

# [AI2613 Lecture 4] Metropolis Algorithm, Countable Infinite Markov Chain

March 18, 2022

## 1 Reversible Markov Chains

A Markov chain  $P$  over state space  $[n]$  is (*time*) *reversible* if there exists some distribution  $\pi$  satisfying

$$\forall i, j \in [n], \pi(i)P(i, j) = \pi(j)P(j, i).$$

This family of identities is called *detailed balance conditions*. Moreover, the distribution  $\pi$  must be a stationary distribution of  $P$ . To see this, note that

$$\pi^\top P(j) = \sum_{i \in [n]} \pi(i)P(i, j) = \sum_{i \in [n]} \pi(j)P(j, i) = \pi(j).$$

The name *reversible* comes from the fact that for any sequence of variables  $X_0, X_1, \dots, X_t$  following the chain which start from the stationary distribution, the distribution of  $(X_0, X_1, \dots, X_{t-1}, X_t)$  is identical to the distribution of  $(X_t, X_{t-1}, \dots, X_1, X_0)$ , namely for all  $x_0, x_1, \dots, x_t \in [n]$ ,

$$\begin{aligned} & \Pr_{X_0 \sim \pi} [X_0 = x_0, X_1 = x_1, \dots, X_t = x_t] \\ &= \pi(x_0)P(x_0, x_1) \cdots P(x_{t-1}, x_t) \\ &= \pi(x_t)P(x_t, x_{t-1}) \cdots P(x_1, x_0) \\ &= \Pr_{X_0 \sim \pi} [X_0 = x_t, X_1 = x_{t-1}, \dots, X_t = x_0] \end{aligned}$$

We will study reversible chains since their transition matrices are essentially *symmetric* in some sense, so many powerful tools in linear algebra apply. We will also see that reversible chains are general enough for most of our (algorithmic) applications. You can verify that the the random walks on the hypercube is reversible Markov chains with respect to uniform distribution.

Recall the two conditions of FTMC: irreducibility and aperiodicity. Since the transition graph is undirected if we only consider the connectivity, irreducibility is equivalent to the connectivity of the transition graph. Aperiodicity, on the other hand, is equivalent to that the graph is *not* bipartite.

## 2 The Metropolis Algorithm

Given a distribution  $\pi$  over a state space  $\Omega$ , how can we design a Markov chain  $P$  so that  $\pi$  is the stationary distribution of  $P$ ? The *Metropolis algorithm* provides a way to achieve the goal as long as the transition graph  $G$  is connected and undirected.

Let  $\Delta$  be the maximum degree of the transition graph except selfloop (that is  $\Delta \triangleq \max_{u \in [n]} \sum_{v \neq u \in [n]} \mathbb{1}[(u, v) \in E]$ ). We describe the following

process to construct a transition matrix  $P$ : Choose  $k \in [\Delta + 1]$  uniformly at random. For any  $i \in [n]$ , let  $\{j_1, j_2, \dots, j_d\}$  be the  $d$  neighbours of  $i$ . We consider the transition at state  $i$ :

- If  $d + 1 \leq k \leq \Delta + 1$ , do nothing.
- If  $k \leq d$ ,
  - propose to move from  $i$  to  $j_k$ .
  - accept the proposal with probability  $\min \left\{ \frac{\pi(j_k)}{\pi(i)}, 1 \right\}$ .

Then the transition matrix is, for  $i, j \in [n]$ ,

$$P(i, j) = \begin{cases} \frac{1}{\Delta+1} \min \left\{ \frac{\pi(j)}{\pi(i)}, 1 \right\}, & \text{if } i \neq j; \\ 1 - \sum_{k \neq i} P(i, k), & \text{if } i = j. \end{cases}$$

We can verify that  $P$  is reversible with respect to  $\pi$ :

$\forall i, j \in \Omega :$

$$\pi(i)P(i, j) = \pi(i) \cdot \frac{1}{\Delta + 1} \min \left\{ \frac{\pi(j)}{\pi(i)}, 1 \right\} = \frac{\min \{ \pi(i), \pi(j) \}}{\Delta + 1} = \pi(j)P(j, i).$$

**Example 1** We give a toy example to show how the algorithm works. Consider a graph with 3 vertices  $\{a, b, c\}$ . There are undirected edges between  $(a, b)$ ,  $(b, c)$  and  $(a, c)$  and selfloops for each vertex. In this situation,  $\Delta = 2$ . If we want to design a transition matrix  $P$  with stationary distribution  $(\frac{1}{2}, \frac{1}{3}, \frac{1}{6})$ , by Metropolis algorithm we have

$$\begin{aligned} P(a, b) &= \frac{1}{2+1} \cdot \frac{2}{3} = \frac{2}{9}, \\ P(a, c) &= \frac{1}{2+1} \cdot \frac{1}{3} = \frac{1}{9}, \\ P(a, a) &= 1 - \frac{1}{9} - \frac{2}{9} = \frac{2}{3}. \end{aligned}$$

The advantage of the Metropolis algorithm is that we do not need to know  $\pi$  in order to implement the algorithm. We only need to know the quantity  $\frac{\pi(j)}{\pi(i)}$ , which is much easier to compute in many applications.

### 3 Sampling Proper Colorings

Let's consider the problem of sampling proper colorings. Given a graph  $G = (V, E)$ , we want to color the vertices using  $q$  colors under the condition that no two adjacent vertices share the same color. More formally, a coloring of  $G$  is a mapping  $c : V \mapsto [q]$ , and we call it *proper* iff  $\forall \{u, v\} \in E, c(u) \neq c(v)$ . The problem is NP-hard in general. However, for  $q > \Delta$  there's always at least one suitable solution and can be easily obtained by a greedy algorithm, where  $\Delta$  is the maximum degree of the graph.

Let  $\mathcal{C}$  be the set of all proper colorings. We want to sample uniformly on  $\mathcal{C}$ . Consider the following Markov chain (assume that the start state is a proper coloring):

- Pick  $v \in V$  and  $c \in [q]$  uniformly at random.
- Recolor  $v$  with  $c$  if the modified coloring is still proper.

The chain is aperiodic since selfloops exist in the walk. For  $q \geq \Delta + 2$ , the chain is irreducible.

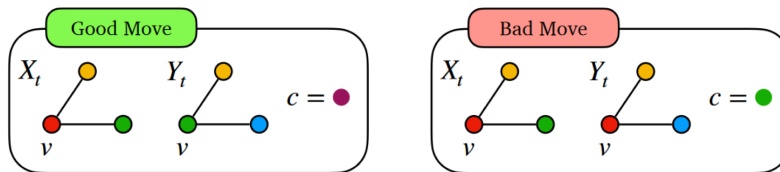
We verify that this Markov chain is reversible, that is, there exist a distribution  $\pi$  that for any  $\sigma, \sigma' \in \mathcal{C}$ ,  $\pi(\sigma)P(\sigma, \sigma') = \pi(\sigma')P(\sigma', \sigma)$ .

- If  $\sigma = \sigma'$ , it is obvious that  $P(\sigma, \sigma') = P(\sigma', \sigma)$ .
- If  $\sigma$  and  $\sigma'$  differ at more than 1 vertices, then  $P(\sigma, \sigma') = P(\sigma', \sigma) = 0$ .
- If  $\sigma$  and  $\sigma'$  differ at exactly 1 vertex, then  $P(\sigma, \sigma') = \frac{1}{n} \cdot \frac{1}{q} = P(\sigma', \sigma)$ .

Thus, it is reversible with respect to the uniform distribution and furthermore, uniform distribution on  $\mathcal{C}$  is the stationary distribution of the Markov chain defined above.

Suppose  $X_t, Y_t$  are two proper colorings. We define the distance  $d(X_t, Y_t)$  as their Hamming distance, i.e. the number of vertices colored differently in two colorings. Our coupling of two chains is that we always choose the same  $v, c$  in each step. The distance between two colorings can change at most 1 since only  $v$  is affected. The possible changes can be divided into two cases:

- Good move:  $X_t(v) \neq Y_t(v)$ , and both change into  $c$  successfully. This will decrease distance by 1.
- Bad move:  $X_t(v) = Y_t(v)$ , and exactly one change succeeds. This will increase distance by 1.



Consider the probabilities of two types of moves. For good moves, w.p.  $\frac{d(X_t, Y_t)}{n}$ ,  $X_t(v) \neq Y_t(v)$ , and there are at least  $q - 2\Delta$  choices of  $c$  to make it a good move. So

$$\begin{aligned} \Pr [d(X_{t+1}, Y_{t+1}) = d(X_t, Y_t) - 1] &= \Pr [(v, c) \text{ is a good move}] \\ &\geq \frac{d(X_t, Y_t)}{n} \cdot \frac{q - 2\Delta}{q}. \end{aligned}$$

For bad moves, there exists a neighbor  $w$  of  $v$  such that its color is different in two colorings, and in one coloring  $w$  is of color  $c$ . Note that there are at

most  $2\Delta$  choices of  $c$  to make it a bad move. So we have

$$\begin{aligned} \Pr [d(X_{t+1}, Y_{t+1}) = d(X_t, Y_t) + 1] &= \Pr_{(v,c) \in V \times [q]} [(v, c) \text{ is a bad move}] \\ &\leq \frac{d(X_t, Y_t)}{n} \cdot \frac{2\Delta}{q}. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbf{E} [d(X_{t+1}, Y_{t+1}) | (X_t, Y_t)] &= d(X_t, Y_t) + \Pr [d(X_{t+1}, Y_{t+1}) = d(X_t, Y_t) + 1] \\ &\quad - \Pr [d(X_{t+1}, Y_{t+1}) = d(X_t, Y_t) - 1] \\ &\leq d(X_t, Y_t) + \frac{d(X_t, Y_t)}{n} \cdot \frac{2\Delta}{q} - \frac{d(X_t, Y_t)}{n} \cdot \frac{q - 2\Delta}{q} \\ &\leq d(X_t, Y_t) \left( 1 - \frac{q - 4\Delta}{nq} \right). \end{aligned}$$

In the case  $q > 4\Delta$ ,

$$D_{\text{TV}}(X_{t+1}, Y_{t+1}) \leq \mathbf{E} [d(X_{t+1}, Y_{t+1})] \leq \left( 1 - \frac{1}{nq} \right)^t n \leq \varepsilon.$$

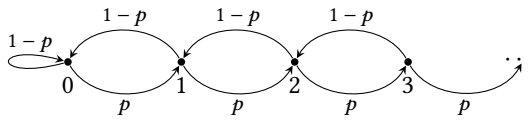
The mixing time is therefore bounded by

$$\tau_{\text{mix}}(\varepsilon) \leq nq \log \frac{n}{\varepsilon}.$$

#### 4 Countably Infinite Markov Chains

We have proved that finite Markov chain must have a stationary distribution using Perron-Frobenius Theorem. However, when the Markov chain has infinite states, even it's countable infinite, there is something going wrong.

Consider the following random walk on  $\mathbb{N}$ . The state space is  $\mathbb{N}$  and at each state  $i$ , go to  $i + 1$  w.p.  $p$  and go to  $i - 1$  w.p.  $1 - p$  (if  $i = 0$ , w.p.  $1 - p$  stay put).



Let  $\pi$  be the stationary distribution of this Markov chain (if there exists a stationary distribution). We have that

$$\pi(0) = \pi(0)(1 - p) + \pi(1)p \quad \implies \pi(1) = \frac{p}{1 - p} \pi(0),$$

$$\pi(1) = \pi(0)p + \pi(2)(1 - p) \quad \implies \pi(2) = \frac{p}{1 - p} \pi(1),$$

...

$$\pi(i) = \pi(i - 1)p + \pi(i + 1)(1 - p) \quad \implies \pi(i + 1) = \frac{p}{1 - p} \pi(i).$$

...

Note that  $\pi$  is a distribution, so  $\sum_{i=0}^{\infty} \pi(i) = 1$ . Then, we have

- If  $p < \frac{1}{2}$ , that is,  $\frac{p}{1-p} < 1$ , then  $\sum_{i=0}^{\infty} \left(\frac{p}{1-p}\right)^i \pi(0) = 1$ . By direct calculation we have  $\pi(0) = \frac{1-2p}{1-p}$  and  $\pi(i) = \left(\frac{p}{1-p}\right)^i \frac{1-2p}{1-p}$  for  $i \in \mathbb{N}$ .
- If  $p > \frac{1}{2}$ , then  $\frac{p}{1-p} > 1$ . When  $i \rightarrow \infty$ , if  $\pi(0) \neq 0$ ,  $\pi(i) \rightarrow \infty$ . This yields that  $\pi(0) = \pi(1) = \dots = \pi(i) = \dots = 0$ . The Markov chain doesn't have a stationary distribution in this case.
- If  $p = \frac{1}{2}$ ,  $\frac{p}{1-p} = 1$ . Then  $\pi(0) = \pi(1) = \dots = \pi(i) = \dots$  and  $\sum_{i=0}^{\infty} \pi(0) = 1$ . This yields that  $\pi(0) = 0$  and there is no stationary distribution in this case.

#### 4.1 Recurrence

**Definition 1** For  $i \in \Omega$ , let  $T_i > 0$  be the first hitting time of state  $i$ . Let  $\mathbf{P}_i = \Pr[\cdot | X_0 = i]$ . We say a state  $i$  is recurrent if  $\mathbf{P}_i[T_i < \infty] = 1$ , o.w. we say the state is transient.

Let  $N_i \triangleq \sum_{t=0}^{\infty} \mathbb{1}[X_t = i]$ , then we have the following propositions.

**Proposition 2** If  $i$  is recurrent, then  $\mathbf{P}_i[N_i = \infty] = 1$ .

*Proof.* Assume that  $\mathbf{P}_i[N_i = \infty] < 1$ . Then there exists  $\Omega' \subseteq \hat{\Omega}$  such that  $N_i < \infty$  on  $\Omega'$  and  $\mathbf{P}_i[\Omega'] > 0$ . This means that with probability larger than 0, we will never reach state  $i$  after the last time we visit it. This is in contradiction with the fact that  $i$  is recurrent. □

**Proposition 3** If  $i$  is recurrent and there exists a finite path from  $i$  to  $j$ , then

- $\mathbf{P}_i[T_j < \infty] = 1$ .
- $\mathbf{P}_j[T_i < \infty] = 1$ .
- $j$  is recurrent.

*Proof.*

- Let  $q \triangleq \mathbf{P}_i[\text{reach } j \text{ before returning to } i]$ . Since there is a finite path from  $i$  to  $j$ , we have  $q > 0$  and  $\mathbf{P}_i[\text{visit } i \text{ } n \text{ times before reaching } j] = (1 - q)^n$ .

Assume that  $\mathbf{P}_i[T_j = \infty] = \alpha > 0$ . Then we have  $\mathbf{P}_i[T_j = \infty | N_i = \infty] = \alpha$  since  $\mathbf{P}_i[N_i = \infty] = 1$ . Let  $T_i^n$  be the  $n^{\text{th}}$  time that the chain visits state  $i$ . Then

$$\forall n > 0, \mathbf{P}_i[T_j > T_i^n | N_i = \infty] \geq \mathbf{P}_i[T_j = \infty | N_i = \infty] = \alpha$$

On the otherhand, we have  $\lim_{n \rightarrow \infty} \mathbf{P}_i[T_j > T_i^n | N_i = \infty] = \lim_{n \rightarrow \infty} \mathbf{P}_i[T_j > T_i^n] = \lim_{n \rightarrow \infty} (1 - q)^n = 0$ . This is a contradiction. Thus,  $\mathbf{P}_i[T_j = \infty] = 0$ .

$$T_i \triangleq \min\{t > 0 | X_t = i\}.$$

Recall the probability space of a stochastic process. One can view the outcomes of the probability space is the set of infinite sequence of real numbers between  $[0, 1]$ , namely  $\hat{\Omega} = [0, 1]^{\mathbb{N}}$ . The sigma-algebra can be defined in a way similar to the problem 1 in our first homework. Therefore, the random variable  $T_i$  is therefore a function  $\hat{\Omega} \rightarrow \mathbb{R}$ .

$$\begin{aligned} \mathbf{P}_i[T_j = \infty] &= \mathbf{P}_i[T_j = \infty | N_i = \infty] \cdot \mathbf{P}_i[N_i = \infty] + \mathbf{P}_i[T_j = \infty | N_i < \infty] \cdot \mathbf{P}_i[N_i < \infty] \\ &= \mathbf{P}_i[T_j = \infty | N_i = \infty] \cdot \mathbf{P}_i[N_i = \infty] + \mathbf{P}_i[T_j = \infty | N_i < \infty] \cdot \mathbf{P}_i[N_i < \infty] \end{aligned}$$

- If  $\mathbf{P}_j[T_i = \infty] = p > 0$ , then we have that  $\mathbf{P}_i[T_i = \infty] \geq q \cdot p > 0$ . This is in contradiction with the fact that  $i$  is recurrent.
- If  $\mathbf{P}_j[T_j = \infty] = r > 0$ , then  $\mathbf{P}_i[T_j = \infty] \geq q \cdot r > 0$ . This is in contradiction with the first item of this proposition.

□