient and null recurrent states*)

Let $\{\pi(x), \forall x \in \mathcal{X}\}$ be stationary distribution. If $x \in \mathcal{X}$ is **transient** or **null recurrent** state, then

$$\pi(x) = 0.$$

- **Theorem 1.7** [Ross, 2014]

An **irreducible recurrent** Markov Chain has a **unique stationary distribution** $\{\pi(x)\}$, given

$$\pi(x) = \frac{1}{\mu_x}, \quad \forall x \in \mathcal{X} \quad (15)$$

where $\mu_x := \mathbb{E}[T_x | X_0 = x]$ is the **expected first return time** of state x .

It implies that as $n \rightarrow \infty$, for any state $x \in \mathcal{X}$, the fraction of time that Markov Chain stays at x is unchanged and is the reciprocal of the expected first return time.

- **Definition** The **periodicity** of a state $x \in \mathcal{X}$ is defined as

$$d(x) = \text{g.c.d.} \{t \geq 0 : K^t(x, x) > 0\} \quad (16)$$

where g.c.d. is the **greatest common divisor**.

- **Definition** If $d(x) \geq 2$, then state x is **periodic**. If $d(x) = 1$, then state x is **aperiodic**

The periodicity property is *closed* under the equivalence class C .

- **Definition** A Markov Chain is **irreducible, positive recurrent** and **aperiodic**, then it is called **ergodic**.

- **Theorem 1.8** A Markov Chain is **ergodic** having stationary distribution π , then

$$\lim_{t \rightarrow \infty} K^t(x, y) = \pi(y), \quad \forall x, y \in \mathcal{X} \quad (17)$$

That is, when a Markov Chain is ergodic, its marginal state distribution will converge to the stationary distribution.

- **Definition** A Markov Chain $(X_t)_t$ is called **time-reversible**, if it has stationary distribution π and the detailed balance equation is satisfied:

$$\pi(i)K(i, j) = \pi(j)K(j, i), \quad \forall i, j \in \mathcal{X} \quad (18)$$

- The reversed process $(Y_k)_k := (X_{t-k})_k$ is a Markov Chain and its transition probability

$$Q(i, j) = \frac{\pi(j)K(j, i)}{\pi(i)} \quad (19)$$

Note that $(Y_k)_k$ and $(X_t)_t$ are statistically equivalent since $Q(i, j) = K(i, j)$.

- **Theorem 1.9** *An ergodic Markov Chain $(X_t)_t$ for which $K(i, j) = 0$ whenever $K(j, i) = 0$ is **time-reversible** if and only if starting from any state i , any path back to i has the **same probability** as its reverse path. That is, for path $i \rightarrow i_1 \rightarrow i_2 \rightarrow \dots \rightarrow i_k \rightarrow i$ and its reverse path $i \leftarrow i_1 \leftarrow i_2 \leftarrow \dots \leftarrow i_k \leftarrow i$*

$$K(i, i_1) K(i_1, i_2) \dots K(i_k, i) = K(i, i_k) \dots K(i_2, i_1) K(i_1, i), \quad \forall i, i_1, \dots, i_k \in \mathcal{X} \quad (20)$$

- **Theorem 1.10 (Reversal Test)**

Let \mathbf{K} be a stochastic matrix indexed by a countable set \mathcal{X} and let π be a probability distribution on \mathcal{X} . Let \mathbf{Q} be a stochastic matrix indexed by \mathcal{X} such that

$$\pi(i)Q(i, j) = \pi(j)K(j, i), \quad \forall i, j \in \mathcal{X}. \quad (21)$$

Then π is a stationary distribution of \mathbf{K}

- **Theorem 1.11 (Ergodic Theorem)** [Robert and Casella, 1999]

If $(X_t)_t$ is Harris recurrent with a σ -finite invariant measure π , then for any $f, g \in L_1(\pi)$ with $\mathbb{E}_\pi[g] \neq 0$,

$$\lim_{T \rightarrow \infty} \frac{S_T(f)}{S_T(g)} = \frac{\mathbb{E}_\pi[f]}{\mathbb{E}_\pi[g]} = \frac{\int f(x)d\pi(x)}{\int g(x)d\pi(x)} \quad (22)$$

It can be shown that if $(X_t)_t$ is **Harris positive** with **stationary distribution** π and if $S_T(h)$ converges μ_0 -almost surely (μ_0 a.s.) to $\mathbb{E}_\pi[h]$, for an initial distribution μ_0 , this convergence occurs for **every initial distribution** μ .

- **Theorem 1.12 (Central Limit Theorem for reversible chains)**[Robert and Casella, 1999]

If $(X_t)_t$ is aperiodic, irreducible, and reversible with stationary distribution π , the **Central Limit Theorem** applies when

$$0 < \sigma^2 = \mathbb{E}_\pi[g^2(X_t)] + 2 \sum_{s=1}^{\infty} \mathbb{E}_\pi[g(X_t)g(X_{t+s})] < \infty. \quad (23)$$

2 Markov Chain Monte Carlo

2.1 From Vanilla Monte Carlo to MCMC

- **Definition** A Markov Chain Monte Carlo (MCMC) method for the simulation of a distribution f is any method producing an ergodic Markov chain $(X_t)_t$ whose stationary distribution is f .

- Compared to vanilla Monte Carlo (e.g. inverse transform, reject sampling, importance sampling), MCMC has the following **characteristics**:

- Unlike vanilla Monte Carlo methods, which rely on i.i.d samples, **Markov Chain Monte Carlo (MCMC) methods** generate **dependent samples** via Markov chain.
- The MCMC updates **preserve the probability measure** π when convergence is attained. That is, *when the Markov chain converges*, the distribution of X_t is the same as the distribution of X_{t+1}, X_{t+2}, \dots . Thus we have obtained a sequence of **identically distributed (but dependent) samples**. When Markov chain converges (**mixing**), we can use samples the same way as we did in vanilla Monte Carlo to approximate the expectation. In particular, an MCMC estimator is

$$J_T = \widehat{\mathbb{E}}_{\pi} [h(X)] = \frac{1}{T} \sum_{t=0}^T h(X_t). \quad (24)$$

The ergodic theorem guarantees the (almost sure) convergence of the empirical average to $\mathbb{E}_{\pi} [h(X)]$ where π is the stationary distribution. A sequence $(X_t)_t$ produced by a Markov chain Monte Carlo algorithm can thus be employed just as an i.i.d sample.

- Similar to importance sampling, we can approximate the expectation $\mathbb{E}_f [h(X)]$ using an alternative proposal distribution g which is the stationary distribution of an ergodic Markov chain. This is the idea behind *Metropolis-Hastings algorithm*.
- Markov Chain Monte Carlo (MCMC) methods are **preferred** in following situations:
 - The target distribution is **high dimensional**. Due to *the curse of dimensionality*, the variance, which is a function of dimension d , will grow exponentially as the dimensionality increases. Moreover, many high dimensional joint distributions are usually not represented in explicit function form due to their complicated partition functions. In this situation, finding a proposal distribution that is close to the target distribution in high dimensional space is also very challenging.
 - Some stochastic optimization algorithms **naturally produce Markov chain structures**. It is a general fact that the use of Markov chains allows for a greater scope than the methods presented in vanilla Monte Carlo.
 - Vanilla Monte Carlo and MCMC algorithms both satisfy the $O(1/\sqrt{n})$ convergence requirement for the approximation of J . There are thus many instances where a specific MCMC algorithm dominates, variance-wise, the corresponding Monte Carlo proposal.

2.2 Metropolis-Hastings Algorithm

- Consider the following energy-based model

$$\pi(\mathbf{x}) = \frac{1}{Z} \exp(-h(\mathbf{x})),$$

$$\text{where } Z = \int_{\mathcal{X}} \exp(-h(\mathbf{x})) d\mathbf{x}$$

is the partition function in high dimensional space.

- The basic idea for Metropolis algorithm is to simulate π using the stationary distribution g from a Markov chain. Compare to analysis of Markov chain itself, which often starts from a

known transition kernel, the Metropolis algorithm starts from a known stationary distribution g and is interested in how to prescribe an efficient transition kernel to reach the equilibrium.

- Starting from an initial configuration \mathbf{X}_0 , the **Metropolis Algorithm** iteratively repeats the following steps [Liu, 2001]:

1. Propose a random unbiased "perturbation" of \mathbf{X}_t so as to generate a new configuration \mathbf{X}' . \mathbf{X}' is seen as generated by a *symmetric* probability transition kernel (or ***proposal function***) $K(\mathbf{X}_t, \mathbf{X}')$, i.e. $K(x, y) = K(y, x)$. Then calculate the change $\Delta h := h(\mathbf{X}') - h(\mathbf{X}_t)$;
2. Generate a uniform random variable $U \in \mathcal{U}[0, 1]$. Accept $\mathbf{X}_{t+1} = \mathbf{X}'$ if

$$U \leq \pi(\mathbf{X}') / \pi(\mathbf{X}_t) := \exp(-\Delta h);$$

Otherwise, accept $\mathbf{X}_{t+1} = \mathbf{X}_t$.

- Heuristically, the Metropolis Algorithm is based on a "*trial-and-error*" strategy: at each iteration, a random perturbation of \mathbf{X}_t is generated by Markov chain. Then the gain is computed for this new sample. If the gain is large enough, it will be accepted by high probability. Otherwise, we keep using the old sample \mathbf{X}_t .
- The requirement for symmetric transition kernel means that the chance of getting \mathbf{X}' from \mathbf{X}_t and the chance of getting \mathbf{X}_t from \mathbf{X}' are equal. This could avoid the "trend bias" at the proposal stage.
- Hastings generalize the idea of Metropolis Algorithm by removing the requirement for the symmetric proposal function K . In Hastings' generalization, it only requires that $K(x, y) > 0 \Leftrightarrow K(y, x) > 0$.
- The **Metropolis-Hastings Algorithm** is described as below:

1. Given current configuration \mathbf{X}_t , draw \mathbf{Y} from the proposal function $K(\mathbf{X}_t, \mathbf{Y})$.
2. Compute the ***Hastings ratio*** (*acceptance function*):

$$r(\mathbf{X}_t, \mathbf{Y}) := \frac{\pi(\mathbf{Y}) K(\mathbf{Y}, \mathbf{X}_t)}{\pi(\mathbf{X}_t) K(\mathbf{X}_t, \mathbf{Y})} \quad (25)$$

3. (***Metropolis Rejection***) Generate a uniform random variable $U \in \mathcal{U}[0, 1]$. Accept $\mathbf{X}_{t+1} = \mathbf{Y}$ if

$$U \leq \alpha(\mathbf{X}_t, \mathbf{Y}),$$

where

$$\alpha(\mathbf{X}_t, \mathbf{Y}) = \min \{1, r(\mathbf{X}_t, \mathbf{Y})\}$$

4. Otherwise, accept $\mathbf{X}_{t+1} = \mathbf{X}_t$.

Clearly if K is symmetric $K(\mathbf{Y}, \mathbf{X}_t) = K(\mathbf{X}_t, \mathbf{Y})$, the ratio (??) is equal to $\exp(-\Delta h)$. The Metropolis Algorithm is a special case of Metropolis-Hastings Algorithm.

- The intuition behind the ratio $K(\mathbf{y}, \mathbf{x}_t) / K(\mathbf{x}_t, \mathbf{y})$ is that it compensates the "flow bias" from the proposal distribution.

- **Theorem 2.1** Let $(\mathbf{X}_t)_t$ be the chain produced by the Metropolis-Hastings algorithm. For every conditional distribution $q(\mathbf{y}|\mathbf{x}) = K(\mathbf{x}, \mathbf{y})$, whose support includes \mathcal{X} ,

1. the kernel K of the chain satisfies the **detailed balance condition** with π ;
2. π is a stationary distribution of the chain.

Proof: For (1), we see that the transition probability of $(\mathbf{X}_t)_t$ is computed as

$$\begin{aligned}
A(\mathbf{x}, \mathbf{y}) &= K(\mathbf{x}, \mathbf{y}) \alpha(\mathbf{x}, \mathbf{y}) \\
&= K(\mathbf{x}, \mathbf{y}) \min \left\{ 1, \frac{\pi(\mathbf{y}) K(\mathbf{y}, \mathbf{x})}{\pi(\mathbf{x}) K(\mathbf{x}, \mathbf{y})} \right\} \\
&= \frac{\min \{ \pi(\mathbf{x}) K(\mathbf{x}, \mathbf{y}), \pi(\mathbf{y}) K(\mathbf{y}, \mathbf{x}) \}}{\pi(\mathbf{x})} := \frac{\delta(\mathbf{x}, \mathbf{y})}{\pi(\mathbf{x})}. \\
&= \pi(\mathbf{y}) \min \left\{ \frac{K(\mathbf{x}, \mathbf{y})}{\pi(\mathbf{y})}, \frac{K(\mathbf{y}, \mathbf{x})}{\pi(\mathbf{x})} \right\}
\end{aligned} \tag{26}$$

We can see that since $\delta(\mathbf{x}, \mathbf{y}) = \delta(\mathbf{y}, \mathbf{x})$ the detailed balance equation holds:

$$\pi(\mathbf{x}) A(\mathbf{x}, \mathbf{y}) = \pi(\mathbf{y}) A(\mathbf{y}, \mathbf{x}). \quad \blacksquare \tag{27}$$

For (2), it follows that when detailed balance equation is satisfied, then π is a stationary distribution for Markov chain $(\mathbf{X}_t)_t$. \blacksquare

The stationarity of π is therefore established for *almost any conditional distribution* q , a fact which indicates the **universality** of Metropolis-Hastings algorithms.

- We can check on the properties of the Markov chain $(\mathbf{X}_t)_t$ from Metropolis-Hastings algorithm:
 - (**Irreducibility**): To make sure the *irreducibility* to hold, the domain \mathcal{X} of target distribution π need to be **connected**. Otherwise, the Markov chain will not be able to reach other connected components. In other words, the stationary distribution π is **not multimodal distribution**.
 - (**Positive recurrent**): A sufficient condition for both the *irreducibility* and *positive recurrence* to hold is the positivity of transition kernel:

$$q(\mathbf{y}|\mathbf{x}) = K(\mathbf{x}, \mathbf{y}) > 0 \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}. \tag{28}$$

This makes sure that every state can be visited in one step. Thus the Markov chain is positive recurrent. In fact, it can be shown in [Robert and Casella, 1999] that if $(\mathbf{X}_t)_t$ is π -irreducible, then it is *Harris recurrent*.

Lemma 2.2 [Robert and Casella, 1999]

If $(\mathbf{X}_t)_t$ from MH algorithm is π -irreducible, then it is *Harris recurrent*.

- (**Aperiodic**): Unlike the Rejection Sampling, the Metropolis Rejection does not simply regenerate a new proposal repeatedly until some of them is accepted. Instead, it reuse the value from the old state. An attempt to make Metropolis Rejection like the Rejection sampling will destroy the property that **this update will preserve the distribution** π

A *sufficient condition* for **aperiodic** property to hold is to have $P\{\mathbf{X}_{t+1} = \mathbf{X}_t\} > 0$, and thus

$$P\{\pi(\mathbf{X}_t) K(\mathbf{X}_t, \mathbf{Y}) \leq \pi(\mathbf{Y}) K(\mathbf{Y}, \mathbf{X}_t)\} < 1 \quad (29)$$

This is **critical** to guarantee the aperiodic property of the Markov chain, which is essential for *Ergodic theorem* to hold. Note that this condition also make sure K is *not the transition kernel* for the Markov chain induced by Metropolis-Hastings algorithm.

- (**Reversibility**): From the theorem, we see that $(\mathbf{X}_t)_t$ is **time-reversible** with π as its **invariant distribution**.
- Note that due to the rejection rule, $A(\mathbf{x}, \mathbf{y}) \neq K(\mathbf{x}, \mathbf{y})$, i.e. K is not the transition kernel for the Markov chain induced by Metropolis-Hastings algorithm.
- In fact, as long as $A(\mathbf{x}, \mathbf{y}) = \pi(\mathbf{y})\delta(\mathbf{x}, \mathbf{y})$ for some symmetric function δ , the detailed balance equation will hold. The challenge, however, lies in the constraint on $\delta(\mathbf{x}, \mathbf{y})$ so that $\int A(\mathbf{x}, \mathbf{y})d\mathbf{y} = \int \pi(\mathbf{y})\delta(\mathbf{x}, \mathbf{y})d\mathbf{y} = 1$.

For instance, We can replace the acceptance function (25) with a more general form

$$r(\mathbf{x}, \mathbf{y}) := \frac{\delta(\mathbf{x}, \mathbf{y})}{\pi(\mathbf{x}) K(\mathbf{x}, \mathbf{y})}$$

where $\delta(\mathbf{x}, \mathbf{y})$ is any symmetric function that makes $r(\mathbf{x}, \mathbf{y}) \leq 1$. In this case, the probability of generating \mathbf{y} given \mathbf{x} is

$$A(\mathbf{x}, \mathbf{y}) := K(\mathbf{x}, \mathbf{y})r(\mathbf{x}, \mathbf{y}) = \frac{\delta(\mathbf{x}, \mathbf{y})}{\pi(\mathbf{x})}.$$

Since $\delta(\mathbf{x}, \mathbf{y}) = \delta(\mathbf{y}, \mathbf{x})$, we have $\pi(\mathbf{x})A(\mathbf{x}, \mathbf{y}) = \pi(\mathbf{y})A(\mathbf{y}, \mathbf{x})$.

- Finally, we see that the convergence of estimator (24) based on ergodic theorem.

Theorem 2.3 [Robert and Casella, 1999]

Suppose that the Metropolis-Hastings Markov chain $(\mathbf{X}_t)_t$ is π -irreducible.

1. If $h \in L_1(\pi)$, then

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T h(\mathbf{X}_t) = \mathbb{E}_\pi[h(\mathbf{X})] \quad (30)$$

2. If, in addition, $(\mathbf{X}_t)_t$ is **aperiodic**, then

$$\lim_{t \rightarrow \infty} \left\| \int_{\mathcal{X}} K^t(\mathbf{x}, \cdot) d\mu(\mathbf{x}) - \pi \right\|_{TV} = 0 \quad (31)$$

for every **initial distribution** μ , where $K^t(\mathbf{x}, \cdot)$ denotes the kernel for t transitions.

- **Corollary 2.4** [Robert and Casella, 1999]

The conclusions of Theorem 2.3 hold if the Metropolis-Hastings Markov chain $(\mathbf{X}_t)_t$ has conditional density $q(\mathbf{y}|\mathbf{x}) = K(\mathbf{x}, \mathbf{y})$ that satisfies (28) and (29).

3 Some Special Metropolis-Hastings Algorithms

3.1 Random-walk Metropolis-Hastings

- A natural approach for the practical construction of a Metropolis-Hastings algorithm is to take into account the value previously simulated to generate the following value; that is, to consider a **local exploration** of the **neighborhood** of the current value of the Markov chain. This idea is already used in algorithms such as the *simulated annealing algorithm* and the *stochastic gradient method*.
- The *Random-walk Metropolis-Hastings algorithm* is very commonly used when **no prior knowledge** is available. In this case, the preferred proposal is just a *uniformly local* move, since it is **uninformative** and **unbiased**.
- For Euclidean space $\mathcal{X} \subset \mathbb{R}^d$, the *random walk* model is [Robert and Casella, 1999]:

$$\mathbf{Y}_t = \mathbf{X}_t + \boldsymbol{\epsilon}_t, \quad (32)$$

where the noise $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$ are independent from the starting position \mathbf{X}_t . The transition probability $q(\mathbf{y}|\mathbf{X}_t) = \mathcal{N}(\mathbf{X}_t, \sigma^2 \mathbf{I}_d) := g(\|\mathbf{y} - \mathbf{X}_t\|_2)$, where $g(\cdot)$ is a symmetric function $g(-r) = g(r)$.

- In general, define a **symmetric distribution** g_σ so that $q(\mathbf{y}|\mathbf{x}) = g_\sigma(\|\mathbf{y} - \mathbf{x}\|_2) > 0$ for all $\|\mathbf{y} - \mathbf{x}\|_2 < \delta$. Here σ controls the "**range**" of the exploration. The random walk with transition kernel as g_σ will produce an ergodic Markov chain.

The most common distributions g in this setup are the **uniform distributions on spheres** centered at the origin or standard distributions like the **normal** and the **Student's t distributions**. *All these distributions usually need to be scaled.*

- The **Random-walk Metropolis-Hastings** is described as below:

1. Draw $\boldsymbol{\epsilon}_t$ from $g_\sigma(\boldsymbol{\epsilon})$ and set $\mathbf{Y} := \mathbf{X}_t + \boldsymbol{\epsilon}_t$;
2. Compute the ratio:

$$r(\mathbf{X}_t, \mathbf{Y}) := \frac{\pi(\mathbf{Y})}{\pi(\mathbf{X}_t)} \quad (33)$$

3. Accept $\mathbf{X}_{t+1} = \mathbf{Y}$ with probability

$$\alpha(\mathbf{X}_t, \mathbf{Y}) = \min \{1, r(\mathbf{X}_t, \mathbf{Y})\}$$

4. Otherwise, accept $\mathbf{X}_{t+1} = \mathbf{X}_t$.

- Random-walk Metropolis is not only simple to implement, it also has a particularly nice **intuition**. The **proposal distribution** is biased towards **large volumes** ($g_\sigma(\cdot)$), and hence the *tails* of the target distribution, while the *Metropolis correction* **rejects** those proposals that jump into neighborhoods where **the density is too small** ($\pi(\mathbf{x}^{(t)})$). The combined procedure then preferentially selects out those proposals that fall into neighborhoods of high probability mass, concentrating towards the typical set as desired.

- **Theorem 3.1** [Robert and Casella, 1999]

Consider a **symmetric** density π which is **α -log-concave** with associated constant α , i.e.

$\pi(\mathbf{x}) = \exp(-h(\mathbf{x}))$ for $h(\cdot)$ is a convex function, and

$$\log \pi(\mathbf{x}) - \log \pi(\mathbf{y}) \geq \alpha \|\mathbf{y} - \mathbf{x}\| \quad (34)$$

for $\|\mathbf{x}\|$ large enough. If the density g is **positive** and **symmetric**, the chain $(\mathbf{X}_t)_t$ of Random-walk Metropolis-Hastings is **geometrically ergodic**. If g is not symmetric, a sufficient condition for geometric ergodicity is that $g(r)$ be bounded by $b \exp(-\alpha |r|)$ for a sufficiently large constant b .

3.2 Independent Metropolis-Hastings (IMH) Algorithm

- **Independent Metropolis-Hastings (IMH) Algorithm:**

1. Generate \mathbf{Y}_t from $g(\mathbf{y})$;
2. Compute ratio:

$$r(\mathbf{X}_t, \mathbf{Y}) := \frac{\pi(\mathbf{Y}_t) g(\mathbf{X}_t)}{\pi(\mathbf{X}_t) g(\mathbf{Y}_t)} = \frac{w(\mathbf{Y}_t)}{w(\mathbf{X}_t)} \quad (35)$$

where $w(\mathbf{x}) = \frac{\pi(\mathbf{x})}{g(\mathbf{x})}$ is the *importance weight*.

3. Accept $\mathbf{X}_{t+1} = \mathbf{Y}_t$ with probability

$$\alpha(\mathbf{X}_t, \mathbf{Y}) = \min \{1, r(\mathbf{X}_t, \mathbf{Y})\}$$

4. Otherwise, accept $\mathbf{X}_{t+1} = \mathbf{X}_t$.

- *Independent Metropolis-Hastings (IMH)* choose the proposal transition function $K(\mathbf{x}, \mathbf{y})$ as an **independent density** $g(\mathbf{y})$. That is, the proposal move \mathbf{Y}_t is generated **independently** from previous sample \mathbf{X}_t .
- This method is an **alternative** to the *Rejection Sampling* and *Importance Sampling*. As with Rejection sampling, the performance of *IMH* depends on how close the proposal distribution g to the target distribution π . It is also suggested that g has longer tail than π for robustness of performance.

It can be shown that if the envelop condition $\pi(\mathbf{x}) \leq M g(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$ holds, then the Metropolis-Hastings Markov chain $(\mathbf{X}_t)_t$ is ergodic, thus the convergence to $\mathbb{E}_\pi[h]$ is guaranteed.

Theorem 3.2 *The Independent Metropolis-Hastings algorithm produces a **uniformly ergodic** chain if there exists a constant M such that*

$$\pi(\mathbf{x}) \leq M g(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}. \quad (36)$$

In this case,

$$\|K^t(\mathbf{x}, \cdot) - \pi\|_{TV} \leq 2 \left(1 - \frac{1}{M}\right)^t$$

where $\|\cdot\|_{TV}$ denotes the total variation norm. On the other hand, if for every M , there exists a set of positive measure where (36) does not hold, $(\mathbf{X}_t)_t$ is not even geometrically ergodic.

- Compared to *Rejection Sampling* and *Importance Sampling*, the *Independent Metropolis-Hastings (IMH)* has **larger expected acceptance probability** (at least $1/M$ when the chain is stationary) [Robert and Casella, 1999]. Thus IMH is more **sample efficient**.

3.3 Configurational bias Monte Carlo

- The *Configurational bias Monte Carlo (CBMC)* can be seen as a **SIS-based IMH**. Assume that we have auxiliary distributions $\pi_s(\mathbf{x}_{1:s})$ where $\mathbf{x}_{1:s} = [x_1, \dots, x_s]$ at each time s . $\mathbf{x} := \mathbf{x}_{1:d}$. The proposal sampling probability $g_s(x_s | \mathbf{x}_{1:(s-1)})$ for all s . Here the proposal joint distribution can be computed sequentially

$$g(\mathbf{x}) = g_1(x_1) \prod_{s=2}^d g_s(x_s | \mathbf{x}_{1:(s-1)}) \quad (37)$$

- The *Configurational bias Monte Carlo (CBMC)* is described as below:

1. Generate \mathbf{Y} according to $g(\mathbf{y})$ via the sequential process in (37).
2. Compute the *importance weight*:

$$w(\mathbf{Y}) = \frac{\pi(\mathbf{Y})}{g(\mathbf{Y})} = \frac{\pi_1(Y_1)}{g_1(Y_1)} \prod_{s=2}^d \frac{\pi_s(\mathbf{Y}_{1:s})}{\pi_s(\mathbf{Y}_{1:(s-1)}) g_s(Y_s | \mathbf{Y}_{1:(s-1)})} \quad (38)$$

Similarly compute the importance weight for \mathbf{X}_t as $w(\mathbf{X}_t)$.

3. Accept $\mathbf{X}_{t+1} = \mathbf{Y}$ with probability

$$\alpha(\mathbf{X}_t, \mathbf{Y}) = \min \left\{ 1, \frac{w(\mathbf{Y})}{w(\mathbf{X}_t)} \right\}$$

4. Otherwise, accept $\mathbf{X}_{t+1} = \mathbf{X}_t$.

- We can modify this process to make **stage-wise acceptance decision**. Suppose that, at time $t-1$, $\mathbf{X}_t = (X_1^{(t)}, \dots, X_d^{(t)})$ is accepted. Then at s -th stage, we accept $\mathbf{Y}_{1:s}$ with probability

$$\alpha_s(\mathbf{X}_{1:s}^{(t)}, \mathbf{Y}_{1:s}) = \min \left\{ 1, \frac{u_s(\mathbf{Y}_{1:s})}{u_s(\mathbf{X}_{1:s}^{(t)})} \right\} \quad (39)$$

where

$$u_s(\mathbf{Y}_{1:s}) = \frac{\pi_s(\mathbf{Y}_{1:s})}{\pi_s(\mathbf{Y}_{1:(s-1)}) g_s(Y_s | \mathbf{Y}_{1:(s-1)})}.$$

That is, the acceptance rate at each stage is equal to the *ratio of incremental weights* between proposed state and the old state.

If **rejected**, we go back to the first stage to rebuild the whole configuration.

Note that it is **not necessary** to make accept-reject decision **for each stage**.

- Note that the **acceptance rate** for multi-stage CBMC is *less than* the original CBMC, since $\prod_s \min \{1, r_s\} \leq \min \{1, \prod_s u_s\}$.

References

Jun S Liu. *Monte Carlo strategies in scientific computing*, volume 10. Springer, 2001.

Christian P Robert and George Casella. *Monte Carlo statistical methods*, volume 2. Springer, 1999.

Sheldon M Ross. *Introduction to probability models*. Academic press, 2014.