

# [CS3958: Lecture 10] Discrete Markov Chains, Fundamental Theorem of Markov Chains, Coupling

Instructor: Chihao Zhang, Scribed by Yulin Wang

December 20, 2022

We already met Markov chains in previous lectures, e.g. the random walk on  $\mathbb{Z}$ . We will formally introduce the model and study its properties. This is an important probabilistic tool for modeling in computer science as well as a powerful tool for designing efficient algorithms.

## 1 Basics of Markov Chains

Consider the random walk on  $\mathbb{Z}$ . One starts at 0 and in each round, he tosses a fair coin to determine the direction of moving: with probability 50% to the left and 50% to the right. We use  $X_t \in \{-1, 1\}$  to denote his movement at time  $t$ , and  $Z_t = Z_{t-1} + X_t$  to denote his position at time  $t$ .

**Definition 1 (Markov Chain)** A sequence of random variables  $X_0, X_1, \dots$  is a Markov chain if and only if for any  $t$  and any  $a_0, a_1, \dots, a_{t+1}$ ,

$$\Pr [X_{t+1} = a_{t+1} \mid X_0 = a_0, X_1 = a_1, \dots, X_t = a_t] = \Pr [X_{t+1} = a_{t+1} \mid X_t = a_t].$$

The  $\{Z_t\}_t \geq 0$  above is a simple Markov chain, since the position at time  $t$  only depends on the position at time  $t - 1$ .

We usually use  $\Omega$  to denote the state space, meaning all possible values that  $X_t$  can take. Today we only consider the case when  $\Omega = [n]$  is finite and the Markov chain is *time-homogeneous*. A time-homogeneous Markov chain can be characterized by a  $n \times n$  transition matrix  $P = (p_{ij})_{i,j \in [n]}$  where  $p_{ij} = \Pr [X_{t+1} = j \mid X_t = i]$  for all  $t \geq 0$  since the transition probabilities do not depend on time. In general, a Markov chain can be equivalently viewed as a random walk on a weighted directed graph where the edge weight from  $i$  to  $j$  means the probability of moving to vertex  $j$  when one is standing at vertex  $i$ . Consider the three-vertex graph on the right.

It corresponds to the Markov chain with transition matrix  $P = (p_{ij}) = \begin{bmatrix} 1/2 & 3/8 & 1/8 \\ 1/3 & 0 & 2/3 \\ 1/4 & 3/4 & 0 \end{bmatrix}$ . We sometimes call the graph the *transition graph* of  $P$ .

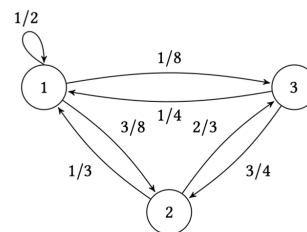
At any time  $t \geq 0$ , we use  $\mu_t \in \Delta_{n-1}$  to denote the distribution of  $X_t$  meaning

$$\mu_t(i) \triangleq \Pr [X_t = i].$$

By the law of total probability,  $\mu_{t+1}(j) = \sum_i \mu_t(i) \cdot p_{ij}$ , we have  $\mu_t^\top P = \mu_{t+1}^\top$ . As a result, we have

**Proposition 2**  $\mu_t^\top = \mu_0^\top P^t$ .

The transition matrix must be a stochastic matrix since  $\sum_j p_{ij} = 1$  for each  $i$ , i.e. the row sum of  $P$  is 1.



This is a useful formula as we can compute the distribution at any time given the initial distribution and the transition matrix.

Sometimes, we will simply denote the transition matrix  $P$  as the Markov chain for convenience.

### 1.1 Stationary Distribution

**Definition 3** A distribution  $\pi$  is a stationary distribution of  $P$  if it remains unchanged in the Markov chain as time progresses, i.e.,

$$\pi^\top P = \pi^\top.$$

One of the major algorithmic applications of Markov chains is the *Markov chain Monte Carlo (MCMC)* method. It is a general method for designing an algorithm to sample from a certain distribution  $\pi$ . The idea of MCMC is

- First design a Markov Chain of which the stationary distribution is the desired  $\pi$ ;
- Simulate the chain from a certain initial distribution for a number of steps and output the state.

Therefore, we hope that the distribution  $\mu_t$  is close to  $\pi$  when  $t$  is large enough.

We already met, and implemented the MCMC method many time – shuffling a deck of cards. After a few operations, a good shuffling rule would produce a card order that is close to the uniform.

One of the main purposes of the course is to understand the MCMC method. Therefore, the following four basic questions regarding stationary distributions are important.

- Does each Markov chain have a stationary distribution?
- If a Markov chain has a stationary distribution, is it unique?
- If the chain has a unique stationary distribution, does  $\mu_t$  always converge to it from any  $\mu_0$ ?
- If  $\mu_t$  always converges to the stationary distribution, what is the rate of convergence?

We will settle the first three questions today. The fourth question, the rate of convergence, will be the topic of coming lectures.

## 2 Fundamental Theorem of Markov Chains

### 2.1 The Existence of Stationary Distribution

We will show that, for every finite Markov chain  $P$ , there exists some  $\pi$  such that  $\pi^\top P = \pi^\top$ . Observe that this is equivalent to "1 is an eigenvalue of

$P^\top$  with a nonnegative eigenvector ( $P^\top \pi = \pi$ ).

We use the following theorem in linear algebra.

**Theorem 4 (Perron-Frobenius Theorem)** . Each nonnegative matrix  $A$  has a nonnegative real eigenvector  $x$  with eigenvalue  $\lambda = \rho(A) = \max \{|\lambda_i|\}$ , where  $\{\lambda_1, \dots, \lambda_n\}$  are eigenvalues of  $A$ .

We will prove the Perron-Frobenius theorem in Section 2.3.

Since  $P$  is a stochastic matrix, we have

$$P \cdot \mathbf{1} = \mathbf{1}.$$

Thus,  $P$  has an eigenvalue 1. Since every eigenvalue of  $P$  is no larger than the row sum, 1 is the largest eigenvalue. Also,  $P^\top$  shares the same characteristic polynomial with  $P$ , which implies the eigenvalues of  $P^\top$  and  $P$  are the same. As a result,  $\rho(P^\top)$  also equals to 1. According to Perron-Frobenius theorem, there exists a nonnegative eigenvector  $\pi$  such that

$$P^\top \pi = \pi,$$

which is equivalent to

$$\pi^\top P = \pi^\top.$$

It then follows that  $\frac{\pi}{\|\pi\|_1}$  is a stationary distribution of  $P$ .

### 2.2 Uniqueness and Convergence

Consider the Markov chain with two states on the right. Clearly, the transition matrix of this Markov chain is

$$P = \begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix}$$

It is easy to verify that

$$\pi = \left( \frac{q}{p+q}, \frac{p}{p+q} \right)^\top$$

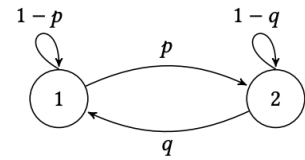
is a stationary distribution of  $P$ . We are going to check whether starting from any  $\mu_0$ , the distribution  $\mu_t$  will always converge to  $\pi$ , i.e.,

$$\lim_{t \rightarrow \infty} \|\mu_0^\top P^t - \pi^\top\| = 0.$$

In our example, the distribution has only two dimensions and the sum of the two components equals to 1, so we only need to check whether the first dimension converges, i.e.,

$$|\mu_0^\top P^t(1) - \pi(1)| \rightarrow 0.$$

Let  $A = (a_{ij})_{i \in [n], j \in [m]}$ . We say  $A$  is nonnegative (resp. positive) if every  $a_{ij} \geq 0$  (resp.  $> 0$ ).



Now we define

$$\begin{aligned}
 \Delta_t &\triangleq |\mu_t(1) - \pi(1)| \\
 &= \left| \mu_{t-1}^T \cdot P(1) - \pi(1) \right| \\
 &= \left| (1-p) \cdot \mu_{t-1}(1) + q \cdot (1 - \mu_{t-1}(1)) - \frac{q}{p+q} \right| \\
 &= \left| (1-p-q) \cdot \mu_{t-1}(1) + q \cdot \left(1 - \frac{1}{p+q}\right) \right| \\
 &= |1-p-q| \cdot \Delta_{t-1}
 \end{aligned}$$

Therefore, we can see that  $\Delta_t \rightarrow 0$  except in the two cases:  $p = q = 0$ ,  $p = q = 1$ .

In fact, the two cases prevent convergence for different reasons.

Let us first consider the case when  $p = q = 0$ . The Markov chain looks like:



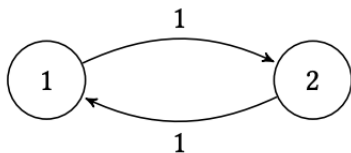
The transition walk graph is disconnected, so it can be partitioned into two disjoint components. Since each component is still a Markov chain, each of them has its own stationary distribution. Notice that any convex combination of these small distributions is a stationary distribution for the whole Markov chain. It immediately follows that in this case the stationary distribution is not unique. It gives a negative answer to the second question.

This observation motivates us to define the following:

**Definition 5 (Irreducibility)** *A finite Markov chain is irreducible if its transition graph is strongly connected.*

If the transition graph of  $P$  is not strongly connected, we say  $P$  is *reducible*.

When  $p = q = 1$ , the Markov chain looks like this:



This transition graph is bipartite. It is easy to see that  $(\frac{1}{2}, \frac{1}{2})$  is the unique stationary distribution of it. However, for  $\mu_0 = (1, 0)$ , one can see that

$\mu_t$  oscillates between "left" and "right". Therefore, the answer to the third question is no.

This phenomenon is captured by the following notion:

**Definition 6 (Aperiodicity)** *A Markov chain is aperiodic if for any state  $v$ , it holds that*

$$\gcd \{ |c| \mid c \in C_v \} = 1,$$

where  $C_v$  denotes the set of the directed cycles containing  $v$  in the transition graph.

Otherwise, we say the chain *periodic*.

We have the following important theorem.

**Theorem 7 (Fundamental theorem of Markov chains)** *If a finite Markov chain  $P \in \mathbb{R}^{n \times n}$  is irreducible and aperiodic, then it has a unique stationary distribution  $\pi \in \mathbb{R}^n$ . Moreover, for any distribution  $\mu \in \mathbb{R}^n$ ,*

$$\lim_{t \rightarrow \infty} \mu^\top P^t = \pi^\top.$$

Although there are many ways to prove the theorem, we will present one based on the so called *coupling* argument, which will be quite useful in answering the fourth question.

### 2.3 Proof of Perron-Frobenius Theorem

Most proofs in the section are from [Mey00]. We first prove the Perron-Frobenius theorem for positive matrices. Then we use this theorem and Lemma 9 to prove Theorem 4.

In the following statement, we use  $|\cdot|$  to denote a matrix or vector of absolute values, i.e.,  $|A|$  is the matrix with entries  $|a_{ij}|$ . We say a vector or matrix is larger than  $\mathbf{0}$  if all its entries are larger than 0 and denote it by  $A > \mathbf{0}$ . We define the operation  $\geq$ ,  $\leq$  and  $<$  for vectors and matrices similarly.

**Theorem 8 (Perron-Frobenius Theorem for Positive Matrices)** *Each positive matrix  $A > \mathbf{0}$  has a positive real eigenvalue  $\rho(A)$ , and  $\rho(A)$  has a corresponding positive eigenvector.*

*Proof.* We first prove that  $\rho(A) > 0$ . If  $\rho(A) = 0$ , then all the eigenvalues of  $A$  is 0 which is equivalent to that  $A$  is nilpotent. This is impossible since every  $a_{ij} > 0$ . Thus  $\rho(A) > 0$  for positive matrix  $A$ .

Assume that  $\lambda$  is the eigenvalue of  $A$  that  $|\lambda| = \rho(A)$ . Then we have

$$|\lambda||x| = |\lambda x| = |Ax| \leq |A||x| = A|x|.$$

Then we show that  $|\lambda||x| < A|x|$  is impossible. Let  $z = A|x|$  and  $y = z - \rho(A)|x|$ . Assume that  $y \neq \mathbf{0}$ , We have that  $Ay > \mathbf{0}$ . There must exist some

$\epsilon > 0$  such that  $Ay > \epsilon \rho(A) \cdot z$  or equivalently,  $\frac{A}{(1+\epsilon)\rho(A)}z > z$ . Successively multiply both sides of  $\frac{A}{(1+\epsilon)\rho(A)}z > z$  by  $\frac{A}{(1+\epsilon)\rho(A)}$  and we have

$$\left(\frac{A}{(1+\epsilon)\rho(A)}\right)^k z > \dots > \frac{A}{(1+\epsilon)\rho(A)}z > z, \quad \text{for } k = 1, 2, \dots$$

Note that  $\lim_{k \rightarrow \infty} \left(\frac{A}{(1+\epsilon)\rho(A)}\right)^k \rightarrow \mathbf{0}$  because  $\rho\left(\frac{A}{(1+\epsilon)\rho(A)}\right) = \frac{\rho(A)}{(1+\epsilon)\rho(A)} < 1$ . Then, in the limit,  $z < \mathbf{0}$ . This conflicts the fact that  $z > \mathbf{0}$ . The assumption that  $y \neq \mathbf{0}$  is invalid

Thus we have  $y = \mathbf{0}$  which means  $\rho(A)$  is a positive eigenvalue of  $A$  and  $|x|$  is the corresponding eigenvector. Since  $\rho(A)|x| = A|x| > \mathbf{0}$ , we have  $|x| > \mathbf{0}$ . □

**Lemma 9** For  $A, B \in \mathbb{C}^{n \times n}$ , if  $|A| \leq B$ , then  $\rho(A) \leq \rho(B)$ .

*Proof.* By spectral radius formula, we have that for any sub-multiplicative norm  $\|\cdot\|$ ,  $\rho(A) = \lim_{k \rightarrow \infty} \|A^k\|^{\frac{1}{k}}$  and  $\rho(B) = \lim_{k \rightarrow \infty} \|B^k\|^{\frac{1}{k}}$ .

Note that since  $|A| \leq B$ , we have  $|A|^k \leq B^k$  for  $k \in \mathbb{N} \setminus \{0\}$ . Then  $\|A^k\|_\infty \leq \| |A|^k \|_\infty \leq \|B^k\|_\infty$  and sequentially  $\|A^k\|_\infty^{\frac{1}{k}} \leq \|B^k\|_\infty^{\frac{1}{k}}$ . Thus,  $\rho(A) \leq \rho(B)$ . □

**Theorem 10** (Theorem 4 restated). Each nonnegative matrix  $A$  has a nonnegative real eigenvalue with spectral radius  $\rho(A) = a$ , and  $a$  has a corresponding nonnegative eigenvector.

*Proof.* Construct a matrix sequence  $\{A_k\}_{k=1}^\infty$  by letting  $A_k = A + \frac{E}{k}$  where  $E$  is the matrix of all 1's. Let  $a_k = \rho(A_k) > 0$  and  $x_k > \mathbf{0}$  is the corresponding eigenvector.<sup>1</sup> Without loss of generality, let  $\|x_k\|_1 = 1$ . Since  $\{x_k\}_{k=1}^\infty$  is bounded, by **BolzanoWeierstrass theorem**, there exists a subsequence of  $\{x_k\}_{k=1}^\infty$  in  $\mathbb{R}^n$  that is convergent. Denote this convergent subsequence by  $\{x_{k_i}\}_{i=1}^\infty$  and  $\{x_{k_i}\}_{i=1}^\infty \rightarrow z$  where  $z \geq \mathbf{0}$  and  $z \neq \mathbf{0}$  (for each  $x_{k_i}$  satisfies that  $\|x_{k_i}\|_1 = 1$ ). Since  $\{A_k\}_{k=1}^\infty$  is monotone decreasing, by Lemma 9, we have that  $a_1 \geq \dots \geq a_k \geq a$ . Sequence  $\{a_k\}_{k=1}^\infty$  is nonincreasing and bounded, so  $\lim_{k \rightarrow \infty} a_k \rightarrow a^*$  exists and  $\lim_{i \rightarrow \infty} a_{k_i} \rightarrow a^* \geq a$ . Then we have

$$Az = \lim_{i \rightarrow \infty} A_{k_i} x_{k_i} = \lim_{i \rightarrow \infty} a_{k_i} x_{k_i} = a^* z.$$

Thus,  $a^*$  is an eigenvalue of  $A$  and  $a^* \leq a$ . Then we have  $a^* = a$ . So  $A$  has a nonnegative real eigenvalue  $a$  and  $z$  is the corresponding nonnegative eigenvector. □

<sup>1</sup> The existence of such  $x_k$  is guaranteed by Theorem 8.

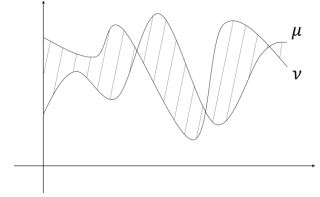
### 3 Coupling

#### 3.1 Total Variation Distance

**Definition 11** The total variation distance between two distributions  $\mu$  and  $\nu$  on a countable state space  $\Omega$  is given by

$$D_{TV}(\mu, \nu) = \frac{1}{2} \sum_{x \in \Omega} |\mu(x) - \nu(x)|.$$

We can look at the figure of two distributions on the sample space. The total variation distance is half the area enclosed by the two curves.



The total variation distance can be equivalently viewed in the following way.

**Lemma 12** We define  $\mu(A) = \sum_{x \in A} \mu(x)$ ,  $\nu(A) = \sum_{x \in A} \nu(x)$ , then we have

$$D_{TV}(\mu, \nu) = \max_{A \subseteq \Omega} |\mu(A) - \nu(A)|.$$

*Proof.* Let  $\Omega^+ \subseteq \Omega$  be the set of states such that  $\mu(x) \geq \nu(x)$ , and let  $\Omega^- \subseteq \Omega$  be the set of states such that  $\nu(x) > \mu(x)$ . It can be easily verified that

$$\max_{A \subseteq \Omega} \mu(A) - \nu(A) = \mu(\Omega^+) - \nu(\Omega^+),$$

$$\max_{A \subseteq \Omega} \nu(A) - \mu(A) = \nu(\Omega^-) - \mu(\Omega^-).$$

By  $\mu(\Omega) = \nu(\Omega) = 1$ ,

$$\mu(\Omega^+) + \mu(\Omega^-) = \nu(\Omega^+) + \nu(\Omega^-) = 1,$$

which implies that

$$\mu(\Omega^+) - \nu(\Omega^+) = \nu(\Omega^-) - \mu(\Omega^-).$$

We derive that

$$\max_{A \subseteq \Omega} |\nu(A) - \mu(A)| = \nu(\Omega^-) - \mu(\Omega^-) = \mu(\Omega^+) - \nu(\Omega^+).$$

Therefore,

$$\begin{aligned} D_{TV}(\mu, \nu) &= \sum_{x \in \Omega} \frac{1}{2} |\mu(x) - \nu(x)| \\ &= \frac{1}{2} (|\mu(\Omega^+) - \nu(\Omega^+)| + |\mu(\Omega^-) - \nu(\Omega^-)|) \\ &= \max_{A \subseteq \Omega} |\nu(A) - \mu(A)|. \end{aligned}$$

□

### 3.2 The Coupling Lemma

The coupling of two distributions is simply a joint distribution of them.

**Definition 13 (Coupling)** Let  $\mu$  and  $\nu$  be two distributions on the same space  $\Omega$ . Let  $\omega$  be a distribution on the space  $\Omega \times \Omega$ . If  $(X, Y) \sim \omega$  satisfies  $X \sim \mu$  and  $Y \sim \nu$ , then  $\omega$  is called a coupling of  $\mu$  and  $\nu$ .

In other words, the marginal probabilities of the disjoint distribution  $\omega$  are  $\mu$  and  $\nu$  respectively.

A special case is when  $X$  and  $Y$  are independent. However, in many applications, we want  $X$  and  $Y$  to be correlated while keeping their respect marginal probabilities correct.

prob \ y	HEAD	TAIL
x \ HEAD	1/3	1/6
x \ TAIL	0	1/2

prob \ y	HEAD	TAIL
x \ HEAD	1/6	1/3
x \ TAIL	1/6	1/3

We now give a toy example about how to construct different couplings on two fixed distributions. There are two coins: the first coin has probability  $\frac{1}{2}$  for head in a toss and  $\frac{1}{2}$  for tail, and the second coin has probability  $\frac{1}{3}$  and  $\frac{2}{3}$  respectively. We now construct two couplings as follows.

The table defines a joint distribution and the sum of a certain row/column equal to the corresponding marginal probability. It is clear that both table are couplings of the two coins. Among all the possible couplings, sometimes we are interested in the one who is "mostly coupled".

**Lemma 14 (Coupling Lemma)** *Let  $\mu$  and  $\nu$  be two distributions on a sample space  $\Omega$ . Then for any coupling  $\omega$  of  $\mu$  and  $\nu$  it holds that,*

$$\Pr_{(X,Y) \sim \omega} [X \neq Y] \geq D_{\text{TV}}(\mu, \nu).$$

And furthermore, there exists a coupling  $\omega^*$  of  $\mu$  and  $\nu$  such that

$$\Pr_{(X,Y) \sim \omega^*} [X \neq Y] = D_{\text{TV}}(\mu, \nu).$$

Let us prove the coupling lemma. For finite  $\Omega$ , designing a coupling is equivalent to filling a  $\Omega \times \Omega$  matrix in the way that the marginals are correct.

Clearly we have

$$\begin{aligned} \Pr [X = Y] &= \sum_{t \in \Omega} \Pr [X = Y = t] \\ &\leq \sum_{t \in \Omega} \min \{ \mu(t), \nu(t) \}. \end{aligned}$$

Thus,

$$\begin{aligned} \Pr [X \neq Y] &\geq 1 - \sum_{t \in \Omega} \min (\mu(t), \nu(t)) \\ &= \sum_{t \in \Omega} (\mu(t) - \min \{ \mu(t), \nu(t) \}) \\ &= \max_{A \subseteq \Omega} \{ \mu(A) - \nu(A) \} \\ &= D_{\text{TV}}(\mu, \nu). \end{aligned}$$

To construct  $\omega^*$  achieving the equality, for every  $t \in \Omega$ , we let

$$\Pr_{(X,Y) \sim \omega^*} [X = Y = t] = \min \{ \mu(t), \nu(t) \}.$$

I leave the construction of the off-diagonal entries of  $\omega^*$  as an exercise.



The coupling lemma provides a way to upper bound the distance between two distributions: For any two distributions  $\mu$  and  $\nu$  and any coupling  $\omega$  of  $\mu$  and  $\nu$ , an upper bound for  $\Pr_{(X,Y)\sim\omega} [X \neq Y]$  is an upper bound for  $D_{\text{TV}}(\mu, \nu)$ . This is a quite useful approach to bound the total variation distance since the convergence of a Markov chain can be implied by that of the total variance  $D_{\text{TV}}(\mu_t, \pi)$ . The coupling lemma also tells us that the upper bound obtained in this way can be tight, as long as you are able to find the optimal coupling. We will examine this in detail in the next lecture.

### *References*

[Mey00] Carl D Meyer. *Matrix analysis and applied linear algebra*, volume 71. SIAM, 2000. 5