

# [AI2613 Lecture 1] Review of Probability Theory

February 24, 2023

## 1 Probability Space

We start with the notion of probability space. A standard reference for the probability theory is [1].

**Definition 1** (Probability Space). A probability space is a tuple  $(\Omega, \mathcal{F}, P(\cdot))$  satisfying the following requirements.

- The universe  $\Omega$  is a set of “outcomes” (which can be either countable or uncountable).
- The set  $\mathcal{F} \subseteq 2^\Omega$  is a  $\sigma$ -algebra (the set of all possible “events”). Here we say  $\mathcal{F}$  is a  $\sigma$ -algebra if  $\mathcal{F}$  satisfies:
  - $\emptyset, \Omega \in \mathcal{F}$ ;
  - $\forall A \in \mathcal{F}$ , it holds  $A^c \in \mathcal{F}$ ;
  - for any finite or countable sequence of sets  $A_1, \dots, A_n, \dots \in \mathcal{F}$ , it holds that  $\bigcup_{i=1}^\infty A_i \in \mathcal{F}$ .
- The probability function  $P(\cdot) : \mathcal{F} \rightarrow [0, 1]$  satisfies
  - $P(\emptyset) = 0, P(\Omega) = 1$ ;
  - $P(A^c) = 1 - P(A)$  for all  $A \in \mathcal{F}$ ;
  - for any finite or countable sequence of disjoint sets  $A_1, \dots, A_n, \dots \in \mathcal{F}$ , it holds that  $P(\bigcup_{i=1}^\infty A_i) = \sum_{i=1}^\infty P(A_i)$ .

$$A^c := \Omega \setminus A.$$

Let  $\mathcal{S} \subseteq 2^\Omega$ . We use  $\sigma(\mathcal{S})$  to denote the minimal  $\sigma$ -algebra containing sets in  $\mathcal{S}$ . That is, for any  $\mathcal{F} \subseteq 2^\Omega$ ,  $\mathcal{F} = \sigma(\mathcal{S})$  if and only if (1)  $\mathcal{F}$  is a  $\sigma$ -algebra; (2)  $\mathcal{S} \subseteq \mathcal{F}$ ; (3) For any  $\mathcal{F}' \subseteq \mathcal{F}$  such that  $\mathcal{S} \subseteq \mathcal{F}'$ ,  $\mathcal{F}'$  is not a  $\sigma$ -algebra.

The term “minimal” here is with respect to the set inclusion relation  $\subseteq$ .

For every  $n \in \mathbb{N}$ , we use  $[n]$  to denote the set  $\{1, 2, \dots, n\}$ .

**Example 1** (Tossing  $n$  fair coins). Let  $\Omega = \{0, 1\}^n$ ,  $\mathcal{F} = 2^\Omega$  and for every  $S \in \{0, 1\}^n$ ,  $P(\{S\}) = \frac{1}{2^n}$ .

**Example 2** (Uniform Reals in  $(0, 1)$ ). The uniform distribution on  $(0, 1)$  is defined as follows:

- $\Omega = (0, 1)$ ;
- $\mathcal{F}$  is the  $\sigma$ -algebra consisting of all **Borel sets** on  $(0, 1)$ , namely the collection of subsets of  $(0, 1)$  obtained from open intervals by repeatedly taking countable unions and complements;
- $\forall$  interval  $I = (a, b)$ ,  $P(I) = b - a$  (This is the **Lebesgue measure**).

The definition here, although a bit wired at the first glance, is in fact the simplest way to capture our intuition that the probability that a point is in  $(a, b)$  should be  $b - a$ . We cannot take  $\mathcal{F} = 2^\Omega$  in Example 2 as doing so may include some *non-measurable* sets. In fact,  $\mathcal{F}$  is called the **Borel algebra**, which is the smallest  $\sigma$ -algebra containing all open intervals. One can construct a non-Borel set in  $(0, 1)$  assuming the *axiom of choice*. In fact, the existence of a non-Borel set is independent of **Zermelo-Fraenkel set theory** without the axiom of choice. We use  $\mathcal{B}$  to denote the collection of Borel sets on  $\mathbb{R}$ . For any  $A \subseteq \mathbb{R}$ , we use  $\mathcal{B}(A)$  to denote  $\mathcal{B} \cap 2^A$ .

## 2 Random Variables

**Definition 2** (Measurable Space). Consider a set  $\Omega$  and a  $\sigma$ -algebra  $\mathcal{F}$  on  $\Omega$ . The tuple  $(\Omega, \mathcal{F})$  is called a measurable space.

**Definition 3** (Measurable Function). Let  $(\Omega, \mathcal{F})$  and  $(\Omega', \mathcal{F}')$  be two measurable spaces and  $X : \Omega \rightarrow \Omega'$  be a function. We say  $X$  is a  $\mathcal{F}$ -measurable function if

$$\forall B' \in \mathcal{F}', X^{-1}(B') \in \mathcal{F},$$

For any function, we use  $\sigma(X)$  to denote the minimal  $\sigma$ -algebra  $\mathcal{F}$  such that  $X$  is  $\mathcal{F}$ -measurable.

**Definition 4** (Random Variable). Let  $\Omega'$  and  $\mathcal{F}'$  in Definition 3 be  $\mathbb{R}$  and the Borel algebra  $\mathcal{B}$ , then  $X$  in Definition 3 is a (real-valued) random variable.

We say a random variable  $X$  discrete if its range  $\text{Ran}(X)$  is countable. In other words,  $X$  can only take at most countable many distinct values. Otherwise, we say  $X$  is a continuous random variable.

**Example 3** (Measurable Functions of Tossing a Dice). Let  $\Omega = [6]$ . We have three  $\sigma$ -algebras on  $\Omega$ :  $\mathcal{F}_1 = 2^{[6]}$ ,  $\mathcal{F}_2 = \sigma(\{1, 3, 5\})$  and  $\mathcal{F}_3 = \sigma(\{1, 2\})$ . Consider three random variables  $X_1, X_2, X_3 : \Omega \rightarrow \mathbb{R}$  such that  $X_1 : \omega \mapsto \omega$ ,  $X_2 : \omega \mapsto \omega \bmod 2$  and  $X_3 : \omega \mapsto \mathbf{1}[\omega \leq 2]$ . Then all these three mappings are  $\mathcal{F}_1$ -measurable, only  $X_2$  is  $\mathcal{F}_2$ -measurable and only  $X_3$  is  $\mathcal{F}_3$ -measurable.

$X^{-1}(B') \triangleq \{\omega \in \Omega | X(\omega) \in B'\}$  is the inverse of  $X$ .

The measurability of a random variable  $X$  captures the intuition that we can safely talk about the probability of  $X$  taking some value. Intuitively  $X$  induces a partition of  $\Omega$  where two outcomes  $\omega_1$  and  $\omega_2$  are in the same partition if and only if  $X(\omega_1) = X(\omega_2)$ . If the partition defined by  $X$  is more “coarser” than the partition defined by a  $\sigma$ -algebra  $\mathcal{F}$ , then  $X$  is  $\mathcal{F}$  measurable.

## 3 Distribution

Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : \Omega \rightarrow \mathbb{R}$  be a  $\mathcal{F}$ -measurable random variable. Let  $\mathcal{B}$  be the Borel algebra on  $\mathbb{R}$ . The distribution space  $(\mathbb{R}, \mathcal{B}, \mathbf{Pr})$  induced by  $X$  is defined as

$$\forall A \in \mathcal{B}, \mathbf{Pr}[A] = \mathbf{Pr}[X \in A] \triangleq \mathbf{P}[X^{-1}(A)].$$

The function  $F(x) := \mathbf{Pr}[X \leq x] = \mathbf{P}(X^{-1}(-\infty, x])$  is called the cumulative distribution function (cdf) of  $X$ .

If a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  satisfies for any  $a \leq b$ :

$$\int_a^b f(x) dx = F(b) - F(a),$$

then we call  $f(x)$  a probability density function (pdf) of  $X$ .

**Example 4** (Exponential Distribution). If  $X \sim \text{Exp}(\lambda)$ , or equivalently it follows exponential distribution with rate  $\lambda$  for  $\lambda > 0$ , then its pdf is

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x > 0 \\ 0, & \text{otherwise.} \end{cases}$$

## 4 Expectation and Variance

**Definition 5** (Expectation). . Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $X : \Omega \rightarrow \mathbb{R}$  be a random variable.

- For a discrete random variable  $X$ , its expectation is

$$\mathbf{E}[X] := \sum_{a \in \text{Ran}(X)} a \cdot \mathbf{Pr}[X = a].$$

If  $\Omega$  is at most countable, we can also write

$$\mathbf{E}[X] = \sum_{\omega \in \Omega} \mathbf{P}(\{\omega\}) \cdot X(\omega).$$

- For a continuous random variable  $X$  with pdf  $f$ , its expectation is

$$\mathbf{E}[X] := \int_{-\infty}^{\infty} t \cdot f(t) dt.$$

Sometimes it is more convenient to equivalently write the expectation as

$$\mathbf{E}[X] = \int_{\Omega} X(\omega) \mu(d\omega) = \int_{\Omega} X d\mu.$$

using *Lebesgue integration*.

**Example 5** (Expectation of Exponential Distribution). Let  $X \sim \text{Exp}(\lambda)$  for  $\lambda > 0$ , then

$$\mathbf{E}[X] = \int_0^{\infty} t \cdot \lambda e^{-\lambda t} dt = \frac{1}{\lambda}.$$

**Definition 6** (Variance). The variance of a random variable  $X$  is

$$\mathbf{Var}[X] := \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2] - \mathbf{E}[X]^2.$$

**Proposition 7.** Let  $X_1, \dots, X_n$  be random variables where  $n$  is a finite constant. Then

$$\mathbf{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbf{E}[X_i].$$

## 5 Conditional Probability

**Definition 8** (Conditional Probability). Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Let  $A, B \in \mathcal{F}$  be two events with  $\mathbf{P}(B) > 0$ . The conditional probability of  $A$  given  $B$  is

$$\mathbf{P}(A|B) := \frac{\mathbf{P}(A \cap B)}{\mathbf{P}(B)}.$$

In the following, we define the notion of *conditional expectation* for those *discrete* random variables.

This is well-defined since we know from the definition of  $\sigma$ -algebra that  $A \cap B \in \mathcal{F}$ .

**Definition 9** (Conditional Expectation). Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space. Let  $A \in \mathcal{F}$  be an event with  $\mathbf{P}(A) > 0$ . Let  $X : \Omega \rightarrow \mathbb{R}$  be a discrete random variable. The conditional expectation of  $X$  conditioned on  $A$  is

$$\mathbf{E}[X | A] := \sum_{a \in \text{Ran}(X)} a \cdot \mathbf{Pr}[X = a | A].$$

Let  $Y : \Omega \rightarrow \mathbb{R}$  be another discrete random variable. The conditional expectation of  $X$  conditioned on  $Y$ , written as  $\mathbf{E}[X | Y]$ , is a random variable  $f_Y : \Omega \rightarrow \mathbb{R}$  such that

$$\forall \omega \in \Omega : f_Y(\omega) = \mathbf{E}[X | Y^{-1}(Y(\omega))] = \mathbf{E}[X | Y = Y(\omega)]. \quad (1)$$

**Proposition 10.**

- $\mathbf{E}[X | Y]$  is  $\sigma(Y)$ -measurable.
- $\mathbf{E}[\mathbf{E}[X | Y]] = \mathbf{E}[f_Y] = \mathbf{E}[X]$ .

*Proof.* • Since the value of  $\mathbf{E}[X | Y]$  is determined by  $Y(\omega)$ , it is clearly  $\sigma(Y)$ -measurable.

- We compute  $\mathbf{E}[f_Y]$  by definition.

$$\begin{aligned} \mathbf{E}[f_Y] &= \sum_{y \in \text{Ran}(Y)} \mathbf{Pr}[Y = y] \cdot \mathbf{E}[X | Y = y] \\ &= \sum_{y \in \text{Ran}(Y)} \mathbf{Pr}[Y = y] \cdot \sum_{x \in \text{Ran}(X)} \mathbf{Pr}[X = x | Y = y] \cdot x \\ &= \sum_{x \in \text{Ran}(X)} x \cdot \sum_{y \in \text{Ran}(Y)} \mathbf{Pr}[Y = y] \cdot \mathbf{Pr}[X = x | Y = y] \\ &= \sum_{x \in \text{Ran}(X)} x \cdot \sum_{y \in \text{Ran}(Y)} \mathbf{Pr}[X = x \wedge Y = y] \\ &= \sum_{x \in \text{Ran}(X)} x \cdot \mathbf{Pr}[X = x] \\ &= \mathbf{E}[X]. \end{aligned}$$

□

## 6 Conditional Expectation for General Random Variables

The definition of conditional expectation for continuous random variables is more subtle. For example, if  $X, Y \sim N(0, 1)$  are two independent random variables following standard normal distribution, then intuitively  $\mathbf{E}[X | Y = 0]$  should be identical to  $\mathbf{E}[X]$ , which is zero. However, we cannot directly adopt the definition before since  $\mathbf{Pr}[Y = 0] = 0$ .

**Definition 11.** Let  $(\Omega, \mathcal{F}, P)$  be the probability space. Let  $X$  be a random variable with  $E[|X|] < \infty$ . The conditional expectation  $E[X | Y]$  is a  $\sigma(Y)$ -measurable random variable  $f_Y$  satisfying

$$\forall A \in \sigma(Y), \int_A f_Y dP = \int_A X dP.$$

The existence and uniqueness of  $f_Y$  follow from **Radon-Nikodym theorem**.

## 7 Balls-into-Bins

Balls-into-bins is a simple random process in which a person throws  $m$  balls into  $n$  bins uniformly at random. Many interesting questions can be asked about the process.

### 7.1 Birthday Paradox

*Birthday paradox* refers to the seemingly counter-intuitive fact that some students in the class are very likely to share the same birthday. Viewing bins as dates and balls as students, the event that two students have the same birthday can be modeled as the event that some bin contains more than one ball.

Note that each ball is thrown independently. Condition on there is no collision after the  $k-1$  balls are thrown, the probability that no collision occurs after throwing the  $k^{\text{th}}$  ball is  $\frac{n-k+1}{n}$ . Hence,

$$\begin{aligned} \Pr[\text{no same birthday}] &= \prod_{k=1}^m \frac{n-k+1}{n} \\ &= \prod_{k=1}^{m-1} \left(1 - \frac{k}{n}\right) \\ &\leq \exp\left\{-\frac{\sum_{k=1}^{m-1} k}{n}\right\} \quad (\text{by } 1+x \leq e^x) \\ &= \exp\left\{-\frac{m(m-1)}{2n}\right\}. \end{aligned} \tag{2}$$

For  $m = O(\sqrt{n})$ , the probability can be arbitrarily close to 0.

### 7.2 Coupon Collector

The coupon collector problem asks the following question: If each box of a brand of cereals contains a coupon, randomly chosen from  $n$  different types of coupons, what is the number of boxes one needs to buy to collect all  $n$  coupons? In the language of balls-into-bins, it asks how many balls one needs to throw until each of the  $n$  bins contains at least one ball.

When  $n$  is sufficiently large, Equation (2) is tight because  $\frac{k}{n} \leq \frac{m}{n} = O(\frac{1}{\sqrt{n}}) \rightarrow 0$  and  $1+x \leq e^x$  is tight when  $x$  is small.

The expectation can be easily calculated using the linearity of expectations. Let  $X_i$  be the number of balls to throw to get the  $i$ -th distinct type of coupon while exactly  $i - 1$  distinct types of coupons are already in hand. Then the number of draws  $X$  to collect all coupons satisfies

$$X = \sum_{i=1}^{n-1} X_i.$$

By the linearity of expectations:

$$\mathbf{E}[X] = \sum_{i=1}^n \mathbf{E}[X_i].$$

It is clear that  $X_i \sim \text{Geom}(\frac{n-i+1}{n})$  and therefore  $\mathbf{E}[X_i] = \frac{n}{n-i+1}$ . As a result,

$$\mathbf{E}[X] = \sum_{i=1}^n \frac{n}{n-i+1} = n \cdot H(n),$$

where  $H(n)$  is the harmonic number satisfying  $\lim_{n \rightarrow \infty} H(n) = \log n + \gamma$  for  $\gamma = 0.577 \dots$

$\gamma$  is called the **Euler constant**.

## 8 Concentration Inequalities

In addition to the expectation, we are often interested in how a random variable deviates from certain fixed value. Concentration inequalities are inequalities of this form.

### 8.1 Markov's Inequality

**Theorem 12** (Markov's Inequality). . For any non-negative random variable  $X$  and  $a > 0$ ,

$$\Pr[X \geq a] \leq \frac{\mathbf{E}[X]}{a}.$$

*Proof.* Since  $X$  is non-negative, we have

$$\mathbf{E}[X] \geq a \cdot \Pr[X \geq a] + 0 \cdot \Pr[X < a].$$

This is equivalent to

$$\Pr[X \geq a] \leq \frac{\mathbf{E}[X]}{a}.$$

□

**Example 6** (Concentration for Coupon Collector). . Recall that  $X$  is the number of balls we need. Apply Markov's inequality, for  $c > 0$  we have

$$\Pr[X \geq c] \leq \frac{\mathbf{E}[X]}{c} = \frac{nH_n}{c}.$$

Thus, the probability that we need to draw the coupon for more than  $100 \cdot nH_n$  times is less than 0.01.

## 8.2 Chebyshev's Inequality

A common trick to improve concentration is to consider  $\mathbf{E}[f(X)]$  instead  $\mathbf{E}[X]$  for some increasing function  $f : \mathbb{R} \rightarrow \mathbb{R}$  since

$$\Pr[X \geq a] = \Pr[f(X) \geq f(a)].$$

Concentration inequalities give a sense that how the random variable deviate from its expectation. Then the probability we care about is actually  $\Pr[|X - \mathbf{E}[X]| \geq a]$  for some positive constant  $a$ . Choosing the increasing function  $f(x) = x^2$ , we get the following Chebyshev's inequality.

**Theorem 13** (Chebyshev's Inequality). *For any random variable with bounded  $\mathbf{E}[X]$  and  $a \geq 0$ , it holds that*

$$\Pr[|X - \mathbf{E}[X]| \geq a] \leq \frac{\mathbf{Var}[X]}{a^2}$$

*Proof.* Let  $Y = |X - \mathbf{E}[X]|$ , then clearly  $Y \geq 0$ . Therefore

$$\begin{aligned} \Pr[|X - \mathbf{E}[X]| \geq a] &= \Pr[Y \geq a] = \Pr[Y^2 \geq a^2] \leq \frac{\mathbf{E}[Y^2]}{a^2} \\ &= \frac{\mathbf{E}[(X - \mathbf{E}[X])^2]}{a^2} = \frac{\mathbf{Var}[X]}{a^2}. \end{aligned}$$

□

**Example 7** (Coupon Collector Revisited). *We apply Chebyshev's inequality to the coupon collector problem. Assuming the notation before, we have*

$$\Pr[X \geq nH_n + t] \leq \Pr[|X - \mathbf{E}[X]| \geq t] \leq \frac{\mathbf{Var}[X]}{t^2}.$$

Recall that the variable  $X_i$  indicates the number of draws to get a new coupon while there are  $i$  coupons in hands. For distinct  $i$  and  $j$ ,  $X_i$  and  $X_j$  are independent. Then

$$\mathbf{Var}[X] = \mathbf{Var}\left[\sum_{i=0}^{n-1} X_i\right] = \sum_{i=0}^{n-1} \mathbf{Var}[X_i].$$

For  $i \in \{0, 1, \dots, n-1\}$ ,  $X_i \sim \text{Geom}\left(\frac{n-i}{n}\right)$ , so we have

$$\mathbf{Var}[X_i] = \frac{1 - \frac{n-i}{n}}{\left(\frac{n-i}{n}\right)^2} = \frac{i \cdot n}{(n-i)^2} \leq \frac{n^2}{(n-i)^2}.$$

It remains to bound  $\sum_{i=0}^{n-1} \frac{1}{(n-i)^2} = \sum_{i=1}^n \frac{1}{i^2}$ . Note that

$$\sum_{i=1}^n \frac{1}{i^2} \leq 1 + \int_1^\infty \frac{dx}{x^2} = 2.$$

Therefore, we have  $\mathbf{Var}[X] \leq 2n^2$  and  $\Pr[X \geq nH_n + t] \leq \frac{2n^2}{t^2}$ . The probability that we need to draw the coupon for more than  $\sqrt{200n} + nH_n$  times is less than 0.01.

The bound obtained by Chebyshev's inequality is much tighter than that via Markov's inequality where in order to obtain the same confidence, one needs to choose  $t = \Theta(n \log n)$ .

# [AI2613 Lecture 2] Discrete Markov Chains, Coupling

March 15, 2023

## 1 Discrete Markov Chain

### 1.1 Markov Chain

**Definition 1** (Discrete Markov Chain). Suppose there is a sequence of random variables

$$X_0, X_1, \dots, X_t, X_{t+1}, \dots$$

where the  $\text{Ran}(X_t) \subseteq \Omega$  for some countable  $\Omega$ . Then we call  $\{X_t\}$  a discrete Markov chain if  $\forall t \geq 1$  the distribution of  $X_t$  is only related to  $X_{t-1}$ , that is  $\forall a_0, a_1, \dots, a_t \in \Omega$ ,

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

**Example 1** (Random Walk on  $\mathbb{Z}$ ). 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. If we use  $X_t$  to denote his position at time  $t$ , then we have  $X_0 = 0$  and for every  $t > 0$ ,  $X_t = X_{t-1} + 1$  with probability 50% and  $X_t = X_{t-1} - 1$  with probability 50%. This is a simple Markov chain, since the position at time  $t$  only depends on the position at time  $t - 1$ .

In this lecture, we consider the situation that the state space  $\Omega = [n]$  is finite. Then a (time-homogeneous) Markov chain can be characterized by a  $n \times n$  matrix  $P = (p_{ij})_{i,j \in [n]}$  where  $p_{ij} = \Pr[X_{t+1} = j | X_t = i]$  for all  $t \geq 0$ .

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  $i$ .

**Example 2** (Finite State Random Walk). The following three vertex directed graph corresponds to the Markov chain with transition matrix

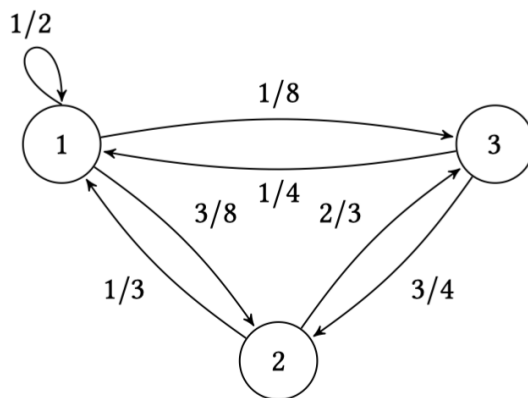
$$P = (p_{ij}) = \begin{bmatrix} 1/2 & 3/8 & 1/8 \\ 1/3 & 0 & 2/3 \\ 1/4 & 3/4 & 0 \end{bmatrix}. \text{ We sometimes call the graph the transition graph of } P.$$

At any time  $t \geq 0$ , we use  $\mu_t$  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  $\mu_t^\top = \mu_0^\top P^t$ . 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.

## 1.2 Stationary Distribution

**Definition 2** (Stationary Distribution). . A distribution  $\pi$  is a stationary distribution of  $P$  if it remains unchanged in the Markov chain as time progresses, i.e.,

$$\pi^T P = \pi^T.$$

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.

**Example 3** (Card Shuffling). Consider a naive “top-to-random” card shuffle: Suppose we have  $n$  cards, every time we take the top card of the deck and insert it into the deck at one of the  $n$  distinct possible places uniformly at random. Thus, there are  $n!$  possible permutations and  $p_{ij} > 0$  only if the  $i^{\text{th}}$  permutation can come to the  $j^{\text{th}}$  through one step “top-to-random” shuffle.

Performing the shuffle repeatedly is a Markov chain. It is not difficult to verify that the uniform distribution  $\left(\frac{1}{n!}, \frac{1}{n!}, \dots, \frac{1}{n!}\right)^T$  over all  $n!$  permutations is a stationary distribution.

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?

## 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  $\pi^T P = \pi^T$ . Observe that this is equivalent to “1 is an eigenvalue of  $P^T$  with a nonnegative eigenvector ( $P^T \pi = \pi$ )”.

We use the following lemma and theorem in linear algebra.

**Lemma 3.** *Every eigenvalue of nonnegative matrix  $P$  is no larger than the maximum row sum of  $P$ .*

*Proof.* Let  $\lambda$  be an eigenvalue of  $P$  and  $x$  is the corresponding eigenvector. We have

$$\|\lambda x\|_\infty = \|Px\|_\infty \leq \|P\|_\infty \cdot \|x\|_\infty.$$

Note that  $\|\lambda x\|_\infty = |\lambda| \|x\|_\infty$  and  $\|x\|_\infty > 0$ . Thus, we have  $\lambda \leq |\lambda| \leq \|P\|_\infty$ , that is  $\lambda$  is no larger than the maximum row sum of nonnegative matrix  $P$ .  $\square$

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

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,  $P^T$  shares the same characteristic polynomial with  $P$ , which implies the eigenvalues of  $P^T$

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

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

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

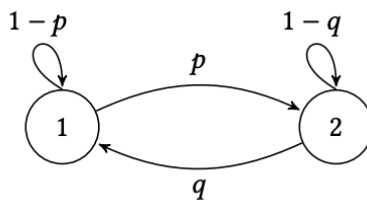
which is equivalent to

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

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

## 2.2 Uniqueness and Convergence

Consider the following Markov chain with two states. 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.,

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

Now we define

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

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 graph is disconnected, so it can be parti-



tioned 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.

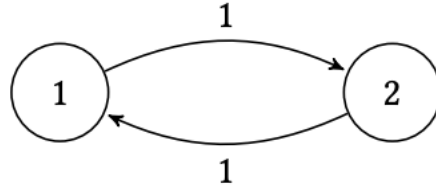
This observation motivates us to define the following:

**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 call 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  $\mu \in \mathbb{R}^n$ ,

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

### 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  $\rho(A)|x| = A|x| > 0$ , we have  $|x| > 0$ .  $\square$

**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,  $\rho(A) \leq \rho(B)$ .  $\square$

**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  $A_k = A + \frac{\mathbf{E}}{k}$  where  $\mathbf{E}$  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 **Bolzano–Weierstrass 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  $z \geq 0$  and  $z \neq 0$  (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.  $\square$

### 3 Coupling

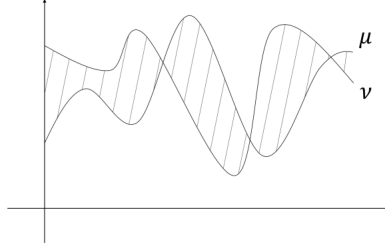
To measure how close the two distributions are, we need to define a distance between them.

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

**Definition 11** (Total Variation Distance). . *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 following figure of two distributions on the sample space. The total variation distance is half the area enclosed by the two curves.



This figure gives us the intuition of the following proposition which states that the total variation distance can be equivalently viewed in another way.

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

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

Our main tool to bound the distance between two distributions is the *coupling*. This is a useful technique in analysis of probabilities. A 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  $\nu$ .*

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”.

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 independently. 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
TAIL	0	1/2

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

**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,

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

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

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

*Proof.* 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,

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

To construct  $\omega^*$  achieving the equality, for every  $t \in \Omega$ , we let  $\Pr_{(X,Y) \sim \omega^*} [X = Y = t] = \min\{\mu(t), \nu(t)\}$ .  $\square$

## References

[Mey00] Carl D Meyer. *Matrix analysis and applied linear algebra*, volume 71. SIAM, 2000. 6

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  $D_{TV}(\mu, \nu)$ . This is a quite useful approach to bound the total variation distance.



# [AI2613 Lecture 3] Proof of FTMC, Mixing Time

March 15, 2023

## 1 Fundamental Theorem of Markov Chains

Recall the fundamental theorem of Markov chains for *finite* chains we introduced in the last lecture.

**Theorem 1** (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  $\mu \in \mathbb{R}^n$ ,*

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

Today we give a proof of the theorem. To this end, we first study the properties of the transition matrix  $P$  of an irreducible and aperiodic chain. Then we introduce the notion of *coupling*, a powerful technique to analyze stochastic processes.

**Claim 2.** *Let  $P \in \mathbb{R}^{n \times n}$  be an irreducible and aperiodic Markov chain. It holds that*

$$\exists t^* : \forall i, j \in [n] : P^{t^*}(i, j) > 0.$$

We use Lemma 3 to prove Claim 2.

**Lemma 3.** *Let  $c_1, c_2, \dots, c_s$  be a group of positive integers satisfying  $\gcd(c_1, \dots, c_s) = 1$ . For any sufficiently large integer  $b$ , there exists  $y_1, y_2, \dots, y_s \in \mathbb{N}$  such that*

$$c_1 y_1 + c_2 y_2 + \dots + c_s y_s = b.$$

That is, there exists some  $b_0 > 0$  such that for any  $b > b_0$ , the diophantine equation  $c_1 y_1 + c_2 y_2 + \dots + c_s y_s = b$  always has non-negative solutions

*Proof.* By **Bézout's identity** there exists  $x_1, x_2, \dots, x_s \in \mathbb{Z}$  such that

$$c_1 x_1 + c_2 x_2 + \dots + c_s x_s = 1.$$

We apply induction on  $s$ . The case  $s = 1$  trivially holds. Assume  $s \geq 2$  and the lemma holds for smaller  $s$ . Let  $g = \gcd(c_1, \dots, c_{s-1})$ . By induction hypothesis, we know that

$$\frac{a_1}{g} \cdot x_1 + \frac{a_2}{g} \cdot x_2 + \dots + \frac{a_{s-1}}{g} \cdot x_{s-1} = b' \iff a_1 \cdot x_1 + a_2 \cdot x_2 + \dots + a_{s-1} x_{s-1} = g \cdot b'$$

has non-negative solutions for sufficiently large  $b'$ . Therefore, we only need to prove that the equation

$$g \cdot b' + a_s \cdot x_s = b \tag{1}$$

has nonnegative solution  $(b', x_s)$  with sufficiently large  $b'$  when  $b$  is sufficiently large. In other words, we need to prove for any  $b_0 > 0$ , eq. (1) has nonnegative solution with  $b' > b_0$  for any sufficiently large  $b$ .

Note that  $\gcd(g, a_s) = 1$ , we can find integers  $(y, x)$  such that

$$g \cdot y + a_s \cdot x = 1 \iff g \cdot (by) + a_s \cdot (bx) = b.$$

Noting that for any  $k \in \mathbb{Z}_{\geq 0}$ , we have  $g \cdot (by + ka_s) + a_s \cdot (bx - kg) = b$ . We need  $by + ka_s > b_0$  and  $bx - kg \geq 0$ , which are equivalent to

$$\frac{bx}{g} \geq k > \frac{b_0 - by}{a_s}.$$

We can always find such an integer  $k$  if  $b \geq g(b_0 + a_s)$ . □

*Proof of Claim 2.* The property of irreducibility implies that

$$\forall i, j : \exists t : P^t(i, j) > 0.$$

Suppose that there are  $s$  loops of length  $c_1, c_2, \dots, c_s$  starting from and ending at state  $i$ . Then by aperiodicity we have

$$\gcd(c_1, c_2, \dots, c_s) = 1.$$

For any sufficiently large  $m$  and any pair of states  $(i, j)$ , by Lemma 3 and irreducibility, there exists a path from  $i$  to  $j$  with exactly  $m$  steps. Thus, there exist  $t^* > 0$  such that for any state pair  $(i, j)$ ,  $P^{t^*}(i, j) > 0$ . Furthermore, for any  $t > t^*$ ,  $P^t(i, j) > 0$  for any  $i, j \in \Omega$ . □

### 1.1 Proof of FTMC

*Proof.* We already know that  $P$  has a stationary distribution  $\pi$ . What we would like to show is that for all starting distribution  $\mu_0$ , it holds that

$$\lim_{t \rightarrow \infty} D_{TV}(\mu_t, \pi) = 0,$$

where  $\mu_t^\top = \mu_0^\top P^t$ .

Suppose that  $\{X_t\}$  and  $\{Y_t\}$  are two identical Markov chains starting from different distribution, where  $Y_0 \sim \pi$  while  $X_0$  is generated from an arbitrary distribution  $\mu_0$ .

Now we have two sequence of random variables:

$$\begin{array}{ccccccccccc} \mu_0 & & \mu_1 & & & & \mu_t & & & & \\ \wr & & \wr & & & & \wr & & & & \\ X_0 & \rightarrow & X_1 & \rightarrow & X_2 & \rightarrow & \dots & \rightarrow & X_t & \rightarrow & X_{t+1} & \rightarrow & \dots \\ & & & & & & & & & & & & \\ Y_0 & \rightarrow & Y_1 & \rightarrow & Y_2 & \rightarrow & \dots & \rightarrow & Y_t & \rightarrow & Y_{t+1} & \rightarrow & \dots \\ \wr & & \wr & & & & \wr & & & & & & \\ \pi & & \pi & & & & \pi & & & & & & \end{array}$$

The coupling lemma establishes the connection between the distance of distributions and the discrepancy of random variables. To show that  $D_{TV}(\mu_t, \pi) \rightarrow 0$ , it is sufficient to construct a coupling  $\omega_t$  of  $\mu_t$  and  $\pi$  and then compute  $\Pr_{(X_t, Y_t) \sim \omega_t} [X_t \neq Y_t]$ .

Here we give a simple coupling. Let  $(X_t, Y_t) \sim \omega_t$  and we construct  $\omega_{t+1}$ . If  $X_t = Y_t$  for some  $t \geq 0$ , then let  $X_{t'} = Y_{t'}$  for all  $t' > t$ , otherwise  $X_{t+1}$  and  $Y_{t+1}$  are independent. Namely,  $\{X_t\}$  and  $\{Y_t\}$  are two independent Markov chains until  $X_t$  and  $Y_t$  reach the same state for some  $t \geq 0$ , and once they meet together then they move together forever.

The coupling lemma tells us that  $D_{TV}(\mu_t, \pi) \leq \Pr_{(X_t, Y_t) \sim \omega_t} [X_t \neq Y_t]$ .

Let  $t^*$  be the same  $t^*$  with Claim 2. Let  $\alpha$  be a positive constant such that  $P^{t^*}(i, j) \geq \alpha > 0$  for any state pair  $(i, j)$ . Define event  $B$  as  $\{\exists t < t^*, X_t = Y_t\}$ . We have that

$$\Pr[X_{t^*} = Y_{t^*}] = \Pr[X_{t^*} = Y_{t^*} \wedge B] + \Pr[X_{t^*} = Y_{t^*} \wedge \bar{B}] \quad (2)$$

Suppose  $\{X'_t\}$  and  $\{Y'_t\}$  are two independent Markov chains with transition matrix  $P$  and  $X'_0 \sim \mu_0$  and  $Y'_0 \sim \pi$ . The only difference between  $(\{X'_t\}, \{Y'_t\})$  and  $(\{X_t\}, \{Y_t\})$  is that  $\{X'_t\}$  and  $\{Y'_t\}$  are independent all the time. Then

$$\begin{aligned} \Pr[X_{t^*} = Y_{t^*} = 1 \wedge \bar{B}] &= \Pr[X'_{t^*} = Y'_{t^*} = 1 \wedge \bar{B}] \\ &= \Pr[X'_{t^*} = 1] \cdot \Pr[Y'_{t^*} = 1] \\ &\quad - \sum_{t=0}^{t^*-1} \sum_{z \in [n]} \Pr[X'_t = Y'_t = z \wedge \forall s < t, X'_s \neq Y'_s] \cdot \Pr[X'_{t^*} = 1 \mid X'_t = z] \cdot \Pr[Y'_{t^*} = 1 \mid Y'_t = z]. \end{aligned}$$

Note that

$$\begin{aligned} \Pr[X_{t^*} = Y_{t^*} \wedge B] &\geq \Pr[X_{t^*} = Y_{t^*} = 1 \wedge B] \\ &= \sum_{t=0}^{t^*-1} \sum_{z \in [n]} \Pr[X_t = Y_t = z \wedge \forall s < t, X_s \neq Y_s] \cdot \Pr[X_{t^*} = 1 \mid X_t = z] \\ &= \sum_{t=0}^{t^*-1} \sum_{z \in [n]} \Pr[X'_t = Y'_t = z \wedge \forall s < t, X'_s \neq Y'_s] \cdot \Pr[X'_{t^*} = 1 \mid X'_t = z]. \end{aligned}$$

Thus, Equation (2)  $\geq \Pr[X'_{t^*} = 1] \cdot \Pr[Y'_{t^*} = 1] \geq \alpha^2$ .

By the coupling and the Markov property, we have

$$\begin{aligned} \Pr[X_{2t^*} \neq Y_{2t^*}] &= \Pr[X_{2t^*} \neq Y_{2t^*} \mid X_{t^*} = Y_{t^*}] \Pr[X_{t^*} = Y_{t^*}] \\ &\quad + \Pr[X_{2t^*} \neq Y_{2t^*} \mid X_{t^*} \neq Y_{t^*}] \Pr[X_{t^*} \neq Y_{t^*}] \\ &\leq \Pr[X_{2t^*} \neq Y_{2t^*} \mid X_{t^*} \neq Y_{t^*}] \Pr[X_{t^*} \neq Y_{t^*}] \\ &\leq (1 - \alpha^2)^2. \end{aligned}$$

Then we have  $\Pr[X_{kt^*} \neq Y_{kt^*}] \leq (1 - \alpha^2)^k$  by recursion. It yields that

$$\Pr[X_t \neq Y_t] = \sum_{x_0, y_0 \in [n]} \mu_0(x_0) \cdot \pi(y_0) \cdot \Pr[X_t \neq Y_t \mid X_0 = x_0, Y_0 = y_0] \rightarrow 0$$

as  $t \rightarrow \infty$ . □

## 2 Mixing Time

We are ready to study the convergence rate of Markov chains. We start with the notion of mixing time. For any  $\varepsilon > 0$ , the mixing time of a Markov chain  $P$  up to error  $\varepsilon$  is the minimum step  $t$  such that if we run the Markov chain from any initial distribution, its total variation distance to the stationary distribution is at most  $\varepsilon$ . Formally,

$$\tau_{\text{mix}}(\varepsilon) := \min_t \max_{\mu_0} D_{\text{TV}}(\mu_t, \pi) \leq \varepsilon.$$

Recalling in our proof of FTMC using the coupling argument, we obtain the following inequality

$$D_{\text{TV}}(\mu_t, \pi) \leq \Pr_{(X_t, Y_t) \sim \omega_t} [X_t \neq Y_t].$$

Therefore, if we can construct a coupling  $\omega_t$  such that for two arbitrary initial distributions,  $\Pr_{(X_t, Y_t) \sim \omega_t} [X_t \neq Y_t] \leq \varepsilon$ , then  $\tau_{\text{mix}}(\varepsilon) \leq t$ .

**Example 1** (Random walk on hypercube). . Consider the random walk on the  $n$ -cube. The state space  $\Omega = \{0, 1\}^n$ , and there is an edge between two state  $x$  and  $y$  iff  $\|x - y\|_1 = 1$ . We start from a point  $X_0 \in \Omega$ . In each step,

- With probability  $\frac{1}{2}$  do nothing.
- Otherwise, pick  $i \in [n]$  uniformly at random and flip  $X(i)$ .

It's equivalent to the following process:

- Pick  $i \in [n], b \in \{0, 1\}$  uniformly at random.
- Change  $X(i)$  to  $b$ .

Now we analyze the mixing time of the process using coupling. We apply the following simple coupling rule:

- We couple two walks  $X_t$  and  $Y_t$  by choosing the same  $i, b$  in every step.

Once a position  $i \in [n]$  has been picked,  $X_t(i)$  and  $Y_t(i)$  will be the same forever. Therefore, the problem again reduces to the coupon collector problem.

For  $t \geq n \log n + cn$ , the probability that the  $i^{\text{th}}$  dimension is not chosen is

$$\left(1 - \frac{1}{n}\right)^{n \log n + cn} \leq \frac{e^{-c}}{n}.$$

Then the probability that there exists at least one dimension which is not chosen is no larger than  $e^{-c}$ . We want this value to be less than  $\varepsilon$ . Then we choose  $c > \log \frac{1}{\varepsilon}$ . Thus,

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

Let's modify the process a bit by changing  $\frac{1}{2}$  into  $\frac{1}{n+1}$ , i.e. w.p.  $\frac{1}{n+1}$  do nothing, to make the lazy walk more active. Note that we add the lazy move in order to make the chain aperiodic.

Now in this case, we describe another coupling of  $X_t, Y_t$ . Without loss of generality, we can reorder the entries of two vectors so that all disagreeing entries come first. Namely there exists an index  $k$  such that  $X_t(i) \neq Y_t(i)$  if  $1 \leq i \leq k$ , and  $X_t(i) = Y_t(i)$  for  $i > k$ . Our coupling is as follows:

- If  $k = 0$ ,  $Y$  acts the same as  $X$ .
- If  $k = 1$ ,  $Y$  acts the same as  $X$  except when  $X$  flips the first entry,  $Y$  does nothing and vice versa.
- For  $k > 2$ , we distinguish between whether  $X$  flip indices in  $[k]$ :
  - If  $X$  did nothing or flipped one of  $i > k$ :  $Y$  acts the same.
  - If  $X$  flipped  $1 \leq i \leq k$ :  $Y$  flips  $(i \bmod k) + 1$ , i.e.  $1 \mapsto 2, 2 \mapsto 3, \dots, k-1 \mapsto k, k \mapsto 1$ .

It's clear that the above is indeed a coupling. In fact, this coupling acts like a doubled speed coupon collector, since in the case  $k > 2$  we can always collect two coupons at a time when lady luck is smiling. It is therefore conceivable that

$$\tau_{\text{mix}} \leq \frac{1}{2} n \log n + O(n).$$

**Example 2** (Shuffling cards). . Given a deck of  $n$  cards, consider the following rule of shuffling

- pick a card uniformly at random;
- put the card on the top.

The shuffling rule can be viewed as a random walk on all  $n!$  permutations of the  $n$  cards and it is easy to verify that the uniform distribution is the stationary distribution. Let us design a coupling for this Markov chain. That is, let  $X_t$  and  $Y_t$  be decks of cards, and we construct  $X_{t+1}$  and  $Y_{t+1}$  by

- picking the same random card and put it on the top.

This is clearly a coupling, and once some card, say  $\heartsuit K$  has been picked, then  $\heartsuit K$  in two decks will be always at the same location. Therefore, if we ask in how many rounds  $T$ ,  $X_T = Y_T$ , then the question is equivalent to the coupon collector problem again. So we have,

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

Note that we are picking the “same card”, not the card at the same location. That is, we draw a random card from  $X_t$ , say  $\heartsuit K$ , and then we pick  $\heartsuit K$  in  $Y_t$  as well.

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

March 23, 2023

## 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 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 Sample Proper Coloring

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 proper coloring problem is NP-hard in general. However, for  $q > \Delta$  there always exists a proper coloring that can be easily obtained by a greedy algorithm, where  $\Delta$  is the maximum degree of the graph.

If we want to count the number of proper colorings, then the problem becomes harder. It is known that for every  $q \geq \Delta$ , the problem is #P-hard. On the other hand, we can use a uniform sampler to obtain an algorithm to

approximately counting the number of proper coloring, at an arbitrarily low cost in the precision.

In fact, it is known that an approximate counting algorithm is equivalent to a uniform sampler in many cases (for example, sampling proper coloring). We only show one direction here: a sampler implies an algorithm for approximate counting. Given a graph  $G = (V, E)$  with  $V = [n]$ , let  $C$  be the set of proper colorings and  $Z = |C|$ . Suppose we have an oracle that can uniformly generate a proper coloring from  $C$ . Fix a proper coloring  $\sigma$ . We have

$$\begin{aligned} \frac{1}{Z} &= \Pr_{x \sim C} [x = \sigma] \\ &= \Pr_{x \sim C} [x(1) = \sigma(1) \wedge x(2) = \sigma(2) \wedge \dots] \\ &= \prod_{i=1}^n \Pr \left[ x(i) = \sigma(i) \mid \bigcap_{j < i} x(j) = \sigma(j) \right]. \end{aligned}$$

The above probability can be estimated by taking a number of samples from the oracle, and computing the ratio between colorings such that  $x(j) = \sigma(j)$  for  $j \leq i$  and ones that  $x(j) = \sigma(j)$  for  $j < i$ . Moreover, the ratio we just estimated is bounded below by an inverse polynomial and therefore polynomial number of sample suffices to estimate ratio accurately. The strategy works even if the sampler is an approximate one. Hence one can approximately compute  $Z$ . See [JV86] for more details.

Now we use MCMC to do sampling. Consider the following Markov chain to sample proper colorings:

- Pick  $v \in V$  and  $c \in [q]$  uniformly at random.
- Recolor  $v$  with  $c$  if possible.

The chain is aperiodic since self-loops exist in the walk. For  $q \geq \Delta + 2$ , the chain is irreducible. The bound  $q \geq \Delta + 2$  is tight for irreducibility since when  $q = \Delta + 1$ , each proper coloring of complete graph is frozen. It is still an open problem if the mixing time of the chain is polynomial in the size of the graph under the condition  $q \geq \Delta + 2$ . The best bound so far requires that  $q \geq (\frac{11}{6} - \varepsilon)\Delta$  for a certain constant  $\varepsilon > 0$ . Here, we shall give a rapid mixing proof when  $q > 4\Delta$  using the method of coupling.

The coupling we used is simple: Both players pick same  $v$  and  $c$  to move. However, we are not able to reduce the analyze the coupling to *coupon collector* as we did before. We introduce a more general method to analyze couplings. We define a certain distance  $d(x, y)$  for any two configurations  $x, y \in \Omega$ . We can assume without loss of generality that if  $x \neq y$  then  $d(x, y) \geq 1$  since  $\Omega$  is finite. Consider a coupling  $\omega_t$  of  $\mu_t, \nu_v$ . Then for every  $t \geq 0$  and  $(X_t, Y_t) \sim \omega_t$ , we try to establish

$$\mathbb{E} [d(X_{t+1}, Y_{t+1}) \mid (X_t, Y_t)] \leq (1 - \alpha)d(X_t, Y_t)$$

It is indeed a Metropolis algorithm. Let

$$\sigma^{v \leftarrow c}(u) = \begin{cases} \sigma(u) & u \neq v \\ c & u = v. \end{cases}$$

$\sigma^{v \leftarrow c}$  is a neighbor of  $\sigma$  on the transition graph, and we accept it if  $\sigma^{v \leftarrow c}$  is a proper coloring, i.e.  $\frac{\pi(\sigma^{v \leftarrow c})}{\pi(\sigma)} = 1$ .



for some  $\alpha \in (0, 1]$ . In other words,  $\{d(X_t, Y_t)\}_{t \geq 0}$  is a supermartingale. This implies that for every  $t \geq 1$ ,

$$\mathbb{E}[d(X_t, Y_t)] \leq (1 - \alpha)\mathbb{E}[d(X_{t-1}, Y_{t-1})] \leq (1 - \alpha)^t d(X_0, Y_0).$$

If we have a universal upper bound for  $d(X_0, Y_0)$ , say  $n$ , then by coupling lemma

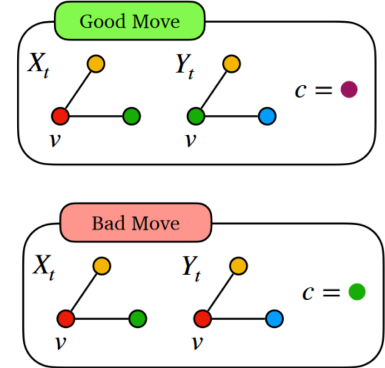
$$\begin{aligned} D_{\text{TV}}(\mu_t, \nu_t) &\leq \Pr_{(X_t, Y_t) \sim \omega_t} [X_t \neq Y_t] \\ &= \Pr[d(X_t, Y_t) \geq 1] \\ &\leq \mathbb{E}[d(X_t, Y_t)] \\ &\leq (1 - \alpha)^t \cdot n. \end{aligned}$$

Now come back to our problem of sampling proper colorings. 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 kinds:

- Good move:  $X_t(v) \neq Y_t(v)$ , and both change into  $c$  successfully. It will decrease distance by 1.
- Bad move:  $X_t(v) = Y_t(v)$ , one succeeds and one fails in the changing. It 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

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



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  $c$ . By a counting argument, we have

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

Therefore,

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

In the case  $q > 4\Delta$ , if we want

$$D_{\text{TV}} \leq \left(1 - \frac{1}{nq}\right)^t n \leq \varepsilon,$$

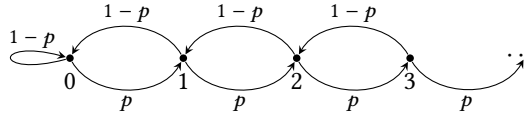
we have the mixing time is bounded by

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

#### 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

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

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

...

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

...

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.  $\square$

**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,  $\mathbf{P}_i[T_j = \infty] = 0$ .

- 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.

$\square$

#### References

[JV86] Mark R Jerrum, Leslie G Valiant, and Vijay V Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoretical computer science*, 43:169–188, 1986. [3](#)

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

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  $\hat{\Omega} \rightarrow \mathbb{R}$ .

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

# [AI2613 Lecture 5] FT of Countably Infinite Markov Chains, Some Applications of Markov Chains

June 10, 2023

## 1 Recurrence and Positive Recurrence

Recall that we say a state  $i$  is *recurrent* if  $\mathbf{P}_i[T_i < \infty] = 1$ . This is equivalent to  $\mathbf{E}_i[N_i] = \infty$ . Otherwise, we say the state is *transient*. A transient state  $j$  will be visited for finite times with probability 1. From Proposition 3 of last lecture, we know that *recurrence* is a class property, that is, given a recurrent state  $i$ , all the other states that  $i$  can reach in finite steps are also recurrent. We are only concerned with irreducible Markov chains in this lecture. So we may say a Markov chain is recurrent or transient in the future.

**Example 1 (Drunk person and drunk bird)** *Imagine a random walk on a grid that we pick a direction uniformly at random at each time step. Can we go back to the original point with probability 1? Or equivalently, is this Markov chain recurrent or transient?*

First we consider the one-dimensional grid. Let  $X_0 = 0$  and  $X_{t+1} = X_t + \Delta$  where  $\Delta$  is uniformly at random picked from  $\{-1, 1\}$ . Then,

$$\mathbf{E}_0[N_0] = \mathbf{E}_0\left[\sum_{t=0}^{\infty} \mathbb{1}[X_t = 0]\right] = \sum_{t=0}^{\infty} \mathbf{P}_0[X_t = 0] = \sum_{m=0}^{\infty} \mathbf{P}_0[X_{2m} = 0].$$

where the last equality follows from the fact that we can not go back within exactly odd steps. Then let's compute  $\mathbf{P}_0[X_{2m} = 0]$  using the *Stirling's formula*. For  $m \geq 1$ ,

$$\mathbf{P}_0[X_{2m} = 0] = \frac{\binom{2m}{m}}{2^{2m}} \approx \frac{\sqrt{4\pi m} \left(\frac{2m}{e}\right)^{2m}}{2\pi m \left(\frac{m}{e}\right)^{2m}} \cdot 2^{-2m} = \frac{1}{\sqrt{\pi m}}.$$

Thus,  $\mathbf{E}_0[N_0] = \sum_{m=0}^{\infty} \mathbf{P}_0[X_{2m} = 0] \approx 1 + \sum_{m=1}^{\infty} \frac{1}{\sqrt{\pi m}}$  which is divergent. This indicates that the Markov chain for random walk on one-dimensional grid is recurrent.

For  $d$ -dimensional grid, we regard the game as independently pick  $\Delta_i$  u.a.r. from  $\{-1, 1\}$  for  $i \in [d]$  at each time step and walk to  $X_{t+1} = X_t + (\Delta_1, \Delta_2, \dots, \Delta_d)$ . So we have that  $\mathbf{P}_i[X_{2m} = \mathbf{0}] = (\mathbf{P}_i[X_{2m}(1) = 0])^d \approx \left(\frac{1}{\sqrt{\pi m}}\right)^d$ . We know that  $1 + \sum_{m=1}^{\infty} \left(\frac{1}{\sqrt{\pi m}}\right)^d$  is divergent if and only if  $d \leq 2$ . Thus, only if the dimension of the grid is 1 or 2, the random walk is recurrent.

**Definition 1 (Positive recurrence)** If a state  $i$  is recurrent and  $\mathbf{E}_i[T_i] < \infty$ , we say it is positive recurrent. If the state is recurrent but with  $\mathbf{E}_i[T_i] = \infty$ , then we say it is null recurrent.

**Example 2 (Drunk person)** We have proved that the Markov chain of drunk person is recurrent. One can show that, even in one-dimension, the chain is null transient (exercise).

In fact  $\mathbf{P}_i[T_i < \infty] = 1 \iff \mathbf{P}_i[N_i = \infty] = 1 \iff \mathbf{E}_i[N_i] = \infty$ . I will leave the proof of this as an exercise.

Here we follow the notations of the last lecture, that is:  $X_0, X_1, \dots, X_t, \dots$  is a sequence of variables that follows the Markov chain  $P$ .  $T_i \triangleq \inf\{t > 0 : X_t = i\}$ ,  $N_i \triangleq \sum_{t=0}^{\infty} \mathbb{1}[X_t = i]$ ,  $\mathbf{P}_i[\cdot] = \Pr[\cdot | X_0 = i]$  and  $\mathbf{E}_i[\cdot] = \mathbf{E}[\cdot | X_0 = i]$ .

Stirling's formula:  $n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n (1 + o(1))$ .

### 1.1 1-D Random Walk

Consider the following one-dimensional random walk:



Let  $X_t$  be the position at time step  $t$ . Let  $T_{i \rightarrow j}$  be the first hitting time of state  $j$  starting from  $i$ , that is,  $T_{i \rightarrow j} = \min \{t > 0 | X_t = j \wedge X_0 = i\}$ . Define event  $\mathcal{A} = [\text{the first step is to the right}]$ . Then we consider the problem that when will this Markov chain be recurrent. Note that

$$\begin{aligned} \Pr [T_{0 \rightarrow 0} < \infty] &= \Pr [T_{0 \rightarrow 0} < \infty | \bar{\mathcal{A}}] \Pr [\bar{\mathcal{A}}] + \Pr [T_{0 \rightarrow 0} < \infty | \mathcal{A}] \Pr [\mathcal{A}] \\ &= (1-p) \cdot 1 + p \cdot \Pr [T_{1 \rightarrow 0} < \infty], \end{aligned} \quad (1)$$

$$\begin{aligned} \Pr [T_{1 \rightarrow 0} < \infty] &= \Pr [T_{1 \rightarrow 0} < \infty | \bar{\mathcal{A}}] \Pr [\bar{\mathcal{A}}] + \Pr [T_{1 \rightarrow 0} < \infty | \mathcal{A}] \Pr [\mathcal{A}] \\ &= (1-p) \cdot 1 + p \cdot \Pr [T_{2 \rightarrow 0} < \infty], \end{aligned} \quad (2)$$

$$\begin{aligned} \Pr [T_{2 \rightarrow 0} < \infty] &= \Pr [T_{2 \rightarrow 1} < \infty \wedge T_{1 \rightarrow 0} < \infty] \\ &= \Pr [T_{2 \rightarrow 1} < \infty] \cdot \Pr [T_{1 \rightarrow 0} < \infty] \\ &= \Pr [T_{1 \rightarrow 0} < \infty]^2. \end{aligned} \quad (3)$$

Let  $y \triangleq \Pr [T_{1 \rightarrow 0} < \infty]$  for brevity. Combine Equation (2) and Equation (3), we have  $y = 1 - p + py^2$  which then yields  $y = 1$  or  $y = \frac{1-p}{p}$ . By Equation (1),  $\Pr [T_{0 \rightarrow 0} < \infty] = 1$  or  $2 - 2p$ .

- When  $p < \frac{1}{2}$ ,  $2 - 2p$  is meaningless as a probability. So  $\Pr [T_{0 \rightarrow 0} < \infty] = 1$  and the Markov chain is recurrent.
- When  $p = \frac{1}{2}$ ,  $2 - 2p = 1$ . The Markov chain is also recurrent in this situation.
- When  $p > \frac{1}{2}$ , we verify that  $\Pr [T_{0 \rightarrow 0} < \infty] < 1$ , and therefore  $\Pr [T_{0 \rightarrow 0} < \infty] = 2 - 2p$ . Let  $\{\Delta_k\}_{k=0}^{\infty}$  be a sequence of i.i.d. random variables with

$$\Delta_k = \begin{cases} +1, & \text{w.p. } p \\ -1, & \text{w.p. } 1-p \end{cases}.$$

Given a sufficiently large  $n \in \mathbb{N}$ , we can walk to  $n$  from 0 in  $n$  steps (i.e.  $X_n = n$ ) with probability  $p^n > 0$ . Assume that we have arrived at  $n$ , consider the probability that we go back to 0 from  $n$  in exactly  $k$  steps. Apparently, this probability is zero when  $k < n$ . For every  $k \geq n$ , we

upper bound the probability  $\Pr [T_{n \rightarrow 0} = k]$ :

$$\begin{aligned} \Pr [T_{n \rightarrow 0} = k] &\leq \Pr \left[ \sum_{t=1}^k \Delta_t = -n \right] \\ &\leq \Pr \left[ \sum_{t=1}^k \Delta_t - \mathbb{E} \left[ \sum_{t=1}^k \Delta_t \right] \leq -n - \mathbb{E} \left[ \sum_{t=1}^k \Delta_t \right] \right] \\ &\leq \exp \left\{ -\frac{2k \left( \frac{n + (2p-1)k}{k} \right)^2}{4} \right\}. \end{aligned}$$

where the third inequality follows from the [Hoeffding's inequality](#).

Then we calculate the probability that we can go back from  $n$  to 0. By union bound,

$$\begin{aligned} \Pr [T_{n \rightarrow 0} < \infty] &= \Pr \left[ \bigcup_{k \geq n} [T_{n \rightarrow 0} = k] \right] \\ &\leq \sum_{k=n}^{\infty} \Pr [T_{n \rightarrow 0} = k] \\ &\leq \exp\{-(2p-1)n\} \sum_{k=n}^{\infty} \exp \left\{ -\frac{n^2}{2k} - \frac{(2p-1)^2 k}{2} \right\}. \end{aligned}$$

Note that

$$\begin{aligned} \sum_{k=n}^{\infty} \exp \left\{ -\frac{n^2}{2k} \right\} \cdot \exp \left\{ -\frac{(2p-1)^2 k}{2} \right\} &\leq \sum_{k=n}^{\infty} \exp \left\{ -\frac{(2p-1)^2 k}{2} \right\} \\ &= \frac{\exp \left\{ -\frac{(2p-1)^2}{2} n \right\}}{1 - \exp \left\{ -\frac{(2p-1)^2}{2} \right\}} \end{aligned}$$

Thus,

$$\Pr [T_{n \rightarrow 0} < \infty] \leq \frac{\exp \left\{ -\frac{(2p-1)^2}{2} n - (2p-1)n \right\}}{1 - \exp \left\{ -\frac{(2p-1)^2}{2} \right\}}. \quad (4)$$

We can find a sufficiently large constant  $n$  such that  $\Pr [T_{n \rightarrow 0} < \infty] < 1$  since the RHS of Equation (4) is exponentially small with regard to  $n$ . So for sufficiently large  $n$ , the probability that we walk to  $n$  and never come back to 0 is larger than  $p^n \cdot \Pr [T_{n \rightarrow 0} = \infty] > 0$ . Thus, this Markov chain is transient.

Now we verify that the Markov chain is positive recurrent when  $p < \frac{1}{2}$  and null recurrent when  $p = \frac{1}{2}$ . Note that

$$T_{0 \rightarrow 0} = \mathbb{1}[\bar{\mathcal{A}}] \cdot 1 + \mathbb{1}[\mathcal{A}](1 + T_{1 \rightarrow 0}) \quad (5)$$

$$T_{1 \rightarrow 0} = \mathbb{1}[\bar{\mathcal{A}}] \cdot 1 + \mathbb{1}[\mathcal{A}](1 + T_{2 \rightarrow 0}) \quad (6)$$

$$T_{2 \rightarrow 0} = T_{2 \rightarrow 1} + T_{1 \rightarrow 0}. \quad (7)$$

Note that  $E[T_{2 \rightarrow 1}] = E[T_{1 \rightarrow 0}]$ . Taking the expectation of Equation (6) and combining with Equation (7), we have

$$E[T_{1 \rightarrow 0}] = 1 - p + p(1 + 2E[T_{1 \rightarrow 0}]),$$

which yields  $E[T_{1 \rightarrow 0}] = \frac{1}{1-2p}$ . Take the expectation of Equation (5), we get  $E[T_{0 \rightarrow 0}] = \frac{1-p}{1-2p}$ . Thus:

- When  $p = \frac{1}{2}$ ,  $E[T_{0 \rightarrow 0}] = \infty$  and the Markov chain is null recurrent.
- When  $p < \frac{1}{2}$ ,  $E[T_{0 \rightarrow 0}] < \infty$  and the Markov chain is positive recurrent.

## 2 Some Applications

### 2.1 Galton-Watson Process

The model was formulated by F. Galton in the study of the survival and extinction of family names. In the nineteenth century, there was concern amongst the Victorians that aristocratic surnames were becoming extinct. In 1873, Galton originally posed the question regarding the probability of such an event, and later H. W. Watson replied with a solution.

Using more modern terms, the process can be defined formally as follows:

**Definition 2 (Galton-Watson Process)** *Suppose that all the individuals reproduce independently of each other and have the same offspring distribution. More precisely, let  $G_t$  denote the number of individuals of  $t$ -th generation:*

- We start from the zero generation. For convenience, let  $G_0 = 1$ .
- Each individual of generation  $t$  gives birth to a random number of children of generation  $t + 1$ . That is,  $\forall t \geq 0$  and  $i \in [G_t]$ , let  $X_{t,i}$  denote the number of children of the  $i$ -th individual in the  $t$ -th generation. Then  $\{X_{t,i}\}$  is an array of i.i.d. random variables with  $\Pr[X_{t,i} = k] = p_k$ .
- All individuals of generation  $t + 1$  are children of individuals of generation  $t$ :

$$G_{t+1} = \sum_{i=1}^{G_t} X_{t,i}$$

It is clear that the process  $\{G_t\}_{t \geq 0}$  is a Markov chain.

Denote by  $\rho$  the probability of extinction, namely

$$\rho \triangleq \Pr[\text{extinction}] = \Pr[\cup_{t \geq 1} \{G_t = 0\}].$$

Then the question is to determine the value of  $\rho$ . First we consider two trivial situations:

- When  $p_0 = 0$ , it is clear that there will be offspring and  $\rho = 0$ .

- When  $p_0 > 0$  and  $p_0 + p_1 = 1$ , we can verify that  $\rho = 1$ . We know that

$$\mathbf{E}[G_{t+1}|G_t] = p_1 \cdot G_t.$$

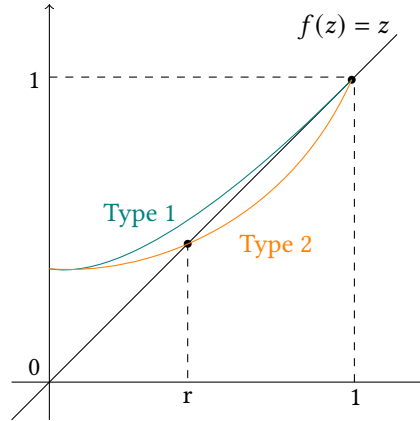
Compute the expectation of both sides, we have  $\mathbf{E}[G_{t+1}] = p_1 \mathbf{E}[G_t]$ .

This yields that when  $t \rightarrow \infty$ ,  $\Pr[G_t \geq 1] \leq \mathbf{E}[G_t] = p_1^t \mathbf{E}[G_0] \rightarrow 0$ .

Then we assume that  $p_0 > 0$  and  $p_0 + p_1 < 1$ . By the independence of each individual and the Markov property, we can calculate  $\rho$  as follows:

$$\begin{aligned} \rho &= \sum_{k=0}^{\infty} \Pr[\text{extinction} \wedge G_1 = k] \\ &= \sum_{k=0}^{\infty} \Pr[\text{extinction} | G_1 = k] p_k \\ &= \sum_{k=0}^{\infty} \rho^k p_k. \end{aligned} \tag{8}$$

Let  $\psi(z) \triangleq \sum_{k=0}^{\infty} p_k z^k$ . Then Equation (8) yields that  $\rho$  is a fixed point of  $\psi$ , i.e.,  $\psi(\rho) = \rho$ . By direct calculation we know  $\psi$  is an increasing and convex function on  $[0, 1]$  with  $\psi(0) = p_0$  and  $\psi(1) = 1$ . Then there can be two types of  $\psi$  depending on whether  $\psi'(1)$  is larger than 1 as the following figure shows.



When  $\psi'(1) = \sum_{k=1}^{\infty} k p_k = \mathbf{E}[X_{t-i}] \leq 1$ ,  $z = 1$  is the only fixed point of  $\psi$  which corresponds to the Type 1 in the figure. That is to say, when  $\mathbf{E}[X_{t-i}] \leq 1$ , we have  $\rho = 1$ .

When  $\mathbf{E}[X_{t-i}] > 1$  (Type 2), although there are two fixed points of  $\psi$ :  $r$  and 1, we claim that  $\rho = r$  rather than 1 by showing that  $\rho \leq r$ . Let  $q_t \triangleq \Pr[G_t = 0]$ . Then  $q_t \leq q_{t+1} < 1$  since  $G_t = 0$  can always yields  $G_{t+1} = 0$ . We induct on  $t$  to show that  $q_t \leq r$ :

- When  $t = 0$ , it is obvious that  $q_0 = 0 < r$ .
- Assume that  $q_t \leq r$ . Since  $q_{t+1} = \sum_{k=0}^{\infty} p_k q_t^k = \psi(q_t)$  and  $\psi$  is an increasing function,  $q_{t+1} = \psi(q_t) \leq \psi(r) = r$ .



We know that  $\rho = \lim_{t \rightarrow \infty} q_t$  and  $q_t \leq r$  for all  $t \geq 0$ . Thus  $\rho \leq r$ . However, we have shown that  $\rho$  is a fixed point of  $\psi$ . So  $\rho = r$  when  $\mathbb{E}[X_{t-i}] > 1$ . In conclusion,  $\rho = 1$  iff  $\mathbb{E}[X_{t-i}] \leq 1$ .

## 2.2 2-SAT

SAT is the problem of determining whether a CNF formula has satisfying assignments.  $k$ -SAT is the special cases of SAT that the clauses of the CNF formula consist of exact  $k$  literals. For example,

$$\phi = (x \vee y) \wedge (y \vee \bar{z}) \wedge (\bar{x} \vee z)$$

is a 2-CNF formula and  $x = y = z = 1$  is one of its satisfying assignments. SAT is NP-complete and we have  $k$ -SAT  $\in$  NP for  $k \geq 3$ . One can use an algorithm for finding strongly connected components to solve 2-SAT problem in linear time. Nevertheless, we introduce a simple randomized algorithm that can also solve this problem in polynomial-time with high probability.

We will extend the algorithm to solving 3-SAT in the homework!

Let  $\phi$  be a 2-CNF formula and  $V = \{v_1, v_2, \dots, v_n\}$  be its set of variables. The algorithm runs as follows:

- Pick an arbitrary assignment  $\sigma_0 : V \rightarrow \{\text{true}, \text{false}\}$ .
- For  $t = 0, 1, 2, \dots, 100n^2$ :
  - If  $\sigma_t$  satisfies  $\phi$ , output  $\sigma_t$ ;
  - Else, pick an arbitrary unsatisfying clause, say  $c = x \vee y$ . Choose from  $\{x, y\}$  uniformly at random and flip the assignment of the chosen literal. Let  $\sigma_{t+1}$  be the flipped assignment.
- Output “ $\phi$  is not satisfiable”.

**Claim 3** *This algorithm outputs the correct answer with probability at least  $1 - \frac{1}{100}$ .*

*Proof.* It is clear that if a 2-SAT instance has no solution then our algorithm will always give the correct answer. So we consider the probability that our algorithm outputs no solution conditioned on that the instance indeed has a satisfying assignment.

Our algorithm produces  $100n^2 + 1$  assignments  $\sigma_0, \sigma_1, \dots, \sigma_{100n^2}$ . We claim that with probability at least  $1 - \frac{1}{100}$ , some of  $\sigma_k$  for  $k \in \{0, \dots, 100n^2 + 1\}$  is a satisfying assignment. The argument here, at first glance, is a bit weird. We fix an arbitrary  $\sigma : V \rightarrow \{\text{true}, \text{false}\}$  satisfying assignment. We in fact prove the following: For large enough  $k$ , conditioned on the event that none of  $\sigma_0, \sigma_1, \dots, \sigma_k$  is a satisfying assignment,  $\sigma_{k+1} = \sigma$  holds with high probability.

Let  $\{X_t\}_{t=0}^{100n^2}$  be a random variable sequence that

$$X_t \triangleq |\{v \in V : \sigma_t(v) = \sigma(v)\}|.$$

First we verify that  $\Pr[X_{t+1} = X_t + 1 \mid \sigma_t] \geq \frac{1}{2}$ <sup>1</sup> and  $\Pr[X_{t+1} = X_t - 1 \mid \sigma_t] \leq \frac{1}{2}$ . WLOG assume we chose the clause  $c = x \vee y$  in round  $t$ . Since  $c$  is not satisfied by  $\sigma_t$ , we have  $\sigma_t(x) = \sigma_t(y) = \text{false}$ . Similarly,  $x \vee y$  is satisfying under  $\sigma$ , so there are three possible assignments of  $\sigma(x)$  and  $\sigma(y)$ :

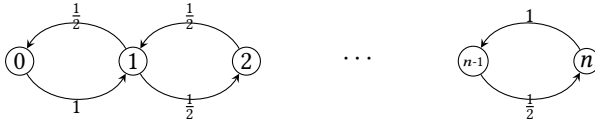
- If  $\sigma(x) = \text{true}$  and  $\sigma(y) = \text{false}$ ,  $\Pr[X_{t+1} = X_t + 1 \mid \sigma_t] = \Pr[\text{flip } x] = \frac{1}{2}$  and  $\Pr[X_{t+1} = X_t - 1 \mid \sigma_t] = \Pr[\text{flip } y] = \frac{1}{2}$ .
- If  $\sigma(x) = \text{false}$  and  $\sigma(y) = \text{true}$ , we have  $\Pr[X_{t+1} = X_t + 1 \mid \sigma_t] = \Pr[X_{t+1} = X_t - 1 \mid \sigma_t] = \frac{1}{2}$  similarly.
- If  $\sigma(x) = \text{true}$  and  $\sigma(y) = \text{true}$ ,  $\Pr[X_{t+1} = X_t + 1 \mid \sigma_t] = \Pr[\text{flip } x \text{ or } y] = 1$ .

Thus we have  $\Pr[X_{t+1} = X_t + 1 \mid \sigma_t] \geq \frac{1}{2}$  on condition that none of  $\sigma_0, \sigma_1, \dots, \sigma_t$  is a satisfying assignment.

Consider the 1-D random walk  $\{Y_t\}_{t \geq 0}$  on  $[n] \cup \{0\}$  that  $Y_0 = X_0$  and for  $Y_t \notin \{0, 1\}$

$$Y_{t+1} = \begin{cases} Y_t + 1, & \text{w.p. } \frac{1}{2} \\ Y_t - 1, & \text{w.p. } \frac{1}{2} \end{cases}.$$

If  $Y_t = 0$ ,  $Y_{t+1} = Y_t + 1$  w.p. 1 and if  $Y_t = n$ , then  $Y_{t+1} = Y_t - 1$  w.p. 1.



Then we have<sup>2</sup>

$$\begin{aligned} \Pr[\text{the algorithm is correct}] &\geq \Pr[\exists t \in [0, 100n^2] \text{ s.t. } X_t = n] \\ &\geq \Pr[\exists t \in [0, 100n^2] \text{ s.t. } Y_t = n]. \end{aligned} \quad (9)$$

Assume that  $Y_0 = X_0 = i$ . Let  $T_{i \rightarrow n}$  be the first hitting time of  $n$  from  $i$ . Then

$$T_{i \rightarrow n} = \sum_{k=i}^{n-1} T_{k \rightarrow k+1}.$$

For  $i > 0$ , we have

$$\begin{aligned} T_{i \rightarrow i+1} &= \mathbb{1}[\mathcal{A}] + \mathbb{1}[\bar{\mathcal{A}}](1 + T_{i-1 \rightarrow i+1}) \\ &= \mathbb{1}[\mathcal{A}] + \mathbb{1}[\bar{\mathcal{A}}](1 + T_{i-1 \rightarrow i} + T_{i \rightarrow i+1}) \end{aligned}$$

Taking the expectation of both sides, we have  $\mathbf{E}[T_{i \rightarrow i+1}] = 2 + \mathbf{E}[T_{i-1 \rightarrow i}]$ .

Note that  $T_{0 \rightarrow 1} = 1$ , then

$$\mathbf{E}[T_{i \rightarrow n}] = \sum_{k=i}^{n-1} \mathbf{E}[T_{k \rightarrow k+1}] = \sum_{k=i}^{n-1} 2k + 1 = n^2 - i^2 \leq n^2.$$

Note that  $\{X_t\}_{t=0}^{100n^2}$  is not a Markov chain since it only contains partial information of  $\sigma_t$  and we cannot determine the distribution of  $X_{t+1}$  given  $X_t$ .

<sup>1</sup> Let  $Y$  be a random variable. Then function  $\Pr[\cdot \mid Y] : \text{Ran}(Y) \rightarrow \mathbb{R}$  is defined by  $\Pr[\cdot \mid Y] = \mathbf{E}[\mathbb{1}[\cdot] \mid Y]$ . Note that  $\Pr[\cdot \mid Y]$  is a random variable. Here we slightly abuse the notations and denote the event " $\forall a \in \text{Ran}(Y), \Pr[\cdot \mid Y = a] \geq \frac{1}{2}$ " as  $\Pr[\cdot \mid Y] \geq \frac{1}{2}$ .

<sup>2</sup> The second inequality can be verified by constructing a coupling which satisfies  $Y_t \geq X_t$  for all  $t \geq 0$ . The existence of such coupling is guaranteed by  $\Pr[X_{t+1} = X_t + 1 \mid \sigma_t] \geq \Pr[Y_{t+1} = Y_t + 1]$ . Specifically, if there is one false and one true in  $\{\sigma(x), \sigma(y)\}$ , then  $Y_{t+1}$  moves the same as  $X_{t+1}$ . If  $\sigma(x) = \sigma(y) = \text{true}$ , then  $Y_{t+1}$  moves +1 or -1 uniformly at random.

Recall  $\mathcal{A} = [\text{the first step is to the right}]$ .

Then we apply the Markov's inequality to give a lower bound for  $\Pr [\exists t \in [0, 100n^2] \text{ s.t. } Y_t = n]$ :

$$\begin{aligned} 1 - \Pr [\exists t \in [0, 100n^2] \text{ s.t. } Y_t = n] &= \Pr [T_{Y_0 \rightarrow n} > 100n^2] \\ &\leq \frac{\mathbf{E} [T_{Y_0 \rightarrow n}]}{100n^2} \leq \frac{1}{100}. \end{aligned}$$

By Equation (9), we know that  $\Pr$  [the algorithm is correct] is lower bounded by  $1 - \frac{1}{100}$ .  $\square$

### 3 Fundamental Theorem

In this section, we develop the fundamental theorem of Markov chains for chains with possibly infinite states. First we introduce some abbreviations to simplify the expression:

- Aperiodicity:[A],
- Irreducibility:[I],
- Recurrence:[R],
- Positive Recurrence: [PR],
- Has a stationary distribution:[S],
- Has a unique stationary distribution:[U],
- Convergence:[C],
- Finiteness:[F].

The finite FTMC can be written as:  $[F]+[A]+[I] \Rightarrow [S]+[U]+[C]$ . For infinite Markov chains, the theorem need to be modified as:  $[PR]+[A]+[I] \Rightarrow [S]+[U]+[C]$ .

Before the proof of the theorem, we need to prepare some mathematical tools.

#### 3.1 Laws of Large Numbers

$X_1, X_2, \dots$  is an infinite sequence of independent and identically distributed Lebesgue integrable random variables with expected value  $\mathbf{E} [X_1] = \mathbf{E} [X_2] = \dots = \mu < \infty$ . Let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  be the sample average. Then we have the following two laws of large numbers.

**Theorem 4 (Weak law of large numbers or Khinchin's law)** *The sample average converge in probability towards the expected value:*

$$\bar{X}_n \xrightarrow{P} \mu \quad \text{when } n \rightarrow \infty.$$

That is, for any positive value  $\epsilon$ ,

$$\lim_{n \rightarrow \infty} \Pr [|\bar{X}_n - \mu| < \epsilon] = 1.$$

**Theorem 5 (Strong law of large numbers or Kolmogorov's law)** *The sample average converges almost surely or with probability 1 to the expected value:*

$$\bar{X}_n \xrightarrow{\text{a.s.}} \mu \quad \text{when } n \rightarrow \infty.$$

That is,

$$\Pr \left[ \lim_{n \rightarrow \infty} \bar{X}_n \rightarrow \mu \right] = 1.$$

As the name of the laws shows, *convergence in probability* is weaker than *convergence with probability 1*. Consider a sequence of independent random variables  $X_1, X_2, \dots$  that  $X_n$  is 1 with probability  $\frac{1}{n}$  and  $X_n$  is 0 with probability  $1 - \frac{1}{n}$ . Then the sequence converges to 0 in probability but not with probability 1 since we cannot find an  $M \in \mathcal{F}$  with measure 1 such that  $\bar{X}_n(\omega) \xrightarrow{n \rightarrow \infty} 0$  for every  $\omega \in M$ .

Let  $(\Omega, \mathcal{F}, P)$  be the probability space. Here  $\bar{X}_n \rightarrow \mu$  means  $\exists M \in \mathcal{F}$  satisfying

- $P(M)=1$ ;
- $\forall \omega \in M, \bar{X}_n(\omega) \xrightarrow{n \rightarrow \infty} \mu$ .

**Theorem 6 (Strong law of large numbers for Markov chains)** *If there is a finite path from state  $i$  to  $j$ , then*

$$P_i \left[ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{1}[X_t = j] = \frac{1}{E_j[T_j]} \right] = 1.$$

*Proof.* If  $j$  is transient, then the random process will visit  $j$  for finite times with probability 1. Thus  $P_i \left[ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{1}[X_t = j] = \frac{1}{E_j[T_j]} \right] = P_i \left[ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{1}[X_t = j] = 0 \right] = 1$ .

If  $j$  is recurrent, we first prove the theorem for  $i = j$ . We call a loop from  $j$  to  $j$  a cycle (we visit  $j$  only at the beginning and end of the loop). Denote  $C_r$  as the length of the  $r^{\text{th}}$  cycle during the process. Let  $S_k = \sum_{r=1}^k C_r$ . Let  $k_n$  be the number of cycles before the  $n + 1$  step, that is,  $k_n = \max \{k | S_k \leq n\}$ . Then we have  $S_{k_n} \leq n < S_{k_n+1}$  and consequently  $\frac{S_{k_n}}{k_n} \leq \frac{n}{k_n} < \frac{S_{k_n+1}}{k_n}$ . Note that with probability 1,  $k_n \rightarrow \infty$  when  $n \rightarrow \infty$ . We have with probability 1 that

$$\lim_{k \rightarrow \infty} \frac{S_k}{k} \leq \lim_{n \rightarrow \infty} \frac{n}{k_n} < \lim_{k \rightarrow \infty} \frac{S_{k+1}}{k}.$$

Note that  $S_k = \sum_{r=1}^k C_r$  where each  $C_r$  is an i.i.d random variable with mean  $E_j[T_j]$ . So by SLLN (Theorem 5), we have  $\lim_{k \rightarrow \infty} \frac{S_k}{k} = E_j[T_j]$  and  $\lim_{k \rightarrow \infty} \frac{S_{k+1}}{k} = \lim_{k \rightarrow \infty} \frac{S_{k+1}}{k+1} \cdot \frac{k+1}{k} = E_j[T_j]$ . As a result, with probability 1,

$$E_j[T_j] = \lim_{n \rightarrow \infty} \frac{n}{k_n} = \lim_{n \rightarrow \infty} \frac{n}{\sum_{t=1}^n \mathbb{1}[X_t = j]}.$$

If  $j$  is recurrent and  $i \neq j$ , let  $T_{i \rightarrow j}$  be the first time visiting  $j$ . Then we have  $\frac{S_{k_n+T_{i \rightarrow j}}}{k_n} \leq \frac{n}{k_n} < \frac{S_{k_n+1+T_{i \rightarrow j}}}{k_n}$ . Since  $P_i[T_j < \infty] = 1$ ,  $P_i[\lim_{k \rightarrow \infty} \frac{T_{i \rightarrow j}}{k} = 0] = 1$ . The remaining proof is the same with the situation that  $i = j$ .  $\square$

**Corollary 7** *Let  $P$  be the transition function of an irreducible Markov chain where  $P^t(i, j) = \Pr[X_t = j | X_0 = i]$ . Then for any states  $i, j$ ,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, j) = \frac{1}{E_j[T_j]}.$$

*Proof.* By the strong law of large numbers for Markov chains, there exists a set  $M \in \mathcal{F}$  such that  $P(M) = 1$  and  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{1}[X_t(\omega) = j] = \frac{1}{\mathbb{E}_j[T_j]}$  for any  $\omega \in M$ . Then,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, j) &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}_i[\mathbb{1}[X_t = j]] \\ &= \lim_{n \rightarrow \infty} \mathbb{E}_i \left[ \frac{1}{n} \sum_{t=1}^n \mathbb{1}[X_t = j] \right] \\ &= \mathbb{E}_i \left[ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{1}[X_t = j] \right] \\ &= \frac{1}{\mathbb{E}_j[T_j]}, \end{aligned}$$

where the third equation follows from the [bounded convergence theorem](#).

□

Bounded Convergence Theorem: If  $X_n \xrightarrow{a.s.} X$  and  $\mathbb{E}[X] < \infty$ , then  $\mathbb{E}[X_n] \rightarrow \mathbb{E}[X]$ .

### 3.2 Proof of the Fundamental Theorem

We will first prove the existence and uniqueness of the stationary distribution in this lecture.(i.e. [S] and [U])

**Theorem 8**  $[I]+[PR] \Rightarrow [S]+[U]$ .

*Proof.* [Proof of [U]] Let  $\mathcal{S}$  be the set of states. Assume  $\pi$  is a stationary distribution of the Markov chain, i.e.,

$$\forall j \in \mathcal{S}, \forall t \geq 0, \sum_{i \in \mathcal{S}} \pi(i) P^t(i, j) = \pi(j).$$

This yields that for  $n \geq 1$ ,

$$\frac{1}{n} \sum_{i \in \mathcal{S}} \pi(i) \sum_{t=1}^n P^t(i, j) = \pi(j).$$

Taking  $n \rightarrow \infty$  and applying the [dominated convergence theorem](#), we have

$$\pi(j) = \sum_{i \in \mathcal{S}} \pi(i) \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, j) = \sum_{i \in \mathcal{S}} \pi(i) \cdot \frac{1}{\mathbb{E}_j[T_j]} = \frac{1}{\mathbb{E}_j[T_j]}.$$

Dominated Convergence Theorem: If  $\int_{\mathcal{S}} |f_n| < \infty$ , then  $\lim_{n \rightarrow \infty} \int_{\mathcal{S}} f_n = \int_{\mathcal{S}} \lim_{n \rightarrow \infty} f_n$ .

□

*Proof.* [Proof of [S]] Then we prove the above  $\pi$  is a stationary distribution.

$\mathcal{S}$  is finite. We first assume  $\mathcal{S}$  is finite, so that we can safely exchange the order of taking limitation and summation in the calculations below.

$$\begin{aligned} \sum_{j \in \mathcal{S}} \pi(j) &= \sum_{j \in \mathcal{S}} \frac{1}{\mathbb{E}_j[T_j]} = \sum_{j \in \mathcal{S}} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, j) \\ &= \lim_{n \rightarrow \infty} \sum_{j \in \mathcal{S}} \frac{1}{n} \sum_{t=1}^n P^t(i, j) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \sum_{j \in \mathcal{S}} P^t(i, j) = 1. \end{aligned}$$

This indicates that  $\pi$  is a legal distribution. We then verify that  $\pi$  is indeed the stationary distribution.

Note that  $P^{t+1}(i, j) = \sum_{k \in S} P^t(i, k)P(k, j)$ . Then

$$\begin{aligned} \frac{1}{\mathbf{E}_j[T_j]} &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, j) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^{t+1}(i, j) \\ &= \lim_{n \rightarrow \infty} \sum_{k \in S} P(k, j) \frac{1}{n} \sum_{t=1}^n P^t(i, k) = \sum_{k \in S} P(k, j) \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, k) \\ &= \sum_{k \in S} P(k, j) \cdot \frac{1}{\mathbf{E}_k[T_k]}. \end{aligned}$$

That is,

$$\pi(j) = \sum_{k \in S} P(k, j)\pi(k).$$

It is worth noting that [PR] is equivalent to [I] when  $S$  is finite.

*S is infinite.* When  $S$  is (countably) infinite, we consider every finite subset  $A$  of  $S$ . Then

$$\begin{aligned} \sum_{j \in A} \pi(j) &= \sum_{j \in A} \frac{1}{\mathbf{E}_j[T_j]} = \sum_{j \in A} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n P^t(i, j) \\ &= \lim_{n \rightarrow \infty} \sum_{j \in A} \frac{1}{n} \sum_{t=1}^n P^t(i, j) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \sum_{j \in A} P^t(i, j) < 1. \end{aligned}$$

Therefore

$$\sum_{j \in S} \pi(j) = \sup_{\text{finite } A \subseteq S} \sum_{j \in A} \pi(j) =: C \leq 1.$$

Since [PR], we know that  $C \neq 0$ . In the following, we will prove that  $\pi/C$  is a stationary distribution. Then  $C = 1$  follows from the uniqueness of the stationary distribution we just proved.

For every finite  $A \subseteq S$ , we have

$$\sum_{k \in A} P(k, j) \cdot \frac{1}{\mathbf{E}_k[T_k]} \leq \frac{1}{\mathbf{E}_j[T_j]}.$$

Therefore,

$$\sum_{k \in S} P(k, j) \cdot \frac{1}{\mathbf{E}_k[T_k]} = \sup_{\text{finite } A \subseteq S} \sum_{k \in A} P(k, j) \cdot \frac{1}{\mathbf{E}_k[T_k]} \leq \frac{1}{\mathbf{E}_j[T_j]}.$$

We show that indeed the equality holds. Assume for a contradiction that

$$\sum_{k \in S} P(k, j) \cdot \frac{1}{\mathbf{E}_k[T_k]} < \frac{1}{\mathbf{E}_j[T_j]}.$$

Summing the both sides over all  $j \in S$ , we obtain

$$\sum_{k \in S} \frac{1}{\mathbf{E}_k[T_k]} < \sum_{j \in S} \frac{1}{\mathbf{E}_j[T_j]},$$

which is a contradiction. As a result, we know

$$\sum_{k \in S} P(k, j) \cdot \frac{1}{\mathbf{E}_k [T_k]} = \frac{1}{\mathbf{E}_j [T_j]},$$

and  $\hat{\pi}(j) = \frac{1}{C \cdot \mathbf{E}_j [T_j]}$  is a stationary distribution. By the uniqueness of the distribution, we have  $C = 1$ .

□

# [AI2613 Lecture 6]: Concentration Inequalities, Martingale

April 6, 2023

## 1 Chernoff Bounds

Recall the Markov inequality and Chebyshev's inequality we introduced before. They are used to prove that a random variable is concentrated around its expectation.

If we apply Markov inequality to

$$\Pr[f(X) \geq f(t)]$$

with  $f(x) = e^{\alpha x}$  where  $\alpha > 0$ , then the bound amounts to bound  $\mathbb{E}[e^{\alpha X}]$  which is the *moment generating function* of  $X$ .

When the random variable  $X$  can be written as the sum of independent Bernoulli variables, its moment generating function is easy to estimate and we obtain sharp concentration bounds.

**Theorem 1 (Chernoff Bound)** . Let  $X_1, \dots, X_n$  be independent random variables such that  $X_i \sim \text{Ber}(p_i)$  for each  $i = 1, 2, \dots, n$ . Let  $X = \sum_{i=1}^n X_i$  and denote  $\mu \triangleq \mathbb{E}[X] = \sum_{i=1}^n p_i$ , we have

$$\Pr[X \geq (1 + \delta)\mu] \leq \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^\mu$$

If  $0 < \delta < 1$ , then we have

$$\Pr[X \leq (1 - \delta)\mu] \leq \left( \frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^\mu$$

*Proof.* We only prove the upper tail bound and the proof of lower tail bound is similar. For every  $\alpha > 0$ , we have

$$\Pr[X \geq (1 + \delta)\mu] = \Pr[e^{\alpha X} \geq e^{\alpha(1+\delta)\mu}] \leq \frac{\mathbb{E}[e^{\alpha X}]}{e^{\alpha(1+\delta)\mu}}.$$

Therefore, we need to estimate the moment generating function  $\mathbb{E}[e^{\alpha X}]$ . Since  $X = \sum_{i=1}^n X_i$  is the sum of independent Bernoulli variables, we have

$$\mathbb{E}[e^{\alpha X}] = \mathbb{E}\left[e^{\alpha \sum_{i=1}^n X_i}\right] = \mathbb{E}\left[\prod_{i=1}^n e^{\alpha X_i}\right] = \prod_{i=1}^n \mathbb{E}[e^{\alpha X_i}].$$

Since  $X_i \sim \text{Ber}(p_i)$ , we can compute  $\mathbb{E}[e^{\alpha X_i}]$  directly:

$$\mathbb{E}[e^{\alpha X_i}] = p_i e^\alpha + (1 - p_i) = 1 + (e^\alpha - 1)p_i \leq e^{((e^\alpha - 1)p_i)}.$$

Therefore,

$$\mathbb{E}[e^{\alpha X}] \leq \prod_{i=1}^n e^{((e^\alpha - 1)p_i)} = e^{((e^\alpha - 1) \sum_{i=1}^n p_i)} = e^{((e^\alpha - 1)\mu)}.$$



Therefore,

$$\Pr [X \leq (1 + \delta)\mu] \leq \frac{\mathbb{E} [e^{\alpha X}]}{e^{\alpha(1+\delta)\mu}} \leq \left( \frac{e^{(e^\alpha - 1)}}{e^{(\alpha(1+\delta))}} \right)^\mu$$

Note that above holds for any  $\alpha > 0$ . Therefore, we can choose  $\alpha$  so as to minimize  $\frac{e^{(e^\alpha - 1)}}{e^{(\alpha(1+\delta))}}$ . To this end, we let  $\left( \frac{e^{(e^\alpha - 1)}}{e^{(\alpha(1+\delta))}} \right)' = 0$ . This gives  $\alpha = \log(1 + \delta)$ . Therefore

$$\Pr [X \leq (1 + \delta)\mu] \leq \left( \frac{e^{(e^\alpha - 1)}}{e^{(\alpha(1+\delta))}} \right)^\mu = \left( \frac{e^\delta}{(1 + \delta)^{(1+\delta)}} \right)^\mu.$$

□

The following form of Chernoff bound is more convenient to use (but weaker):

**Corollary 2** For any  $0 < \delta < 1$ ,

$$\begin{aligned} \Pr [X \geq (1 + \delta)\mu] &\leq \exp\left\{\left(-\frac{\delta^2}{3}\mu\right)\right\} \\ \Pr [X \leq (1 - \delta)\mu] &\leq \exp\left\{\left(-\frac{\delta^2}{2}\mu\right)\right\} \end{aligned}$$

*Proof.* We only prove the upper tail. It suffices to verify that for  $0 < \delta < 1$ , we have

$$\frac{e^\delta}{(1 + \delta)^{(1+\delta)}} \leq \exp\left\{\left(-\frac{\delta^2}{3}\right)\right\}$$

Taking logarithm of both sides, this is equivalent to

$$\delta - (1 + \delta) \ln(1 + \delta) \leq -\frac{\delta^2}{3}$$

Let  $f(\delta) = \delta - (1 + \delta) \ln(1 + \delta) + \frac{\delta^2}{3}$  and note that

$$f'(\delta) = -\ln(1 + \delta) + \frac{2}{3}\delta, \quad f''(\delta) = -\frac{1}{1 + \delta} + \frac{2}{3}.$$

Then for  $0 < \delta < 1/2$ ,  $f''(\delta) < 0$ , and for  $1/2 < \delta < 1$ ,  $f''(\delta) > 0$ . Therefore,  $f'(\delta)$  first decreases and then increases in  $[0, 1]$ . Also note that  $f'(0) = 0$ ,  $f'(1) < 0$  and  $f'(\delta) \leq 0$  when  $0 \leq \delta \leq 1$ . Therefore  $f(\delta) \leq f(0) = 0$ . □

**Example 1 (Tossing  $p$ -coins)** . Consider a  $p$ -coin that we get a head with probability  $p$  when tossing it. If we toss a  $p$ -coin  $n$  times, the average number of heads is  $pn$ . We want to determine the value  $\delta$  such that with high probability (say 99%), the total number of heads is in the interval of  $[(1 - \delta)pn, (1 + \delta)pn]$ . We use Chernoff bound to determine  $\delta$ .

Let  $X$  denote the total number of heads, and  $X_i \sim \text{Ber}(p)$  be the indicator of whether the  $i$ -th toss gives a head. Then by Chernoff bound, we have

$$\Pr [|X - pn| \geq \delta \cdot pn] \leq 2 \exp\left\{\left(-\frac{\delta^2}{3} \cdot pn\right)\right\} \leq 0.01$$

So if  $p$  is a constant, it suffices to choose

$$\delta = \Omega\left(\frac{1}{\sqrt{n}}\right).$$

## 2 Hoeffding's Inequality

One of annoying restrictions of Chernoff bound is that each  $X_i$  needs to be a Bernoulli random variable. We first relax this requirement by introducing Hoeffding's inequality which allows  $X_i$  to follow any distribution, provided its value is almost surely bounded.

**Theorem 3 (Hoeffding's Inequality)** *Let  $X_1, \dots, X_n$  be independent random variables where each  $X_i \in [a_i, b_i]$  for certain  $a_i \leq b_i$  with probability 1. Let  $X = \sum_{i=1}^n X_i$  and  $\mu \triangleq \mathbb{E}[X] = \sum_{i=1}^n \mathbb{E}[X_i]$ , then*

$$\Pr[|X - \mu| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

for all  $t \geq 0$ .

It is instructive to compare Hoeffding and Chernoff when  $X_i$ 's are independent Bernoulli variables. Formally, let  $X_1, \dots, X_n$  be i.i.d. random variables where  $X_i \sim \text{Ber}(p)$  for all  $i = 1, \dots, n$ . Set  $X = \sum_{i=1}^n X_i$  and denote  $\mathbb{E}[X] = np$  by  $\mu$ . By Hoeffding's inequality, we have

$$\Pr[|X - \mu| \geq t] \leq 2 \exp\left(-\frac{2t^2}{n}\right).$$

By Chernoff Bound, we have

$$\Pr[|X - \mu| \geq t] \leq 2 \exp\left(-\frac{t^2}{3pn}\right).$$

Comparing the exponent, it is easy to see that for  $p > 1/6$ , Hoeffding's inequality is tighter up to a certain constant factor. However, for smaller  $p$ , Chernoff bound is significantly better than Hoeffding's inequality.

We consider the balls-in-a-bag problem. There are  $g$  green balls and  $r$  red balls in a bag. These balls are the all same except for the color. We want to estimate the ratio  $\frac{r}{r+g}$  by drawing balls. There are two scenarios.

- Draw balls with replacement. Let  $X_i = \mathbf{1}[\text{the } i\text{-th ball is red}]$ . Let  $X = \sum_{i=1}^n X_i$ . Then clearly each  $X_i \sim \text{Ber}\left(\frac{r}{r+g}\right)$  and  $\mathbb{E}[X] = n \cdot \frac{r}{r+g}$ .

Since all  $X_i$ 's are independent, we can directly apply Hoeffding's inequality and obtain

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{n}\right).$$

- Draw balls without replacement. Again we let  $Y_i = \mathbf{1}[\text{the } i\text{-th ball is red}]$ , then unlike the case of drawing with replacement, variables in  $\{Y_i\}$  are dependent. Let  $Y = \sum_{i=1}^n Y_i$ . We first calculate  $\mathbb{E}[Y]$ .

For every  $i \geq 1$ ,  $\mathbb{E}[Y_i]$  is the probability that the  $i$ -th draw is a red ball. Note that drawing without replacement is equivalent to first drawing a

uniform permutation of  $r + g$  balls and drawing each ball one by one in that order. Therefore, the probability of  $Y_i = 1$  is  $\frac{r \cdot (r+g-1)!}{(r+g)!} = \frac{r}{r+g}$ . So we have  $E[Y] = n \cdot \frac{r}{r+g}$ .

However, since  $\{Y_i\}$  are dependent, we cannot apply Hoeffding's inequality directly. This motivate us to generalize it by removing the requirement of independence.

### 3 Martingale

We develop the theory of martingale, which is a core concept in probability theory. We use martingale to get rid of the independence requirement in the concentration inequalities mentioned above.

Consider a fair gambling game in which the expected gain in each round is zero. As a result, regardless of how much one bets in each round, the money in expectation remains the same. The balances after each round form a *martingale*.

**Definition 4 (Martingale)** Let  $\{X_n\}_{n \geq 0}$  and  $\{Z_n\}_{n \geq 0}$  be two sequences of random variables. Let  $Z_n = \sum_{t=0}^n X_t$ .<sup>1</sup> We say  $\{Z_n\}_{n \geq 0}$  is a martingale w.r.t.  $\{X_n\}_{n \geq 0}$  if

$$E[Z_{n+1} \mid X_0, X_1, \dots, X_n] = Z_n.$$

Sometimes we say a single sequence  $\{X_n\}_{n \geq 0}$  is a martingale if it is a martingale w.r.t. itself. Formally, if for every  $n \geq 0$ , it holds that

$$E[X_{n+1} \mid X_0, \dots, X_n] = X_n.$$

For convenience, from now on we use  $\bar{X}_{i,j} = (X_i, X_{i+1}, \dots, X_j)$  to simplify the notations. The conditional expectation  $E[Z_{n+1} \mid \bar{X}_{0,n}]$  is equivalent to  $E[Z_{n+1} \mid \sigma(\bar{X}_{0,n})]$  where  $\sigma(\bar{X}_{0,n})$  is the  $\sigma$ -algebra generated by  $X_0, \dots, X_n$ . This motivates us to define martingale in a more general way.

**Definition 5 (Martingale (defined by filtration))** Let  $\{\mathcal{F}_n\}_{n \geq 0}$  be a sequence of  $\sigma$ -algebras. We call such  $\sigma$ -algebra sequence a filtration if it satisfies

$$\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots \subseteq \mathcal{F}_n \subseteq \mathcal{F}_{n+1} \subseteq \dots$$

Given a filtration  $\{\mathcal{F}_n\}_{n \geq 0}$ , let  $\{Z_n\}_{n \geq 0}$  be a stochastic process that  $Z_n$  is  $\mathcal{F}_n$ -measurable for every  $n \geq 0$ . Then we say  $\{Z_n\}_{n \geq 0}$  is a martingale w.r.t.  $\{\mathcal{F}_n\}_{n \geq 0}$  if for every  $n \geq 0$

$$E[Z_{n+1} \mid \mathcal{F}_n] = Z_n.$$

**Example 2 (1-D Random Walk)** Consider a random walk on  $\mathbb{Z}$  starting from 0. The probability to the left and the probability to the right are both  $\frac{1}{2}$  at each step. Denote the action at the  $n$ -th step by a uniform random variable

<sup>1</sup> Consider the problem of fair gambling where  $X_n$  is the gain of  $n$ -th round and  $Z_n = \sum_{t=1}^n X_t$ .  $\{Z_n\}_{n \geq 0}$  is not necessarily a Markov chain. The value  $X_n$  may depend on information before round  $n - 1$ .

If  $E[Z_{n+1} \mid \mathcal{F}_n] \leq Z_n$  in Definition 5, we call  $\{Z_n\}_{n \geq 0}$  a supermartingale w.r.t.  $\{\mathcal{F}_n\}_{n \geq 0}$ . Similarly, if  $E[Z_{n+1} \mid \mathcal{F}_n] \geq Z_n$ , we call it a submartingale.

$X_n \in \{-1, +1\}$ . Let  $S_n = \sum_{k=0}^n X_k$ . Then we can verify  $\{S_n\}_{n \geq 0}$  is a martingale w.r.t.  $\{X_n\}_{n \geq 0}$  (or w.r.t.  $\{S_n\}_{n \geq 0}$ ) by noticing that

$$\mathbb{E} \left[ S_{n+1} \mid \bar{X}_{0,n} \right] = \mathbb{E} \left[ S_n + X_{n+1} \mid \bar{X}_{0,n} \right] = S_n + \mathbb{E} \left[ X_{n+1} \mid \bar{X}_{0,n} \right] = S_n.$$

More generally, if  $\mathbb{E} \left[ X_k \mid \bar{X}_{0,n} \right] = \mu$ , we define  $Y_k = X_k - \mu$  and  $S'_n \triangleq \sum_{k=0}^n Y_k = S_n - (n+1)\mu$ . Then  $S'_n$  is a martingale w.r.t.  $\{Y_n\}_{n \geq 0}$ .

**Example 3** Consider a sequence of random variables  $\{X_n\}_{n \geq 0}$  where  $\mathbb{E} \left[ X_n \mid \bar{X}_{0,n-1} \right] = 1$  for all  $n \geq 1$ . Let  $P_n = \prod_{k=0}^n X_k$ . Then we can verify  $\{P_n\}_{n \geq 0}$  is a martingale w.r.t.  $\{X_n\}_{n \geq 0}$  by verifying that

$$\mathbb{E} \left[ P_{n+1} \mid \bar{X}_{0,n} \right] = \mathbb{E} \left[ P_n \cdot X_{n+1} \mid \bar{X}_{0,n} \right] = P_n \cdot \mathbb{E} \left[ X_{n+1} \mid \bar{X}_{0,n} \right] = P_n.$$

**Example 4 (Galton-Watson Process)** Recall the Galton-Watson process we discussed in the last lecture. Suppose that all the individuals reproduce independently of each other and have the same offspring distribution. Let  $G_t$  be the number of individuals of the  $t$ -th generation. Each individual of generation  $t$  gives birth to a random number of children of generation  $t+1$ . Denote by  $X_{t,k}$  the number of children of the  $k$ -th individual in the  $t$ -th generation. Assume that  $X_{t,k}$  are i.i.d. and let  $\mu \triangleq \mathbb{E} \left[ X_{t,k} \right]$ . Then we have  $G_{t+1} = \sum_{k=1}^{G_t} X_{t,k}$ . Thus,

$$\mathbb{E} \left[ G_{t+1} \mid G_t \right] = \mathbb{E} \left[ \sum_{k=1}^{G_t} X_{t,k} \mid G_t \right] = G_t \cdot \mathbb{E} \left[ X_{t,1} \right] = \mu G_t.$$

Define  $M_t = \mu^{-t} G_t$ . Then

$$\mathbb{E} \left[ M_{t+1} \mid G_t \right] = \mu^{-t-1} \mathbb{E} \left[ G_{t+1} \mid G_t \right] = \mu^{-t} G_t = M_t.$$

That is,  $\{M_t\}_{t \geq 0}$  is a martingale w.r.t.  $\{G_t\}_{t \geq 0}$ .

**Example 5 (Pólya's urn)** Suppose there are some white balls and black balls in an urn. All of these balls are identical except the colors. Consider the following stochastic process: each round we pick a ball uniformly at random and observe its color; then we return the ball, and add an additional ball of the same color into the urn. We repeat the process, and our goal is to study the sequence of colors of the selected balls.

W.l.o.g. assume that we start from only one white ball and one black ball in the urn, and the index of rounds starts from 2. Then after round  $n$ , there are exactly  $n$  balls in the urn. Let  $X_n$  be the number of black balls after round  $n$ , and  $Z_n = \frac{X_n}{n}$  is the ratio of black balls after round  $n$ . Clearly  $Z_2 = \frac{1}{2}$ . Then we have

$$\begin{aligned} \mathbb{E} \left[ Z_{n+1} \mid \bar{X}_{2,n} \right] &= \frac{1}{n+1} \mathbb{E} \left[ X_{n+1} \mid \bar{X}_{2,n} \right] \\ &= \frac{1}{n+1} (Z_n(X_n + 1) + (1 - Z_n)X_n) = \frac{Z_n + X_n}{n+1} = Z_n. \end{aligned}$$

That is,  $\{Z_n\}_{n \geq 2}$  is a martingale w.r.t.  $\{X_n\}_{n \geq 2}$ .

Example 5 shows that  $X_n$  does not have to be i.i.d..

# [AI2613 Lecture 7] Doob Martingale, Azuma-Hoeffding, McDiarmid

April 6, 2023

## 1 Hoeffding's Inequality

We introduced the following Hoeffding's inequality to bound the concentration for the sum of a sequence independent random variables.

**Theorem 1 (Hoeffding's Inequality)** Let  $X_1, \dots, X_n$  be independent random variables where each  $X_i \in [a_i, b_i]$  for certain  $a_i \leq b_i$  with probability 1. Let  $X = \sum_{i=1}^n X_i$  and  $\mu \triangleq \mathbf{E}[X] = \sum_{i=1}^n \mathbf{E}[X_i]$ , then

$$\Pr[|X - \mu| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

for all  $t \geq 0$ .

Before proving Theorem 1 in Section 3, we see a practical application of Hoeffding's inequality.

**Example 1 (Meal Delivery)** During the quarantine of our campus, the professors deliver meals for students using their private cars or trikes. Then a practical problem is how to estimate the amount of meals on a trike conveniently (See the [news](#)).

Assume there are  $n$  boxes of meal on the trike ( $n \geq 200$  and is unknown for us). Let  $X_i$  be the weight of the  $i$ -th box of meal. Assume that  $X_1, X_2, \dots, X_n$  are i.i.d. random variables, each  $X_i \in [250, 350]$  (unit: gram) and  $\mu = \mathbf{E}[X_i] = 300$ . Let  $S$  be the total weight of the meal boxes on the trike, that is,  $S = \sum_{i=1}^n X_i$ . We can weigh the meal boxes and use  $\hat{n} = \frac{S}{\mu}$  as an estimator for  $n$ . Now we compute how accurate this estimator is.

Let  $\delta \geq 0$  be a constant. By Hoeffding's inequality,

$$\Pr[|\hat{n} - n| > \delta n] = \Pr[|S - \mu n| > \delta \mu n] \leq 2 \exp\left\{-\frac{2\delta^2 \mu^2 n^2}{\sum_{i=1}^n (350 - 250)^2}\right\}. \quad (1)$$

Plugging  $\mu = 300$ ,  $\delta = 0.05$  and  $n \geq 200$  into Equation (1), by direct calculation, we have

$$\Pr[\hat{n} \in [0.95n, 1.05n]] \geq 1 - 2.4682 \times 10^{-4}.$$

## 2 Concentration on Martingale

We consider the balls-in-a-bag problem. There are  $g$  green balls and  $r$  red balls in a bag. These balls are the all same except for the color. We want to estimate the ratio  $\frac{r}{r+g}$  by drawing balls. There are two scenarios.

- Draw balls with replacement. Let  $X_i = \mathbf{1}[\text{the } i\text{-th ball is red}]$ . Let  $X = \sum_{i=1}^n X_i$ . Then clearly each  $X_i \sim \text{Ber}\left(\frac{r}{r+g}\right)$  and  $\mathbb{E}[X] = n \cdot \frac{r}{r+g}$ . Since all  $X_i$ 's are independent, we can directly apply Hoeffding's inequality and obtain

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{n}\right).$$

- Draw balls without replacement. Again we let  $Y_i = \mathbf{1}[\text{the } i\text{-th ball is red}]$ , then unlike the case of drawing with replacement, variables in  $\{Y_i\}$  are dependent. Let  $Y = \sum_{i=1}^n Y_i$ . We first calculate  $\mathbb{E}[Y]$ .

For every  $i \geq 1$ ,  $\mathbb{E}[Y_i]$  is the probability that the  $i$ -th draw is a red ball. Note that drawing without replacement is equivalent to first drawing a uniform permutation of  $r + g$  balls and drawing each ball one by one in that order. Therefore, the probability of  $Y_i = 1$  is  $\frac{r \cdot (r+g-1)!}{(r+g)!} = \frac{r}{r+g}$ . So we have  $\mathbb{E}[Y] = n \cdot \frac{r}{r+g}$ .

However, since  $\{Y_i\}$  are dependent, we cannot apply Hoeffding's inequality directly. This motivate us to generalize it by removing the requirement of independence.

### 2.1 Azuma-Hoeffding's Inequality

**Theorem 2 (Azuma-Hoeffding's Inequality)** *Let  $\{Z_n\}_{n \geq 0}$  is a martingale with respect to a filtration  $\{\mathcal{F}_n\}$ . If for every  $i \geq 1$ ,  $|Z_i - Z_{i-1}| \leq c_i$  with probability 1, then*

$$\Pr[|Z_n - Z_0| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

Azuma-Hoeffding's inequality generalizes Hoeffding's inequality by letting  $Z_n = \sum_{i=1}^n (X_i - \mathbb{E}[X_i])$  and  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ .

The proof of this theorem is in Section 3. The requirement of martingale in Theorem 2 seems to be even harder to satisfy than the requirement of independence. However, in many cases, we can construct a doob martingale to apply the Azuma-Hoeffding's inequality.

**Definition 3 (Doob Martingale, Doob Sequence)** *Let  $X_1, \dots, X_n$  be a sequence of (unnecessarily independent) random variables and  $f(\bar{X}_{1,n}) = f(X_1, \dots, X_n) \in \mathbb{R}$  be a function. For  $i \geq 0$ , Let  $Z_i \triangleq \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i}\right]$ . Then we call  $\{Z_n\}_{n \geq 0}$  a Doob martingale or a Doob sequence.*

It is easy to verify that  $\{Z_n\}_{n \geq 0}$  in Definition 3 is indeed a martingale w.r.t.  $\{X_n\}$  by seeing

$$\mathbb{E}\left[Z_i \mid \bar{X}_{1,i-1}\right] = \mathbb{E}\left[\mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i}\right] \mid \bar{X}_{1,i-1}\right] = \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}\right] = Z_{i-1}.$$

Let  $\mathcal{F} = \sigma(\bar{X}_{1,i})$ . We can see that  $Z_i$  is  $\mathcal{F}_i$  measurable by definition. Moreover, we know that  $Z_0 = \mathbf{E} \left[ f(\bar{X}_{1,n}) \right]$  and  $Z_n = f(\bar{X}_{1,n})$ .

Recall the balls-in-a-bag problem we discussed above. Define  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  by letting  $f(y_1, y_2, \dots, y_n) = \sum_{i=1}^n y_i$ . Then in the drawing without replacement scenario,  $Y = \sum_{i=1}^n Y_i = f(Y_1, Y_2, \dots, Y_n)$ . Now we construct the Doob martingale for  $f$ .

Let  $Z_i = \mathbf{E} \left[ f(\bar{Y}_{1,n}) \mid \bar{Y}_{1,i} \right]$ . We know that  $Z_0 = \mathbf{E} \left[ f(\bar{Y}_{1,n}) \right] = \mathbf{E} [Y] = n \cdot \frac{r}{r+g}$  and  $Z_n = f(\bar{Y}_{1,n})$ . In order to apply Azuma-Hoeffding, we need to bound the *width* of the martingale  $|Z_i - Z_{i-1}|$ .

By definition,

$$Z_i - Z_{i-1} = \mathbf{E} \left[ f(\bar{Y}_{1,n}) \mid \bar{Y}_{1,i} \right] - \mathbf{E} \left[ f(\bar{Y}_{1,n}) \mid \bar{Y}_{1,i-1} \right].$$

If we use  $S_i$  to denote the number of red balls among the first  $i$  balls, namely  $S_i = \sum_{j=1}^i Y_j$ , then

$$\mathbf{E} \left[ f(\bar{Y}_{1,n}) \mid \bar{Y}_{1,i} \right] = \mathbf{E} \left[ f(\bar{Y}_{1,n}) \mid S_i \right] = S_i + (n-i) \cdot \frac{r - S_i}{g + r - i}.$$

Therefore  $S_i = S_{i-1} + Y_i$  and

$$\begin{aligned} Z_i - Z_{i-1} &= \left( S_i + (n-i) \cdot \frac{r - S_i}{g + r - i} \right) - \left( S_{i-1} + (n-i+1) \cdot \frac{r - S_{i-1}}{g + r - i + 1} \right) \\ &= \frac{g + r - n}{g + r - i} \left( Y_i + \frac{S_{i-1} - r}{g + r - i + 1} \right). \end{aligned}$$

Note that  $r \geq S_{i-1}$  and  $g \geq (i-1) - S_{i-1}$ , we have

$$\begin{aligned} Z_i - Z_{i-1} &\leq \frac{g + r - n}{g + r - i} \left( 1 + \frac{S_{i-1} - r}{g + r - i + 1} \right) \leq \frac{g + r - n}{g + r - i} \leq 1, \\ Z_i - Z_{i-1} &\geq \frac{g + r - n}{g + r - i} \left( \frac{S_{i-1} - r}{g + r - i + 1} \right) \geq -\frac{g + r - n}{g + r - i} \geq -1. \end{aligned}$$

Therefore  $-1 \leq X_i \leq 1$  and we can apply Azuma-Hoeffding to  $Z_n - Z_0$  to obtain

$$\Pr [|Y - \mathbf{E} [Y]| \geq t] \leq 2 \exp \left( -\frac{t^2}{2n} \right).$$

## 2.2 McDiarmids Inequality

The Doob sequence we used in the balls-in-a-bag example is a very powerful and general tool to obtain concentration bounds. For a model defined by  $n$  random variables  $X_1, \dots, X_n$  and any quantity  $f(X_1, \dots, X_n)$  that we want to estimate, we can apply the Azuma-Hoeffding inequality to the Doob sequence of  $f$ . As shown in the previous example, the quality of the bound relies on the *width* of the martingale, that is, the magnitude of  $|Z_i - Z_{i-1}|$ . To determine the width of each  $|Z_i - Z_{i-1}|$  is relatively easy if the function  $f$  and the variables  $\{X_i\}_{1 \leq i \leq n}$  enjoy certain nice properties.

**Definition 4 (c-Lipschitz Function)** A function  $f(x_1, \dots, x_n)$  satisfies  $c$ -Lipschitz condition if

$$\forall i \in [n], \forall x_1, \dots, x_n, \forall y_i : |f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, y_i, \dots, x_n)| \leq c.$$

The McDiarmid's inequality is the application of Azuma-Hoeffding inequality to Lipschitz  $f$  and independent  $\{X_i\}$ .

**Theorem 5 (McDiarmid's Inequality)** Let  $f$  be a function on  $n$  variables satisfying  $c$ -Lipschitz condition and  $X_1, \dots, X_n$  be  $n$  independent variables. Then we have

$$\Pr[|f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)]| \geq t] \leq 2e^{-\frac{t^2}{nc^2}}.$$

*Proof.* We use  $f$  and  $\{X_i\}_{i \geq 1}$  to define a Doob martingale  $\{Z_i\}_{i \geq 1}$ :

$$\forall i : Z_i = \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i}\right].$$

Then

$$Z_i - Z_{i-1} = \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i}\right] - \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}\right].$$

Next we try to determine the width of  $Z_i - Z_{i-1}$ . Clearly

$$Z_i - Z_{i-1} \geq \inf_x \left\{ \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}, X_i = x\right] - \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}\right] \right\},$$

and

$$Z_i - Z_{i-1} \leq \sup_y \left\{ \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}, X_i = y\right] - \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}\right] \right\}.$$

The gap between the upper bound and the lower bound is

$$\sup_{x,y} \left\{ \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}, X_i = y\right] - \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bar{X}_{1,i-1}, X_i = x\right] \right\}.$$

For every  $x, y$  and  $\sigma_1, \dots, \sigma_{i-1}$ ,

$$\begin{aligned} & \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bigwedge_{1 \leq j \leq i-1} X_j = \sigma_j, X_i = y\right] - \mathbb{E}\left[f(\bar{X}_{1,n}) \mid \bigwedge_{1 \leq j \leq i-1} X_j = \sigma_j, X_i = x\right] \\ &= \sum_{\sigma_{i+1}, \dots, \sigma_n} \left( \Pr\left[\bigwedge_{i+1 \leq j \leq n} X_j = \sigma_j \mid \bigwedge_{1 \leq j \leq i-1} X_j = \sigma_j, X_i = y\right] \cdot f(\sigma_1, \dots, \sigma_{i-1}, y, \sigma_{i+1}, \dots, \sigma_n) \right. \\ & \quad \left. - \Pr\left[\bigwedge_{i+1 \leq j \leq n} X_j = \sigma_j \mid \bigwedge_{1 \leq j \leq i-1} X_j = \sigma_j, X_i = x\right] \cdot f(\sigma_1, \dots, \sigma_{i-1}, x, \sigma_{i+1}, \dots, \sigma_n) \right) \\ &\stackrel{(\heartsuit)}{=} \sum_{\sigma_{i+1}, \dots, \sigma_n} \Pr\left[\bigwedge_{i+1 \leq j \leq n} X_j = \sigma_j\right] \cdot (f(\sigma_1, \dots, \sigma_{i-1}, y, \sigma_{i+1}, \dots, \sigma_n) - f(\sigma_1, \dots, \sigma_{i-1}, x, \sigma_{i+1}, \dots, \sigma_n)) \\ &\stackrel{(\clubsuit)}{\leq} c. \end{aligned}$$

where  $(\heartsuit)$  uses independence of  $\{X_i\}$  and  $(\clubsuit)$  uses the  $c$ -Lipschitz property of  $f$ .



Applying Azuma-Hoeffding, we have

$$\Pr [|Z_n - Z_0| \geq t] = \Pr [|f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)]| \geq t] \leq 2e^{-\frac{2t^2}{nc^2}}.$$

□

Then we examine two applications of McDiarmid's inequality.

**Example 2 (Pattern matching)** Let  $P \in \{0, 1\}^k$  be a fixed string. For a random string  $X \in \{0, 1\}^n$ , what is the expected number of occurrences of  $P$  in  $X$ ?

The expectation of occurrence times can be easily calculated using the linearity of expectations. We define  $n$  independent random variables  $X_1, \dots, X_n$ , where  $X_i$  denotes  $i$ -th character of  $X$ . Let  $Y = f(X_1, \dots, X_n)$  be the number of occurrences of  $P$  in  $X$ . Note that there are at most  $n - k + 1$  occurrences of  $P$  in  $X$ , and we can enumerate the first position of each occurrence. By the linearity of expectation, we have

$$\mathbb{E}[f] = \frac{n - k + 1}{2^k}.$$

We can then use McDiarmid's inequality to show that  $f$  is well-concentrated. To see this, we note that variables in  $\{X_i\}$  are independent and the function  $f$  is  $k$ -Lipschitz: If we change one bit of  $X$ , the number of occurrences changes at most  $k$ .

Therefore

$$\Pr [|Z_n - Z_0| \geq t] = \Pr [|f - \mathbb{E}[f]| \geq t] \leq 2e^{-\frac{2t^2}{nk^2}}.$$

Another application of McDiarmid's Inequality is to establish the concentration of chromatic number for Erdős-Rényi random graphs  $\mathcal{G}(n, p)$ .

**Example 3 (Chromatic Number of  $\mathcal{G}(n, p)$ )** Recall the notation  $\mathcal{G}(n, p)$  specifies a distribution over all undirected simple graphs with  $n$  vertices. In the model, each of the  $\binom{n}{2}$  possible edges exists with probability  $p$  independently.

For a graph  $G \sim \mathcal{G}(n, p)$ , we use  $\chi(G)$  to denote its chromatic number, the minimum number  $q$  so that  $G$  can be properly colored using  $q$  colors. There are different ways to represent  $G$  using random variables.

The most natural way is to introduce a variable  $X_e$  for every pair of vertices  $e = \{u, v\} \subseteq V$  where  $X_e = \mathbf{1}[\text{the edge } e \text{ exists in } G]$ . Then  $\{X_e\}$  are independent and the chromatic number can be written as a function  $\chi(X_{e_1}, X_{e_2}, \dots, X_{e_{\binom{n}{2}}})$ . It is easy to see that  $\chi$  is 1-Lipschitz as removing to adding one edge can only change the chromatic number by at most one. So by McDiarmid's inequality, we have

$$\Pr [|\chi - \mathbb{E}[\chi]| \geq t] \leq 2e^{-2t^2 \binom{n}{2}^{-1}}.$$

However, this bound is not satisfactory as we need to set  $t = \Theta(n)$  in order to upper bound the RHS by a constant.

We can encode the graph  $G$  in a more efficient way while reserving the Lipschitz and the independence property. Suppose the vertex set of  $G$  is  $\{v_1, \dots, v_n\}$ . We define  $n$  random variables  $Y_1, \dots, Y_n$ , where  $Y_i$  encodes the edges between  $v_i$  and  $\{v_1, \dots, v_{i-1}\}$ . Once  $Y_1, \dots, Y_n$  are given, the graph is known, so the chromatic number can be written as a function  $\chi(Y_1, \dots, Y_n)$ . Since  $Y_i$  only involves the connections between  $v_i$  and  $v_1, \dots, v_{i-1}$ , the  $n$  variables are independent.

It is also easy to see that if  $X_i$  changes, the chromatic number changes at most one. Hence  $\chi$  is 1-Lipschitz as well. Applying McDiarmid's inequality we have

$$\Pr[|\chi - \mathbb{E}[\chi]| \geq t] \leq 2e^{-\frac{2t^2}{n}}.$$

In this way, we only need  $t = \Theta(\sqrt{n})$  to bound the RHS.

### 3 Proof

#### 3.1 Proof of Theorem 1

First, we prove the following Hoeffding's lemma which will be the main technical ingredient to prove the inequality.

**Lemma 6** Let  $X$  be a random variable with  $\mathbb{E}[X] = 0$  and  $X \in [a, b]$ . Then it holds that

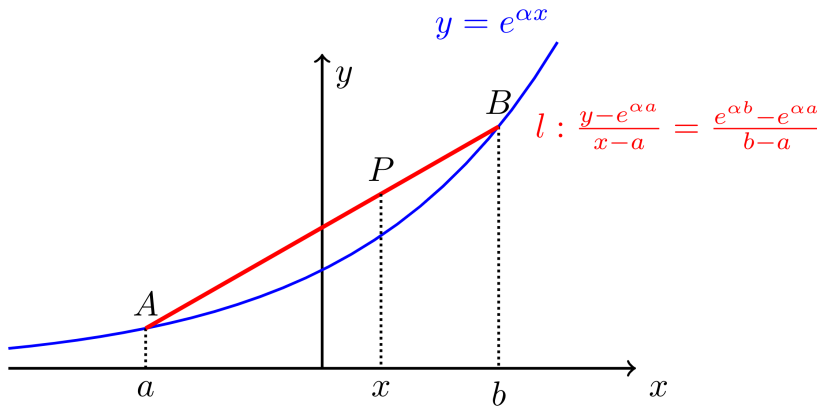
$$\mathbb{E}[e^{\alpha X}] \leq \exp\left(\frac{\alpha^2(b-a)^2}{8}\right) \text{ for all } \alpha \in \mathbb{R}.$$

*Proof.*

We first find a linear function to upper bound  $e^{\alpha x}$  so that we could apply the linearity of expectation to bound  $\mathbb{E}[e^{\alpha X}]$ . By the convexity of the exponential function and as illustrated in the figure below, we have

$$e^{\alpha x} \leq \frac{e^{\alpha b} - e^{\alpha a}}{b - a}(x - a) + e^{\alpha a}, \text{ for all } a \leq x \leq b.$$

Thus,



$$\begin{aligned}
\mathbf{E}[e^{\alpha x}] &\leq \frac{e^{\alpha b} - e^{\alpha a}}{b - a}(-a) + e^{\alpha a} = \frac{-a}{b - a}e^{\alpha b} + \frac{b}{b - a}e^{\alpha a} \\
&= e^{\alpha a} \left( \frac{b}{b - a} - \frac{a}{b - a}e^{\alpha(b-a)} \right) \\
&= e^{-\theta t} (1 - \theta + \theta e^t) \quad (\theta = -\frac{a}{b-a}, t = \alpha(b-a)) \\
&\triangleq e^{g(t)},
\end{aligned}$$

where

$$g(t) = -\theta t + \log(1 - \theta + \theta e^t).$$

By Taylor's theorem, for every real  $t$  there exists a  $\delta$  between 0 and  $t$  such that,

$$g(t) = g(0) + tg'(0) + \frac{1}{2}g''(\delta)t^2$$

Note that,

$$\begin{aligned}
g(0) &= 0; \\
g'(0) &= -\theta + \frac{\theta e^t}{1 - \theta + \theta e^t} \Big|_{t=0} \\
&= 0; \\
g''(\delta) &= \frac{\theta e^t (1 - \theta + \theta e^t) - \theta e^t}{(1 - \theta + \theta e^t)^2} \\
&= \frac{(1 - \theta)\theta e^t}{(1 - \theta + \theta e^t)^2} \\
&= \frac{(1 - \theta)\theta}{\theta^2 z + 2(1 - \theta)\theta + \frac{(1 - \theta)^2}{z}} \quad (z = e^t) \\
&\leq \frac{(1 - \theta)\theta}{2\theta(1 - \theta) + 2(1 - \theta)\theta} \quad (z > 0) \\
&= \frac{1}{4}.
\end{aligned}$$

Thus

$$g(t) \leq 0 + t \cdot 0 + \frac{1}{2}t^2 \cdot \frac{1}{4} = \frac{1}{8}t^2 = \frac{1}{8}\alpha^2(b-a)^2.$$

Therefore,  $\mathbf{E}[e^{\alpha x}] \leq \exp\left(\frac{\alpha^2(b-a)^2}{8}\right)$  holds.  $\square$

Armed with Hoeffding's lemma, it is routine to prove Hoeffding's inequality.

*Proof.* [Proof of Theorem 1]

First note that we can assume  $\mathbf{E}[X_i] = 0$  and therefore  $\mu = 0$  (if not so, replace  $X_i$  by  $X_i - \mathbf{E}[X_i]$ ). By symmetry, we only need to prove that  $\Pr[X \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$ . Since

$$\Pr[X \geq t] \stackrel{\alpha > 0}{=} \Pr[e^{\alpha X} \geq e^{\alpha t}] \leq \frac{\mathbf{E}[e^{\alpha X}]}{e^{\alpha t}}$$

and

$$\mathbf{E}[e^{\alpha X}] = \mathbf{E}\left[e^{\alpha \sum_{i=1}^n X_i}\right] = \prod_{i=1}^n \mathbf{E}[e^{\alpha X_i}],$$

applying Hoeffding's lemma for each  $\mathbf{E} [e^{\alpha X_i}]$  yields

$$\mathbf{E} [e^{\alpha X_i}] \leq \exp \left( \frac{\alpha^2 (b_i - a_i)^2}{8} \right).$$

Let  $\alpha = \frac{4t}{\sum_{i=1}^n (b_i - a_i)^2}$ , we have,

$$\begin{aligned} \Pr [X \geq t] &\leq \frac{\prod_{i=1}^n \mathbf{E} [e^{\alpha X_i}]}{e^{\alpha t}} \leq \exp \left( -\alpha t + \frac{\alpha^2}{8} \sum_{i=1}^n (b_i - a_i)^2 \right) \\ &= \exp \left( -\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right). \end{aligned}$$

□

### 3.2 Proof of Theorem 2

Now we will sketch a proof of the Azuma-Hoeffding, which is quite similar to our proof of the Hoeffding inequality.

*Proof.* [Proof of Theorem 2]

Recall when we were trying to prove the Hoeffding inequality, the most difficult part is to estimate the term

$$\mathbf{E} [e^{\alpha Z_n}] = e^{\alpha Z_0} \cdot \mathbf{E} \left[ \prod_{i=1}^n e^{\alpha (Z_i - Z_{i-1})} \right].$$

In the case of Azuma-Hoeffding, we can use the property of martingales instead of independence to obtain a bound of this term. To see this, we have

$$\begin{aligned} \mathbf{E} \left[ \prod_{i=1}^n e^{\alpha Z_i - Z_{i-1}} \right] &= \mathbf{E} \left[ \mathbf{E} \left[ \prod_{i=1}^n e^{\alpha Z_i - Z_{i-1}} \middle| \mathcal{F}_{n-1} \right] \right] \\ &= \mathbf{E} \left[ \prod_{i=1}^{n-1} e^{\alpha Z_i - Z_{i-1}} \mathbf{E} [e^{\alpha Z_n - Z_{n-1}} \mid \mathcal{F}_{n-1}] \right]. \end{aligned}$$

The bounds then follows by an induction argument and a conditional expectation version of Hoeffding lemma:

$$\mathbf{E} [e^{\alpha (Z_n - Z_{n-1})} \mid \mathcal{F}_{n-1}] \leq e^{-\frac{\alpha c_i^2}{8}}.$$

The proof is almost the same as our proof of Hoeffding lemma in the last lecture. □

# [AI2613 Lecture 8] Optional Stopping Theorem

April 17, 2023

## 1 Stopping Time

Suppose  $Z_0, Z_1, \dots, Z_n, \dots$  is a martingale with respect a certain filtration  $\{\mathcal{F}_t\}$ . We know that for any  $t$ ,  $E[Z_t] = E[Z_0]$ . However, does  $E[Z_\tau] = E[Z_0]$  still hold if  $\tau$  is a random variable?

Consider the following gambling strategy in a fair game. At the first round, the gambler bet \$1. Then he simply double his stake until he wins

Let  $\tau$  be the first time he wins. Then expected money he win at time  $\tau$  is 1, which is not equal to 0, his initial money. In order to understand the phenomenon, let us first formally introduce *stopping time*.

**Definition 1 (Stopping Time)** Let  $\tau \in \mathbb{N} \cup \{\infty\}$  be a random variable. We say  $\tau$  is a stopping time if for all  $t \geq 0$ , the event " $\tau \leq t$ " is  $\mathcal{F}_t$ -measurable.

For example, the first time that a gambler wins five games in a row is a stopping time, since for a given  $t$ , this can be determined by looking at the outcomes of all the previous games, and therefore the time is  $\mathcal{F}_t$ -measurable. However, the *last* time the gambler wins five games in a row is *not* a stopping time, since determining whether the time is  $t$  cannot be done without knowing  $X_{t+1}, X_{t+2}, \dots$

### 1.1 Optional Stopping Theorem(OST)

The optional stopping theorem provides sufficient condition for  $E[Z_\tau] = E[Z_0]$  to hold.

**Theorem 2 (Optional Stopping Theorem)** Let  $\{X_t\}_{t \geq 0}$  be a martingale and  $\tau$  be a stopping time with respect to  $\{\mathcal{F}_t\}_{t \geq 0}$ . Then  $E[X_\tau] = E[X_0]$  if at least one of the following conditions holds:

1.  $\tau$  is bounded almost surely, that is,  $\exists n \in \mathbb{N}$  such that  $\Pr[\tau \leq n] = 1$ ;
2.  $\Pr[\tau < \infty] = 1$ , and there is a finite  $M$  such that  $|X_t| \leq M$  for all  $t < \tau$ ;
3.  $E[\tau] < \infty$ , and there is a constant  $c$  such that  $E[|X_{t+1} - X_t| \mid \mathcal{F}_t] \leq c$  for all  $t < \tau$ .

We will prove the theorem next time. Let us look back at the boy-or-girl example mentioned in the first class.

**Example 1 (Boy or Girl)** Suppose there is a country in which people only want boys. What is the sex ratio of the country in the following three scenarios?

- Each family continues to have children until they have a boy.

The strategy was called [martingale](#)!

- If  $\tau = 1$ , he wins 1 dollar.
- If  $\tau = 2$ , he wins  $-1 + 2 = 1$  dollar.
- If  $\tau = 3$ , he wins  $-1 - 2 + 4 = 1$  dollar.
- ...

- Each family continues to have children until there are more boys.
- Each family continues to have children until there are more boys or there are 10 children.

We can model the problem as a random walk. Suppose there is a man walking randomly on a one-dimensional axis. Let  $\{X_t\}_{t \geq 0}$  be the positions of the man at each time where  $X_t$  stands for the number of boys minus the number of girls in the first  $t$  children of a family. Starting at  $X_0 = 0$ , at time 0, the man takes a step  $c_t \in_{\mathbb{R}} \{-1, 1\}$  and reach  $X_{t+1}$ , i.e.,  $X_{t+1} = X_t + c_t$ . It is easy to verify that  $\{X_t\}_{t \geq 0}$  is a martingale. The three scenarios mentioned correspond to the following three different definitions of a stopping time  $\tau$ . The identity  $\mathbb{E}[X_\tau] = \mathbb{E}[X_0]$  means that the sex ratio is balanced. We will check respectively whether it is the case using OST.

- Let  $\tau$  be the first time  $t$  such that  $c_t = 1$ . Then  $\mathbb{E}[\tau] < \infty$  since by definition  $\tau \sim \text{Geom}(\frac{1}{2})$ , and  $|X_{t+1} - X_t| \leq 1$  for all  $t < \tau$ . Therefore from the 3rd condition of OST we have  $\mathbb{E}[X_\tau] = \mathbb{E}[X_0] = 0$ . In other words, if the man stops at the first time of  $c_t = 1$ , then the expected final position is 0.
- Let  $\tau$  be the first time  $t$  such that  $X_t = 1$ , then of course  $\mathbb{E}[X_\tau] = 1 \neq \mathbb{E}[X_0]$ . This process is called “1-d random walk with one absorbing barrier” and it is well-known that  $\mathbb{E}[\tau] = \infty$ . No condition in OST is satisfied.
- Let  $\tau$  be the minimum between 10 and the first time  $t$  such that  $X_t = 1$ . In this case,  $\tau$  is at most 10, which satisfies the first condition of OST. Therefore we have  $\mathbb{E}[X_\tau] = \mathbb{E}[X_0] = 0$ .

The property  $\mathbb{E}[\tau] = \infty$  of the random work is called “null recurrent”. You can find more on this from my lecture on stochastic processes.

## 2 Applications of OST

### 2.1 Doob's martingale inequality

With OST, we can obtain concentration property of the maximum element in a sequence of random variables.

**Claim 3** Let  $\{X_t\}_{t \geq 0}$  be a martingale with respect to itself where  $X_t \geq 0$  for every  $t$ . Prove that for every  $n \in \mathbb{N}$ ,

$$\Pr \left[ \max_{0 \leq t \leq n} X_t \geq \alpha \right] \leq \frac{\mathbb{E}[X_0]}{\alpha}.$$

*Proof.* We define a stopping time  $\tau$  when the first element that is greater than  $\alpha$  occurs, and otherwise set  $\tau = n$ . Formally, define

$$\tau \triangleq \min \left( n, \min_{t \leq n} \{t \mid X_t \geq \alpha\} \right).$$

By definition of  $\tau$ , we have

$$\Pr \left[ \max_{0 \leq t \leq n} X_t \geq \alpha \right] = \Pr [X_\tau \geq \alpha].$$

Since  $\tau$  is bounded, we apply Optional Stopping Theorem to obtain that  $E[X_\tau] = E[X_0]$ . Therefore, by Markov's Inequality,

$$\Pr \left[ \max_{0 \leq t \leq n} X_t \geq \alpha \right] = \Pr [X_\tau \geq \alpha] \leq \frac{E[X_\tau]}{\alpha} = \frac{E[X_0]}{\alpha}$$

□

## 2.2 One-dimensional Random Walk with Two Absorbing Barriers

We consider another problem in one-dimensional random walk. Let  $a, b > 0$  be two integers. A man starts the random walk from 0 and stops when he arrives at  $-a$  or  $b$ . Let  $\tau$  be the time when the man first reaches  $-a$  or  $b$ , i.e., the first time  $t$  that  $X_t = -a$  or  $X_t = b$ . The model is called “one-dimensional random walk with two absorbing barriers”. We want to compute the expected value of  $E[\tau]$ , that is, the average stopping time of the walk.

We want to construct a martingale  $\{Y_t\}_{t \geq 0}$  such that OST can be applied to  $\{Y_t\}_{t \geq 0}$  and  $\tau$  and thereby we can derive an equality related to  $E[\tau]$ . Before calculating  $E[\tau]$ , we first determine  $\Pr[X_\tau = -a]$ , the probability that the man stops at position  $-a$ . Let  $P_a \triangleq \Pr[X_\tau = -a]$ . We want to apply OST to show  $E[X_\tau] = E[X_0]$ . Therefore, we verify that some of conditions in OST is satisfied.

In a time period of length  $T = a + b$ , if the man walks towards the same direction, he must have stopped, either at  $-a$  or  $b$ , which happens with probability  $2^{-(a+b)}$ . Therefore, if we divide the time into consecutive periods in this manner, in expected finite time, we can meet some period when the event happened. Hence,  $E[\tau] < \infty$ . Moreover, we clearly have  $E[|X_{t+1} - X_t| | \mathcal{F}_t] < 1$  for every  $0 \leq t < \tau$ , so the third condition of OST holds, which implies that  $E[X_\tau] = E[X_0]$ . On the other hand, we have  $E[X_\tau] = P_a \cdot (-a) + (1 - P_a) \cdot b$ . These two equalities give  $P_a = \frac{b}{a+b}$ .

Then for all  $t \geq 0$ , we define a new random variable  $Y_t \triangleq X_t^2 - t$  which involves the time  $t$ . The following fact is easy to verify by definition.

**Claim 4**  $\{Y_t\}_{t \geq 0}$  is a martingale.

*Proof.* First we have

$$\begin{aligned} E[Y_{t+1} | \mathcal{F}_t] &= E[X_{t+1}^2 - (t+1) | \mathcal{F}_t] \\ &= E[(X_t + c_t)^2 - (t+1) | \mathcal{F}_t] \\ &= E[X_t^2 | \mathcal{F}_t] + 2E[X_t c_t | \mathcal{F}_t] + E[c_t^2 | \mathcal{F}_t] - (t+1). \end{aligned}$$

Since  $X_t$  is  $\mathcal{F}_t$ -measurable,  $E[c_t | \mathcal{F}_t] = 0$  and  $E[c_t^2 | \mathcal{F}_t] = 1$ , we can further derive that

$$E[Y_{t+1} | \mathcal{F}_t] = X_t^2 + 0 + 1 - (t+1) = X_t^2 - t = Y_t.$$

Hence  $\{Y_t\}_{t \geq 0}$  is a martingale.

□

We've discussed one-dimensional random walk with one absorbing barrier before

Sometimes one can use OST in a reverse way. Consider the random walk with only one barrier at  $-a$ . The fact that  $E[\tau] = \infty$  can be proved in the following way (due to Biaoshuai Tao): If  $E[\tau] < \infty$ , then by (cond 3 of) OST,  $E[X_\tau] = E[X_0] = 0$ . On the other hand, we know  $X_\tau = -a \neq 0$ . Therefore it must be that  $E[\tau] = \infty$ .

Note that  $X_t \in [-a, b]$  for all  $t \geq 0$ . Thus  $|Y_{t+1} - Y_t| = |X_{t+1}^2 - (t+1) - X_t^2 + t| = |X_{t+1}^2 - X_t^2 - 1|$  is bounded by some constant. We can apply OST again to obtain  $E[Y_\tau] = E[Y_0] = 0$ . On the other hand, we have  $E[Y_\tau] = E[X_\tau^2] - E[\tau]$  by definition, and thus

$$E[\tau] = E[X_\tau^2] = a^2 P_a + b^2 (1 - P_a) = a^2 \cdot \frac{b}{a+b} + b^2 \cdot \frac{a}{a+b} = ab.$$

### 2.3 Pattern Matching

Suppose that there is a  $\{H, T\}$ -string  $P$  of length  $\ell$  (H for “head” and T for “tail”). We flip a coin consecutively until the last  $\ell$  results form exactly the same string as  $P$ . How many times do we flip the coin?

Note that if we flip the coin  $N$  times and observe the string  $S$  consisting of  $N$  results. No matter which pattern we choose, by the linearity of expectation, the expected number of occurrence<sup>1</sup> is

$$E[\text{\# of occurrence of } P \text{ in } S] = \sum_{i=1}^{n-\ell+1} E[\mathbb{1}[S_{i,i+1,\dots,i+\ell-1} = P]] = (n - \ell + 1) \cdot 2^{-\ell}.$$

<sup>1</sup> That means the expected number of substrings exactly the same as  $P$  in the resulting string  $S$ .

However, if we would like to compute the first time that pattern  $P$  occurs, the pattern itself has an impact on the expected time. Intuitively, let's consider two patterns HT and HH. Assume that the first flipping result is H. Then we consider what happens if the second result fails. Suppose that the desired pattern is HT and H appears. Although we fail, we obtain an H. However, if the desired pattern is HH and the second flipping result is T, then we obtain nothing and the first two flips are a waste. So we should believe that the expected times of the first occurrence of HT is smaller than HH.

We now use the optional stopping theorem to solve this problem. Let  $P = p_1 p_2 \dots p_\ell$ . For every  $n \geq 0$ , assume that before  $n + 1$ -th flipping there is a new gambler  $G_{n+1}$  coming with 1 unit of money to bet that the following  $\ell$  result (i.e., the  $n + 1$ -th to  $n + \ell$ -th results) are exactly the same as  $P$ . At the  $n + k$ -th flipping,  $G_{n+1}$  will bet that the result is  $p_k$  by an all in strategy, that is, if the  $n + k$ -th result is  $p_k$  then  $G_{n+1}$  will have twice as much money as before; otherwise they will lose all. Suppose that the pattern  $P = \text{HTHTH}$  and the flipping results are  $\text{HTHHTH}$ . The following table shows the total money of each gambler after flipping.

Let  $X_t$  be the result of  $t$ -th flipping,  $M_i(t)$  denote the money that  $G_i$  has after  $t$ -th flipping, and  $Z_t \triangleq \sum_{i=1}^t (M_i(t) - 1)$  be the total income of all gamblers after  $t$ -th flipping. It is easy to verify that  $\{M_i(t)\}_{t \geq 0}$  is a martingale with respect to  $\{X_t\}$  since

$$E[M_i(t+1) \mid \bar{X}_{0,t}] = \frac{1}{2} \cdot 2M_i(t) + \frac{1}{2} \cdot 0 = M_i(t).$$

Then by the linearity of expectation we conclude that  $\{Z_t\}_{t \geq 0}$  is a martingale with respect to the flipping results  $\{X_t\}$  since  $E[M_i(t)] = 1$ . Let



Gambler	H	T	H	H	T	H	T	H	Money	
1	H	T	H	T					0	$1 \rightarrow 2 \rightarrow 4 \rightarrow 8 \rightarrow 0$
2		H							0	$1 \rightarrow 0$
3			H	T					0	$1 \rightarrow 2 \rightarrow 0$
4				H	T	H	T	H	32	$1 \rightarrow 2 \rightarrow 4 \rightarrow 8 \rightarrow 16 \rightarrow 32$
5					H				0	$1 \rightarrow 0$
2						H	T	H	8	$1 \rightarrow 2 \rightarrow 4 \rightarrow 8$
5							H		0	$1 \rightarrow 0$
5								H	2	$1 \rightarrow 2$

$\tau$  be the stopping time defined by the first time that some gambler wins, namely, the first time that P occurs in the flipping results. Applying Condition 2 of OST we obtain that  $\mathbf{E}[Z_\tau] = \mathbf{E}[Z_0] = 0$ . Sequentially we have  $\mathbf{E}[\sum_{i=1}^{\tau} M_i(\tau) - \tau] = 0$  and  $\mathbf{E}[\tau] = \sum_{i=1}^{\tau} \mathbf{E}[M_i(\tau)]$ .

Note that  $M_i(t) = 0$  for  $i \leq \tau - \ell$  and  $M_i(t) = 2^{\tau-i+1} \chi_{\tau-i+1}$  for  $i > \tau - \ell$  where  $\chi_j$  is defined by

$$\chi_j = \mathbb{1}[p_1 p_2 \dots p_j = p_{\ell-j+1} \dots p_{\ell-1} p_\ell].$$

Hence,

$$\mathbf{E}[\tau] = \sum_{i=\tau-\ell+1}^{\tau} \mathbf{E}[M_i(\tau)] = \sum_{i=1}^{\ell} 2^i \chi_i.$$

Recall the example of HH and HT. If P is HH,  $\mathbf{E}[\tau] = 2 + 4 = 6$ . If P is HT,  $\mathbf{E}[\tau] = 4$ . This confirms our hypothesis that  $\mathbf{E}[\tau]$  for HH is larger than  $\mathbf{E}[\tau]$  for HT.

## 2.4 Wald's Equation

In practice, we often need to analyze the (expected) running time of following procedure where both *cond* and *compute()* are random.

```
while cond do
  compute();
end while
```

Assume the  $i$ -th call to *compute()* costs  $X_i$  time and the algorithm terminates after  $T$  iterations. Then the total running time is  $N \triangleq \sum_{i=1}^T X_i$ . Suppose  $X_i$ s are independently and identically distributed as a random variable  $X$ . The Wald's equation gives a formula for  $\mathbf{E}[N]$ .

**Theorem 5 (Wald's Equation)** *If we have*

- $X_1, X_2, \dots$  are non-negative, independent, identically distributed random variables with the same distribution as  $X$ .
- $T$  is a stopping time for  $X_1, X_2, \dots$ .

- $\mathbf{E}[T], \mathbf{E}[X] < \infty$ ,

then

$$\mathbf{E}\left[\sum_{i=1}^T X_i\right] = \mathbf{E}[T] \cdot \mathbf{E}[X].$$

*Proof.* For  $i \geq 1$ , let  $Z_i := \sum_{j=1}^i (X_j - \mathbf{E}[X])$ . Clearly the sequence  $Z_1, Z_2, \dots$  is a martingale with respect to  $X_1, X_2, \dots$  and  $\mathbf{E}[Z_1] = 0$ . And we have

$$\begin{aligned} \mathbf{E}[|Z_{i+1} - Z_i| \mid \mathcal{F}_i] &= \mathbf{E}[|X_{i+1} - \mathbf{E}[X]| \mid \mathcal{F}_i] \\ &\leq \mathbf{E}[X_{i+1} + \mathbf{E}[X] \mid \mathcal{F}_i] \\ &\leq 2\mathbf{E}[X]. \end{aligned}$$

We know that  $\mathbf{E}[T], \mathbf{E}[X] < \infty$ , and therefore applying OST derives  $\mathbf{E}[Z_T] = \mathbf{E}[Z_1] = 0$ . Then

$$\begin{aligned} \mathbf{E}[Z_T] &= \mathbf{E}\left[\sum_{j=1}^T (X_j - \mathbf{E}[X])\right] \\ &= \mathbf{E}\left[\sum_{i=1}^T X_i - T\mathbf{E}[X]\right] \\ &= \mathbf{E}\left[\sum_{i=1}^T X_i\right] - \mathbf{E}[T] \mathbf{E}[X] = 0. \end{aligned}$$

□

*An Application of Wald's Equation: A Routing Problem* Let us consider an application of Wald's equation. There are  $n$  senders and one receiver. In each round, each sender sends a packet to the receiver with probability  $\frac{1}{n}$ . Since all senders share the same channel, if there are multiple packets sent at the same time, all of them will fail. The question is, on average, how many rounds are required so that each sender can successfully send at least one packet.

We let  $X_i$  be the variable indicating how long the receiver needs to get another packet after he has received  $i - 1$  ones (counting packets from repeated sender). And let  $T$  be the number of packets received when first time the receiver receives at least one packet from each sender. The quantity we are interested in is

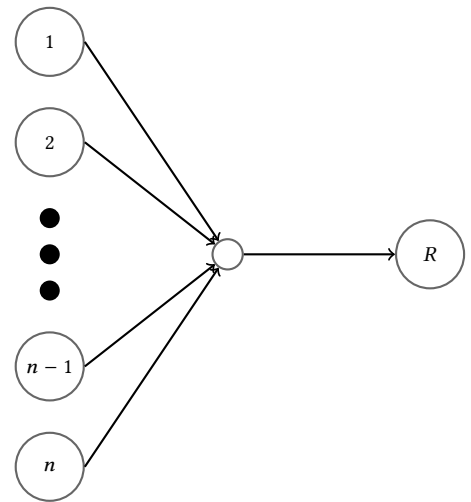
$$N \triangleq \sum_{i=1}^T X_i.$$

Clearly  $X_1, X_2, \dots$  are independently and identically distributed, and  $\mathbf{E}[T]$  is finite. Therefore  $\mathbf{E}[N] = \mathbf{E}[T] \cdot \mathbf{E}[X_1]$  by Wald's equation.

Note that by the definition,  $T$  is the number of coupons in the coupon collector's problem we met before. So  $\mathbf{E}[T] = nH_n = \Theta(n \log n)$ .

On the otherhand,  $X_1 \sim \text{Geom}(p)$  with

$$p = n \cdot \frac{1}{n} \left(1 - \frac{1}{n}\right)^{n-1} \approx e^{-1}$$



which implies  $E[X_1] = e$ . Therefore,

$$E[N] = E[T] \cdot E[X_1] \approx enH_n.$$

### 3 Proof of Optional Stopping Theorem

Let us restate the theorem.

**Theorem 6 (Optional Stopping Theorem)** *Let  $\{X_t\}_{t \geq 0}$  be a martingale and  $\tau$  be a stopping time with respect to  $\{\mathcal{F}_t\}_{t \geq 0}$ . Then  $E[X_\tau] = E[X_0]$  if at least one of the following conditions holds:*

1.  $\tau$  is bounded almost surely, that is,  $\exists n \in \mathbb{N}$  such that  $\Pr[\tau \leq n] = 1$ ;
2.  $\Pr[\tau < \infty] = 1$ , and there is a finite  $M$  such that  $|X_t| \leq M$  for all  $t < \tau$ ;
3.  $E[\tau] < \infty$ , and there is a constant  $c$  such that  $E[|X_{t+1} - X_t| \mid \mathcal{F}_t] \leq c$  for all  $t < \tau$ .

*Proof.* It is obvious that for every  $n \in \mathbb{N}$ ,  $E[X_n] = E[X_0]$ . So first we show that for every  $n \in \mathbb{N}$ ,  $E[X_{\min\{n, \tau\}}] = E[X_0]$ . Define  $Z_n \triangleq X_{\min\{n, \tau\}} = X_0 + \sum_{i=0}^{n-1} (X_{i+1} - X_i) \mathbb{1}[\tau > i]$ . We verify that  $\{Z_n\}_{n \geq 0}$  is a martingale. By definition

$$\begin{aligned} E[Z_{n+1} \mid \mathcal{F}_n] &= E[Z_n + (X_{n+1} - X_n) \mathbb{1}[\tau > n] \mid \mathcal{F}_n] \\ &= Z_n + \mathbb{1}[\tau > n] (E[X_{n+1} \mid \mathcal{F}_n] - X_n) \\ &= Z_n. \end{aligned}$$

So we have  $E[X_{\min\{n, \tau\}}] = E[Z_n] = E[Z_0] = E[X_0]$ .

Therefore, this motivates us to decompose  $X_\tau$  into two terms:

$$\forall n \in \mathbb{N}, X_\tau = X_{\min\{n, \tau\}} + \mathbb{1}[\tau > n] \cdot (X_\tau - X_n).$$

Taking expectation and letting  $n$  tend to infinity, we obtain

$$E[X_\tau] = E[X_0] + \lim_{n \rightarrow \infty} E[\mathbb{1}[\tau > n] \cdot (X_\tau - X_n)].$$

Therefore, we only need to verify that each of the three conditions in the statement guarantee  $\lim_{n \rightarrow \infty} E[\mathbb{1}[\tau > n] \cdot (X_\tau - X_n)] = 0$ .

1. If  $\tau$  is bounded almost surely, then clearly  $E[\mathbb{1}[\tau > n] \cdot (X_\tau - X_n)] = 0$  for sufficiently large  $n$ .
2. In this case,

$$\begin{aligned} E[\mathbb{1}[\tau > n] \cdot (X_\tau - X_n)] &\leq E[\mathbb{1}[\tau > n] \cdot (|X_\tau| + |X_n|)] \\ &\leq 2M \cdot E[\mathbb{1}[\tau > n]] \\ &= 2M \cdot \Pr[\tau > n] \rightarrow 0 \text{ as } n \rightarrow \infty. \end{aligned}$$

3. In order to apply our bounds on the gap between consecutive  $X_t$ , we write

$$\begin{aligned} \mathbb{1}[\tau > n] \cdot (X_\tau - X_n) &= \sum_{t=n}^{\tau-1} (X_{t+1} - X_t) \\ &\leq \sum_{t=n}^{\tau-1} |X_{t+1} - X_t| \\ &= \sum_{t=n}^{\infty} |X_{t+1} - X_t| \cdot \mathbb{1}[\tau > t]. \end{aligned}$$

Taking expectation on both sides, we have

$$\begin{aligned} \mathbb{E}[\mathbb{1}[\tau > n] \cdot (X_\tau - X_n)] &\leq \mathbb{E}\left[\sum_{t=n}^{\infty} |X_{t+1} - X_t| \cdot \mathbb{1}[\tau > t]\right] \\ &= \sum_{t=n}^{\infty} \mathbb{E}[|X_{t+1} - X_t| \cdot \mathbb{1}[\tau > t]] \\ &= \sum_{t=n}^{\infty} \mathbb{E}[\mathbb{E}[|X_{t+1} - X_t| \cdot \mathbb{1}[\tau > t] \mid \mathcal{F}_t]] \\ &= \sum_{t=n}^{\infty} \mathbb{E}[\mathbb{E}[|X_{t+1} - X_t| \mid \mathcal{F}_t] \cdot \mathbb{1}[\tau > t]] \\ &\leq \sum_{t=n}^{\infty} c \cdot \Pr[\tau > t], \end{aligned}$$

where the first equality follows from the monotone convergence theorem.

On the other hand, we know  $\mathbb{E}[\tau] = \sum_{t=0}^{\infty} \Pr[\tau > t] < \infty$ . Therefore, the tail of this sequence,  $\sum_{t=n}^{\infty} \Pr[\tau > t] \rightarrow 0$  as  $n \rightarrow \infty$ .

□

# [AI2613 Lecture 9] Poisson Distribution, Poisson Process

May 4, 2023

## 1 Poisson Distribution

**Example 1** *Let's consider a scenario where there is a restaurant that has had 100, 120, 80, 75, and 110 customers in the past five days. To ensure that they have the right amount of ingredients for tomorrow, it's important to estimate the number of customers that they will have based on the information from the previous days. Although the natural approach would be to calculate the average number of customers (which in this case would be 97), it's worth noting that this method could lead to a shortage of food on three out of the first five days if it were implemented in practice.*

In order to examine the distribution of the number of customers coming to the restaurant, we need to make some assumptions. Let's assume that there are a total of  $n$  equally-sized time slots throughout the day, with each slot being small enough so that no more than one customer can enter the restaurant during a given slot. We'll also assume that the probability of a customer entering the restaurant during a slot is denoted by  $p$ , and that the occurrence of a customer entering during one slot is independent of any other slot.

Formally, let  $X_i \triangleq 1[\text{there is a customer coming in the } i\text{-th slot}]$  for  $i \in [n]$ . Then we know  $X_i \sim \text{Ber}(p)$  and  $X_i$ 's are mutually independent. Let  $Z_n = \sum_{i=1}^n X_i$  and  $\lambda = \mathbb{E}[Z_n] = pn$ . Now let's compute the distribution of the number of customers  $Z_n$ . For any constant  $k \in \mathbb{N}$ ,

$$\begin{aligned} \Pr[Z_n = k] &= \binom{n}{k} p^k (1-p)^{n-k} \\ &= \frac{n(n-1) \cdots (n-k+1)}{k!} \cdot \left(\frac{\lambda}{n}\right)^k \cdot \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{n(n-1) \cdots (n-k+1)}{n^k} \cdot \frac{\lambda^k}{k!} \cdot \left(1 - \frac{\lambda}{n}\right)^n \cdot \left(1 - \frac{\lambda}{n}\right)^{-k}. \end{aligned} \quad (1)$$

Note that  $\lambda$  and  $k$  are constants. Thus, when  $n \rightarrow \infty$ , Equation (1) equals to  $\frac{\lambda^k}{k!} e^{-\lambda}$  and  $Z_n$  follows Poisson distribution with mean  $\lambda$ .

**Definition 1 (Poisson Distribution)** *We say a random variable  $X$  follows Poisson distribution with mean  $\lambda$ , written as  $X \sim \text{Pois}(\lambda)$ , if for any  $k \in \mathbb{Z}$ ,*

$$\Pr[X = k] = \begin{cases} \frac{\lambda^k}{k!} e^{-\lambda} & \text{if } k \geq 0, \\ 0 & \text{if } k < 0. \end{cases}$$

Since we get the distribution of  $Z_n$  by taking the limit, we need to verify that it is a distribution and its mean is indeed  $\lambda$ :

- We have  $\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = 1$ . Thus it is indeed a distribution.
- Since

$$\mathbf{E}[Z_n] = \sum_{k=0}^{\infty} k \frac{\lambda^k}{k!} e^{-\lambda} = \lambda \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} e^{-\lambda} = \lambda \sum_{k=0}^{\infty} \frac{\lambda^k}{(k)!} e^{-\lambda} = \lambda,$$

the expectation of  $Z_n$  indeed equals to  $\lambda$ .

What is the distribution of the number of customers in two days? It follows from the following property of Poisson distributions.

**Proposition 2** Suppose  $X_1 \sim \text{Pois}(\lambda_1)$  and  $X_2 \sim \text{Pois}(\lambda_2)$  are two independent random variables. Then

$$X_1 + X_2 \sim \text{Pois}(\lambda_1 + \lambda_2).$$

*Proof.* For  $n \geq 0$ ,

$$\begin{aligned} \Pr[X_1 + X_2 = n] &= \sum_{m=0}^n \Pr[X_1 = m] \cdot \Pr[X_2 = n - m] \\ &= \sum_{m=0}^n \frac{\lambda_1^m}{m!} e^{-\lambda_1} \cdot \frac{\lambda_2^{n-m}}{(n-m)!} e^{-\lambda_2} \\ &= e^{-(\lambda_1 + \lambda_2)} \cdot \sum_{m=0}^n \binom{n}{m} \frac{\lambda_1^m \lambda_2^{n-m}}{n!} \\ &= e^{-(\lambda_1 + \lambda_2)} \frac{(\lambda_1 + \lambda_2)^n}{n!} \end{aligned}$$

□

It is easy to extend the proposition to more independent Poisson random variables.

**Corollary 3** Suppose that  $X_1, X_2, \dots, X_n$  are  $n$  mutually independent random variables where  $X_i \sim \text{Pois}(\lambda_i)$ . Then

$$\sum_{i=1}^n X_i \sim \text{Pois}\left(\sum_{i=1}^n \lambda_i\right).$$

## 2 Poisson Processes

### 2.1 Definition of a Poisson Process

If we count the number of customers during a period of time rather than a single day, e.g., from day  $t_1$  to day  $t_2$ , it should follow  $\text{Pois}((t_2 - t_1)\lambda)$ . A Poisson process summarizes all the relevant information.

**Definition 4** A Poisson process  $\{N(s) : s \geq 0\}$  with rate  $\lambda$  is a stochastic process that

1.  $N(0) = 0$ ;
2.  $\forall t, s \geq 0, N(t+s) - N(s) \sim \text{Pois}(\lambda t)$ ;
3.  $\forall t_0 \leq t_1 \leq \dots \leq t_n, N(t_1) - N(t_0), N(t_2) - N(t_1), \dots, N(t_n) - N(t_{n-1})$   
are mutually independent.

In fact, we can view the Poisson process in another way by considering the time gaps between arrivals. To see this, we first recall the exponential distribution.

**Definition 5** The probability density function of the exponential distribution with rate  $\lambda > 0$  is

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

The corresponding cumulative distribution function is

$$F(t) = \int_{-\infty}^t f(x) dx = \int_0^t \lambda e^{-\lambda x} dx = 1 - e^{-\lambda t}.$$

Then the following proposition gives another characterization of the Poisson process.

**Proposition 6** Suppose that  $\tau_1, \tau_2, \dots, \tau_n$  is a sequence of independent random variables that  $\tau_i \sim \text{Exp}(\lambda)$ . Let  $T_n = \sum_{i=1}^n \tau_i$ . For  $s \geq 0$ , let  $N(s) = \max \{n \mid T_n \leq s\}$ . Then  $N(s)$  is a Poisson process with rate  $\lambda$ .

Before proving this proposition, we discuss some properties of the exponential distribution.

## 2.2 Properties of Exponential Distribution

**Proposition 7** Let  $X \sim \text{Exp}(\lambda)$ . Then  $\mathbf{E}[X] = \frac{1}{\lambda}$ .

*Proof.*

$$\begin{aligned} \mathbf{E}[X] &= \int_0^{\infty} t \cdot \lambda e^{-\lambda t} dt = \left( -te^{-\lambda t} \right) \Big|_0^{\infty} + \int_0^{\infty} e^{-\lambda t} dt \\ &= -\frac{1}{\lambda} e^{-\lambda t} \Big|_0^{\infty} = \frac{1}{\lambda}. \end{aligned}$$

□

**Proposition 8** Let  $X \sim \text{Exp}(\lambda)$ . Then  $\text{Var}[X] = \frac{1}{\lambda^2}$ .

*Proof.* Note that

$$\text{Var}[X] = \mathbf{E}[X^2] - \mathbf{E}[X]^2 = \mathbf{E}[X^2] - \frac{1}{\lambda^2}.$$

In Proposition 6, we can regard  $\tau_i$  as the time gap between the arrival of the  $i-1$ -th and the  $i$ -th customer. The parameter  $\lambda$  can be understood as the coming rate. Then the CDF of  $\tau_i$   $F(t) = 1 - e^{-\lambda t}$  is the probability that the  $i$ -th customer comes within time  $t$  after the arrival of the  $i-1$ -th person.

Since  $\lambda$  is the arriving rate, we can imagine that the average time between arrivals  $\mathbf{E}[\tau_i]$  is the reciprocal of  $\lambda$ . This gives an intuition of Proposition 7.

And

$$\begin{aligned}\mathbf{E}[X^2] &= \int_0^\infty t^2 \cdot \lambda e^{-\lambda t} dt = \left( -t^2 e^{-\lambda t} \right) \Big|_0^\infty + \int_0^\infty 2te^{-\lambda t} dt \\ &= 2 \int_0^\infty t \cdot e^{-\lambda t} dt = \mathbf{E}[X] \cdot \frac{2}{\lambda} = \frac{2}{\lambda^2}.\end{aligned}$$

Thus we have  $\text{Var}[X] = \frac{1}{\lambda^2}$ .  $\square$

**Proposition 9 (Lack of Memory)** Let  $X \sim \text{Exp}(\lambda)$ . Then for any  $t, s > 0$ ,

$$\Pr[X > t + s \mid X > s] = \Pr[X > t].$$

*Proof.*

$$\begin{aligned}\Pr[X > t + s \mid X > s] &= \frac{\Pr[X > t + s \wedge X > s]}{\Pr[X > s]} = \frac{\Pr[X > t + s]}{\Pr[X > s]} \\ &= \frac{e^{-\lambda(s+t)}}{e^{-\lambda s}} = e^{-\lambda t}.\end{aligned}$$

$\square$

**Proposition 10 (Exponential Races)** Let  $X_1 \sim \text{Exp}(\lambda_1)$  and  $X_2 \sim \text{Exp}(\lambda_2)$  be two independent random variables. Then  $Y \triangleq \min\{X_1, X_2\} \sim \text{Exp}(\lambda_1 + \lambda_2)$ .

*Proof.* By the independence, we have

$$\Pr[Y > t] = \Pr[X_1 > t \wedge X_2 > t] = \Pr[X_1 > t] \cdot \Pr[X_2 > t] = e^{-(\lambda_1 + \lambda_2)t}.$$

$\square$

Proposition 10 describes the distribution of the earliest customer of two restaurants. And we can easily generalize this to the case of more restaurants.

**Corollary 11** Let  $X_1, X_2, \dots, X_n$  be  $n$  mutually independent random variables where  $X_i \sim \text{Exp}(\lambda_i)$ . Then  $Y \triangleq \min\{X_1, X_2, \dots, X_n\}$  has an exponential distribution with rate  $\sum_{i=1}^n \lambda_i$ .

Now we consider the problem “who wins the race?”. That is, the restaurants are racing to see who will first have a customer. We first assume that there are only two random variables. Let  $f_\lambda$  be the probability density function of exponential distribution with rate  $\lambda$ . Using the law of total probability, we can compute the probability that  $X_1$  wins the race as follows:

$$\begin{aligned}\Pr[X_1 < X_2] &= \int_0^\infty f_{\lambda_1}(t) \Pr[X_2 \geq t] dt \\ &= \int_0^\infty \lambda_1 e^{-\lambda_1 t} e^{-\lambda_2 t} dt \\ &= \lambda_1 \int_0^\infty e^{-(\lambda_1 + \lambda_2)t} dt = \frac{\lambda_1}{\lambda_1 + \lambda_2}.\end{aligned}$$

Thus, clearly, the probability that  $X_i$  wins the race among  $n$  random variables is  $\frac{\lambda_i}{\sum_{j=1}^n \lambda_j}$ .



### 2.3 Proof of Proposition 6

**Proposition 12 (Proposition 6 restated)** Suppose that  $\tau_1, \tau_2, \dots, \tau_n$  is a sequence of independent random variables that  $\tau_i \sim \text{Exp}(\lambda)$ . Let  $T_n = \sum_{i=1}^n \tau_i$ . For  $s \geq 0$ , let  $N(s) = \max \{n \mid T_n \leq s\}$ . Then  $N(s)$  is a Poisson process with rate  $\lambda$ .

*Proof.* Note that  $T_n = \sum_{i=1}^n \tau_i$  is the arrival time of the  $n$ -th customer. Let  $g_n$  be the probability density function of  $T_n$ . First we prove that the distribution of  $T_n$  follows the Gamma distribution  $\Gamma(n, \lambda)$ :

$$g_n(t) = \begin{cases} \lambda e^{-\lambda t} \cdot \frac{(\lambda t)^{n-1}}{(n-1)!} & t \geq 0, \\ 0 & t < 0. \end{cases}$$

We prove this by induction. Note that when  $n = 1$ ,  $T_1 = \tau_1 \sim \text{Exp}(\lambda) = \Gamma(1, \lambda)$ . Suppose that  $T_n \sim \Gamma(n, \lambda)$  for some  $n \geq 1$ . By the independence of  $T_n$  and  $\tau_{n+1}$ , for  $t \geq 0$  we have

$$\begin{aligned} g_{n+1}(t) &= \int_0^t g_n(s) \cdot f_\lambda(t-s) ds \\ &= \int_0^t \lambda e^{-\lambda s} \cdot \frac{(\lambda s)^{n-1}}{(n-1)!} \cdot \lambda e^{-\lambda(t-s)} ds \\ &= \lambda e^{-\lambda t} \frac{\lambda^n}{(n-1)!} \int_0^t s^{n-1} ds \\ &= \lambda e^{-\lambda t} \frac{\lambda^n}{(n-1)!} \cdot \frac{t^n}{n} = \lambda e^{-\lambda t} \cdot \frac{(\lambda t)^n}{n!}. \end{aligned}$$

Then we compute the distribution of  $N(t)$ .

$$\begin{aligned} \Pr [N(t) = n] &= \Pr [T_n \leq t \wedge T_{n+1} > t] \\ &= \int_0^t g_n(s) \cdot \Pr [\tau_{n+1} > t-s] ds \\ &= \int_0^t \lambda e^{-\lambda s} \cdot \frac{(\lambda s)^{n-1}}{(n-1)!} \cdot e^{-\lambda(t-s)} ds \\ &= \lambda^n e^{-\lambda t} \frac{t^n}{n!}. \end{aligned}$$

Thus,  $N(t) \sim \text{Pois}(\lambda t)$ . Then we verify that  $\{N(t) : t \geq 0\}$  satisfies the three conditions in Definition 4.

First it is clear that  $N(t) = 0$  when  $t = 0$ . By the lack of memory property, we know that  $N(s+t) - N(s)$  follows the same distribution as  $N(t) - N(0)$ , which  $\text{Pois}(\lambda t)$ . Furthermore, it is easy to see that  $N(s+t) - N(s)$  is independent of  $N(r)$  for all  $r \leq s$  again by the lack of memory property. It implies that  $N(t)$  has independent increments, and hence completes our proof of Proposition 6.  $\square$

Imagine the difference between  $N(s+t) - N(s)$  and  $N(t) - N(0)$ . In the  $N(t) - N(0)$ , we start to wait for the first customer at time 0, while in  $N(s+t) - N(s)$ , at time  $s$ , we might have waited for the first customer in the period for sometime. However, due to the lack of memory property of the waiting time, this is equivalent to start to wait at time  $s$ .

### 2.4 Thinning

In the example of customers coming into the restaurant, sometimes we have a more detailed characterization of customers, such as the gender.

We associate an i.i.d. random variable  $Y_i$  with  $i$ -th arrival, and then use the value of  $Y_i$  to label the arrival and separate the Poisson process into several. Suppose that  $Y_i \in \mathbb{N}$  and let  $p_j = \Pr[Y_i = j]$ . For all  $j \in \text{Range}(Y_i)$ , let  $N_j(t)$  denote the number of arrivals with label  $j$  that have arrived by time  $t$ . Then  $\{N_j(t)\}$  is called a thinning of a Poisson process. We have the following useful and surprising proposition.

**Proposition 13** *For each  $j$ ,  $\{N_j(t) : t \geq 0\}$  is a Poisson process with rate  $p_j\lambda$ . Moreover, the collections of processes  $\{\{N_j(t) : t \geq 0\} : j \in \text{Range}(Y)\}$  are mutually independent.*

*Proof.* For convenience we assume that  $Y_i \in \{0, 1\}$ . Then the following calculation concludes the independence and the distribution of  $N_j(t)$  at the same time.

$$\begin{aligned} \Pr[N_0(t) = j \wedge N_1(t) = k] &= \Pr[N_0(t) = j \wedge N(t) = k + j] \\ &= \Pr[N(t) = k + j] \cdot \Pr[N_0(t) = j \mid N(t) = k + j] \\ &= e^{-\lambda t} \frac{(\lambda t)^{j+k}}{(j+k)!} \cdot \binom{j+k}{j} p_0^j p_1^k \\ &= e^{-p_0 \lambda t} \frac{(p_0 \lambda t)^j}{j!} \cdot e^{-p_1 \lambda t} \frac{(p_1 \lambda t)^k}{k!}. \end{aligned}$$

Thus, when there are  $n$  labels, it is easy to verify that  $N_j(t) \sim \text{Pois}(p_j\lambda)$  and they are mutually independent.  $\square$

Let's see an application of Poisson process.

**Example 2 (Maximum Likelihood of Poisson Process)** *Suppose there are two editors reading a book of 300 pages. Editor A finds 100 typos in the book, and editor B finds 120 typos, 80 of which are in common.*

*Suppose that the author's typos follow a Poisson process with some unknown rate  $\lambda$  per page. The two editors catch typos with unknown probabilities of success  $p_A$  and  $p_B$  respectively. We want to know how many typos there actually are. We can estimate this by determining  $\lambda$ ,  $p_A$  and  $p_B$ . Clearly, there are four types of typos:*

*Type 1 The typo is found by neither of the editors. This happens w.p.  $q_1 = (1 - p_A)(1 - p_B)$ .*

*Type 2 The typo is found only by editor A. This happens w.p.  $q_2 = (1 - p_A)p_B$ .*

*Type 3 The typo is found only by editor B. This happens w.p.  $q_3 = (1 - p_B)p_A$ .*

*Type 4 The typo is found by both editors. This happens w.p.  $q_4 = p_A p_B$ .*

*So the occurrence of type  $i$  typos follows an independent Poisson process with rate  $q_i\lambda$ . That is, letting  $X_1, X_2, X_3$  and  $X_4$  be the occurrence time of the corresponding type of typos in this book, then  $X_i \sim \text{Pois}(300q_i\lambda)$ . Note that*

Here is an example explains why this proposition is surprising. Assume that the customers coming into a restaurant is a Poisson process, and each customer is male or female independently with probability  $1/2$  and  $1/2$  respectively. In fact we can assume that we flip a coin to determine whether the arriving customer is male or female. So intuitively, one might think that a large number of men (such as 50) arriving in one hour would indicate a large volume of business and hence a larger than normal number of women arriving. However this proposition tells us that the number of men arriving and the number of women arriving are independent.

there are 20 typos of type 2, 40 typos of type 3 and 80 typos of type 4. We claim that the most likely values of the rates are

$$\begin{cases} 300(1 - p_A)p_B\lambda = 20, \\ 300(1 - p_B)p_A\lambda = 40, \\ 300p_Ap_B\lambda = 80. \end{cases}$$

This yields that  $p_A = \frac{2}{3}$ ,  $p_B = \frac{4}{5}$  and  $\lambda = \frac{1}{2}$ .

It remains to prove that claim. Suppose  $X \sim \text{Pois}(\theta)$  with some unknown  $\theta$ . Then given  $z$ , our goal is to find  $\arg \max_{\theta} \Pr[X = z \mid X \sim \text{Pois}(\theta)]$ . Note that  $\Pr[X = z \mid X \sim \text{Pois}(\theta)] = e^{-\theta} \frac{\theta^z}{z!}$  and  $\log e^{-\theta} \frac{\theta^z}{z!} = -\theta + z \cdot \log \theta$ . So it is equivalent to find

$$\arg \max_{\theta} -\theta + z \cdot \log \theta \tag{2}$$

Let the derivation of Equation (2) equals to 0. We have  $\theta = z$ , that is,  $\arg \max_{\theta} e^{-\theta} \frac{\theta^z}{z!} = z$ .

# [AI2613 Lecture 10] Poisson Approximation

May 4, 2023

## 1 Coupon Collector Problem with Non-Uniform Coupons

Recall the coupon collector problem we met many times in this course: If each box of a brand of cereals contains a coupon which is chosen from  $n$  different types uniformly at random, then we need to buy  $nH_n$  boxes in expectation to collect all types of coupons.

In this lecture, we generalize the setting by allowing the non-uniformity. Suppose that each purchase yields a coupon of type  $j$  w.p.  $p_j$  for  $j \in [n]$  and the coupon types contained in different boxes are independent, where  $\sum_{j=1}^n p_j = 1$ . Let  $N_j$  be the first time that we get type  $j$ . Then  $N_j$  follows the geometric distribution with parameter  $p_j$ . Let  $N$  be the number of purchases until all  $n$  types of coupons are collected, that is,  $N = \max_{j \in [n]} N_j$ . We would like to compute  $\mathbf{E}[N]$  to see how many times of purchases is needed in expectation. However, it is not easy to compute the expected value of  $\max_{j \in [n]} N_j$  since  $N_j$ 's are not independent.

### 1.1 Coupon Collector Problem with Poisson Draw

We consider a variation of the coupon collector problem where the coupons are collected with Poisson draw. That is, each arrival of the Poisson process with rate 1 brings a coupon and the probability of the coupon being of type  $j$  is  $p_j$ . Note that this process is different from the ordinary coupon collector problem since the arrival time is random.

Recall the thinning of Poisson process we discussed in the last lecture. Let  $X_j(t)$  be the number of type  $j$  coupons we collect in time  $[0, t]$  with Poisson draw. Then  $\{X_j(t)\}$  is a thinning of the process, meaning that  $\{X_j(t)\}$  is a Poisson process with rate  $p_j$  and  $X_j(t)$  is independent of  $X_i(t)$  for  $i \neq j$ . For  $j \in [n]$ , let  $T_j \triangleq \min \{t \mid X_j(t) = 1\}$  be the first time that type  $j$  coupon appears. Obviously,  $T_j$  is the same as  $\tau_j(1)$ <sup>1</sup> and  $T_j \sim \text{Exp}(p_j)$ .

To determine the time of collecting all kinds of coupons, we need to compute  $\mathbf{E}[T]$  where  $T = \max_{j \in [n]} T_j$ . The following proposition to compute expectation is useful.

**Proposition 1** *Let  $X$  be a non-negative random variable.*

- *If  $X$  is discrete and  $X \in \mathbb{N}$ , then  $\mathbf{E}[X] = \sum_{t=1}^{\infty} \Pr[X \geq t]$ .*
- *If  $X$  is continuous, then  $\mathbf{E}[X] = \int_0^{\infty} \Pr[X \geq t] dt$ .*

*Proof.*

<sup>1</sup> Here  $\tau_j(1)$  denotes the time gap between the arrival of the customers with coupon  $j$ .

- When  $X$  is discrete, we apply the double counting trick:

$$\begin{aligned} \mathbf{E}[X] &= \sum_{s=1}^{\infty} s \Pr[X = s] = \sum_{s=1}^{\infty} \sum_{t=1}^s \Pr[X = s] \\ &= \sum_{t=1}^{\infty} \sum_{s=t}^{\infty} \Pr[X = s] = \sum_{t=1}^{\infty} \Pr[X \geq t]. \end{aligned}$$

- When  $X$  is continuous,

$$\begin{aligned} \mathbf{E}[X] &= \mathbf{E}\left[\int_0^X 1 \, dt\right] = \mathbf{E}\left[\int_0^{\infty} \mathbf{1}[X \geq t] \, dt\right] \\ &\stackrel{(\heartsuit)}{=} \int_0^{\infty} \mathbf{E}[\mathbf{1}[X \geq t]] \, dt = \int_0^{\infty} \Pr[X \geq t] \, dt, \end{aligned}$$

where  $(\heartsuit)$  comes from the *Fubini's theorem*.

□

Note that for any  $t \in \mathbb{R}_{\geq 0}$ ,

$$\Pr[T \geq t] = 1 - \Pr[T < t] = 1 - \prod_{j=1}^n \Pr[T_j < t] = 1 - \prod_{j=1}^n (1 - e^{-p_j t}).$$

By the continuous version of Proposition 1, we have

$$\mathbf{E}[T] = \int_0^{\infty} \Pr[T \geq t] \, dt = \int_0^{\infty} 1 - \prod_{j=1}^n (1 - e^{-p_j t}) \, dt.$$

That is, we need a time of  $\int_0^{\infty} (1 - e^{-p_j t}) \, dt$  in expectation to collect all kinds of coupons.

## 1.2 Standard Coupon Collector Problem

Then we relate the result we obtained in the previous section on the coupon collector with Poisson draw to the standard coupon collector problem by the technique of coupling. Specifically, let  $\tau_i$  denote the time gap between the  $i-1$ -th and the  $i$ -th arrival. Imagine the standard version as one customer coming with a coupon in hand with constant time gap between arrivals. We couple the two process by letting the  $i$ -th arrival in the Poisson version carry the same type of coupon with the  $i$ -th arrival in the ordinary version.

Recall that  $N$  is the number of purchases until all  $n$  types of coupons are collected in the standard coupon collector problem. Then we have  $T = \sum_{i=1}^N \tau_i$ . Note that  $\tau_i \sim \text{Exp}(1)$  and  $\mathbf{E}[\tau_i] = 1$ . If  $N$  is a constant, we can deduce  $\mathbf{E}[N] = \mathbf{E}\left[\sum_{i=1}^N \tau_i\right] = \mathbf{E}[T]$  directly. However,  $N$  is a random variable and thus the summation and expectation are not guaranteed to be exchangeable. To show the validity of  $\mathbf{E}[N] \mathbf{E}[\tau_i] = \mathbf{E}\left[\sum_{i=1}^N \tau_i\right]$  in this case, we make use of the Wald's equation introduced before.

**Theorem 2 (Wald's Equation)** Let  $X_1, X_2, \dots$  be  $n$  i.i.d. random variables that  $\mathbf{E}[|X_1|] < \infty$ . Let  $T$  be a stopping time that  $\mathbf{E}[T] < \infty$ . Then we have  $\mathbf{E}[\sum_{t=1}^T X_t] = \mathbf{E}[T] \mathbf{E}[X_1]$ .

It is easy to verify that  $\mathbf{E}[\tau_i] = 1 < \infty$  and  $\mathbf{E}[N] < \infty$  in our case. So applying the Wald's equation, we have  $\mathbf{E}[N] \mathbf{E}[\tau_i] = \mathbf{E}[\sum_{i=1}^N \tau_i]$  and sequentially

$$\mathbf{E}[N] = \mathbf{E}[T] = \int_0^\infty 1 - \prod_{j=1}^n (1 - e^{-p_j t}) dt. \quad (1)$$

Then we go back to the coupon collector problem with uniform coupons for sanity check. Let  $x = e^{-\frac{t}{n}}$ . If  $p_j = \frac{1}{n}$  for any  $j \in [n]$ , we have

$$\begin{aligned} \mathbf{E}[N] &= \int_0^\infty 1 - \prod_{j=1}^n (1 - e^{-p_j t}) dt \\ &= n \int_0^\infty 1 - (1 - x)^n d \log x \\ &= n \int_0^\infty \frac{1}{x} - \frac{(1 - x)^n}{x} dx \\ &= n \int_0^\infty \sum_{k=1}^n \frac{(1 - x)^{k-1}}{x} - \frac{(1 - x)^k}{x} dx \\ &\stackrel{(\heartsuit)}{=} n \sum_{k=1}^n \int_0^\infty (1 - x)^{k-1} dx \\ &= n \sum_{k=1}^n \frac{1}{k} = nH_n, \end{aligned}$$

where the  $(\heartsuit)$  follows from the *Fubini's theorem*. This verifies Equation (1) when the types of coupons are uniform.

## 2 Balls-into-Bins

Recall the balls-into-bins problem where we throw  $m$  identical balls into  $n$  bins. For  $i \in [n]$ , let  $X_i$  be the number of balls in the  $i$ -th bin. Then we have  $X_i \sim \text{Binom}(m, \frac{1}{n})$  and  $\mathbf{E}[X_i] = \frac{m}{n}$ . This model can be used to describe the scheme of the hash table. To avoid frequent collision when mapping the keys into slots, it is natural for us to be concerned about the value of  $\max_{i \in [n]} X_i$ . However, we are faced with the difficulty that  $X_i$ 's are not independent when computing the distribution of  $\max_{i \in [n]} X_i$ . It turns out that one can use independent Poisson variables to approximate the distribution. First we have:

**Theorem 3** The distribution of  $(X_1, X_2, \dots, X_n)$  is the same as that of  $(Y_1, Y_2, \dots, Y_n)$  on condition that  $\sum_{i=1}^n Y_i = m$  where  $Y_i \sim \text{Pois}(\lambda)$  are independent Poisson random variables with an arbitrary rate  $\lambda$ .

*Proof.* Given  $(a_1, a_2, \dots, a_n) \in \mathbb{N}^n$  and  $\sum_{i=1}^n a_i = m$ , we have

$$\Pr[(X_1, X_2, \dots, X_n) = (a_1, a_2, \dots, a_n)] = \frac{1}{n^m} \cdot \frac{m!}{a_1! a_2! \dots a_n!}. \quad (2)$$

And

$$\begin{aligned} & \Pr\left[(Y_1, Y_2, \dots, Y_n) = (a_1, a_2, \dots, a_n) \mid \sum_{i=1}^n Y_i = m\right] \\ &= \frac{\Pr[(Y_1, Y_2, \dots, Y_n) = (a_1, a_2, \dots, a_n) \wedge \sum_{i=1}^n Y_i = m]}{\Pr\left[\sum_{i=1}^n Y_i = m\right]} \\ &= \frac{\prod_{i=1}^n \Pr[Y_i = a_i]}{\Pr\left[\sum_{i=1}^n Y_i = m\right]} \\ &= \frac{\prod_{i=1}^n e^{-\lambda} \frac{\lambda^{a_i}}{a_i!}}{e^{-\lambda n} \frac{(\lambda n)^m}{m!}} = \frac{1}{n^m} \cdot \frac{m!}{a_1! a_2! \dots a_n!}, \end{aligned}$$

which equals to the RHS of Equation (2).  $\square$

Furthermore, we can deduce the following corollary from Theorem 3.

**Corollary 4** Let  $f: \mathbb{N}^n \rightarrow \mathbb{R}$  be an arbitrary function and  $Y_1, Y_2, \dots, Y_n$  be  $n$  independent Poisson random variables with rate  $\lambda = \frac{m}{n}$ . Then we have

$$\mathbb{E}[f(X_1, X_2, \dots, X_n)] \leq e^{\sqrt{m}} \cdot \mathbb{E}[f(Y_1, Y_2, \dots, Y_n)].$$

*Proof.* By the law of total probability, we have

$$\begin{aligned} \mathbb{E}[f(Y_1, Y_2, \dots, Y_n)] &= \sum_{k=0}^{\infty} \mathbb{E}\left[f(Y_1, Y_2, \dots, Y_n) \mid \sum_{i=1}^n Y_i = k\right] \Pr\left[\sum_{i=1}^n Y_i = k\right] \\ &\geq \mathbb{E}\left[f(Y_1, Y_2, \dots, Y_n) \mid \sum_{i=1}^n Y_i = m\right] \Pr\left[\sum_{i=1}^n Y_i = m\right] \\ &= \mathbb{E}[f(X_1, X_2, \dots, X_n)] \Pr\left[\sum_{i=1}^n Y_i = m\right]. \end{aligned}$$

Note that  $\sum_{i=1}^n Y_i \sim \text{Pois}(m)$ , then we have

$$\Pr\left[\sum_{i=1}^n Y_i = m\right] = e^{-m} \frac{m^m}{m!} > \frac{1}{e^{\sqrt{m}}},$$

where the inequality comes from the *Stirling's formula*.  $\square$

Equipped with Corollary 4, we have the following theorem to bound  $X = \max_{i \in [n]} X_i$ .

We can see from the proof of Corollary 4 that the choice of  $\lambda = \frac{m}{n}$  is to maximize  $\Pr\left[\sum_{i=1}^n Y_i = m\right]$ .

**Theorem 5 (Max Load)** When  $m = n$ , we have  $X = \Theta\left(\frac{\log n}{\log \log n}\right)$  w.p.  $1 - o(1)$ .

*Proof.* First we prove the upper bound, that is, there exists a constant  $c_1$  such that  $\Pr\left[X \geq \frac{c_1 \log n}{\log \log n}\right] = o(1)$ . Let  $k = \frac{c_1 \log n}{\log \log n}$  for brevity. By union

bound, we have

$$\begin{aligned}\Pr[X \geq k] &= \Pr[\exists i \in [n], X_i \geq k] \leq \sum_{i=1}^n \Pr[X_i \geq k] \\ &= n \cdot \Pr[X_1 \geq k] \leq n \cdot \binom{n}{k} \frac{1}{n^k} \leq n \cdot \left(\frac{en}{k}\right)^k \frac{1}{n^k} = n \cdot \left(\frac{e}{k}\right)^k.\end{aligned}$$

Note that

$$\begin{aligned}k \log k &= \frac{c_1 \log n}{\log \log n} \cdot (\log \log n - \log \log \log n + \log c_1) \\ &> c_1 \log n \left(1 - \frac{\log \log \log n}{\log \log n}\right) > \frac{c_1}{2} \log n.\end{aligned}$$

Letting  $c = 6$ , we have that

$$\log n + k - k \log k < -\log n.$$

Thus,  $\Pr[X \geq k] \leq n \cdot \left(\frac{e}{k}\right)^k < \frac{1}{n} = o(1)$  for  $c_1 = 6$ .

Then we prove the lower bound. Again let  $g = \frac{c_2 \log n}{\log \log n}$  for a constant  $c_2$ . Let  $f(X_1, X_2, \dots, X_n) \triangleq \mathbf{1}[X < g] = \mathbf{1}[\max_{i \in [n]} X_i < g]$ . Then by Corollary 4,

$$\begin{aligned}\Pr[X < g] &= \mathbf{E}[f(X_1, X_2, \dots, X_n)] \\ &\leq e\sqrt{n} \cdot \mathbf{E}[f(Y_1, Y_2, \dots, Y_n)] \\ &= e\sqrt{n} \cdot \Pr\left[\max_{i \in [n]} Y_i < g\right].\end{aligned}\tag{3}$$

By the definition of  $Y_i$  in Corollary 4, we have

$$\begin{aligned}\Pr\left[\max_{i \in [n]} Y_i < g\right] &= (\Pr[Y_1 \leq g])^n = (1 - \Pr[Y_1 > g])^n \\ &\leq (1 - \Pr[Y_1 = g+1])^n = \left(1 - \frac{1}{(g+1)!e}\right)^n \leq e^{-\frac{n}{(g+1)!e}}\end{aligned}$$

Note that

$$\begin{aligned}\log(g+1)! &= \sum_{i=1}^{g+1} \log i < \int_1^{g+2} \log x \, dx \\ &= (g+2) \log(g+2) - g - 1 \leq (g+2) \log g - g + 3 \\ &= \frac{c_2 \log n + 2 \log \log n}{\log \log n} (\log \log n - \log \log \log n + \log c_2) - \frac{c_2 \log n}{\log \log n} + 3 \\ &\leq c_2 \log n - \log \log n - 2.\end{aligned}$$

Letting  $c_2 = 1$ , we have  $\log(g+1)! \leq \log n - \log \log n - 2$  and sequentially

$$e(g+1)! \leq \frac{n}{e \log n}.$$

Thus,

$$\Pr\left[\max_{i \in [n]} Y_i < g\right] \leq e^{-\frac{n}{(g+1)!e}} \leq e^{-e \log n} = n^{-e}.$$

Combining with Equation (3), we have  $\Pr\left[X < \frac{\log n}{\log \log n}\right] \leq e\sqrt{n} \cdot n^{-e} = o(1)$ .

□



# [AI2613 Lecture 11] Brownian Motion

May 14, 2023

## 1 Brownian Motion

Brownian motion describes the random motion of small particles suspended in a liquid or in a gas. This process was named after the botanist Robert Brown, who observed and studied a jittery motion of pollen grains suspended in water under a microscope. Later, Albert Einstein gave a physical explanation of this phenomenon. In mathematics, Brownian motion is characterized by the *Wiener process*, named after Norbert Wiener, a famous mathematician and the originator of cybernetics.

To motivate the definition of Brownian motion, we start from the 1-D random walk starting from 0. Let  $Z_t$  be our position at time  $t$  and  $X_t$  be the move of the  $t$ -th step. The value of  $X_t$  is chosen from  $\{-1, 1\}$  uniformly at random. Note that  $Z_0 = 0$  and  $Z_{t+1} = Z_t + X_t$ . So  $Z_T = \sum_{t=0}^{T-1} X_t$ . Then we have

$$\mathbf{E}[Z_T] = 0 \text{ and } \mathbf{Var}[Z_T] = \sum_{t=0}^{T-1} \mathbf{Var}[X_t] = T.$$

Suppose now we move with every  $\Delta t$  seconds and with step length  $\delta$ . Then our position at time  $T$  is  $Z(T) = \delta \sum_{t=1}^{\frac{T}{\Delta t}} X_t$ . We are interested in the behavior of the process when  $\Delta t \rightarrow 0$ . We have

$$\mathbf{E}[Z(T)] = 0 \text{ and } \mathbf{Var}[Z(T)] = \delta^2 \sum_{t=1}^{\frac{T}{\Delta t}} \mathbf{Var}[X_t] = \delta^2 \cdot \frac{T}{\Delta t}.$$

We can identify the expectation and the variance of this process with the discrete random walk when  $\Delta t \rightarrow 0$  by choosing  $\delta = \sqrt{\Delta t}$ . It follows from the central limit theorem that

$$Z(T) = \sqrt{\Delta t} \sum_{t=1}^{\frac{T}{\Delta t}} X_t \xrightarrow{\Delta t \rightarrow 0} \sqrt{\Delta t} \mathcal{N}(0, \frac{T}{\Delta t}) = \mathcal{N}(0, T).$$

In other words, the “continuous” version of the 1-D random walk follows  $\mathcal{N}(0, T)$  at time  $T$ . This is the basis of the Wiener process. Now we introduce its formal definition.

**Definition 1** (Standard Brownian Motion / Wiener Process). *We say a stochastic process  $\{W(t)\}_{t \geq 0}$  is a standard Brownian motion or Wiener process if it satisfies*

- $W(0) = 0$ ;

- **Independent increments:**  $\forall 0 \leq t_0 \leq t_1 \leq \dots \leq t_n$ ,  $W(t_1) - W(t_0)$ ,  $W(t_2) - W(t_1)$ ,  $\dots$ ,  $W(t_n) - W(t_{n-1})$  are mutually independent;
- **Stationary increments:**  $\forall s, t > 0$ ,  $W(s+t) - W(s) \sim \mathcal{N}(0, t)$ ;
- $W(t)$  is continuous almost surely.<sup>1</sup>

Recall that the probability density of the Gaussian distribution  $N(\mu, \sigma^2)$  is

$$f_{N(\mu, \sigma^2)}(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

We use  $\Phi(\cdot)$  to denote the CDF of  $N(0, 1)$ , namely  $\Phi(t) = \int_{-\infty}^t f_{N(0,1)}(x) dx$ .

In the following, we use  $f_t(x)$  to denote the probability density of  $N(0, t)$ . For any  $t_1 \leq t_2 \leq \dots \leq t_n$ , the joint density of  $W(t_1), W(t_2), \dots, W(t_n)$  is

$$f(x_1, \dots, x_n) = f_{t_1}(x_1) f_{t_2-t_1}(x_2 - x_1) \dots f_{t_n-t_{n-1}}(x_n - x_{n-1})$$

**Example 1.** Let  $0 \leq s \leq t$ . We can compute the conditional distribution of  $X(s)$  when  $X(t) = y$ . We use  $f_{s|t}(x|y)$  to denote the probability density of  $X(s) = x$  conditioned  $X(t) = y$ . Clearly

$$f_{s|t}(x|y) = \frac{f_s(x) f_{t-s}(y-x)}{f_t(y)} = C \cdot \exp\left(-\frac{(x-ys/t)^2}{2s(t-s)/t}\right),$$

where  $C$  is some universal constant irrelevant to  $x, y, s, t$ . As a result, the conditional distribution is the Gaussian  $N(\frac{s}{t}y, \frac{s}{t}(t-s))$ .

Let  $\{W(t)\}_{t \geq 0}$  be a standard Brownian motion. If  $\{X(t)\}_{t \geq 0}$  satisfies  $X(t) = \mu \cdot t + \sigma W(t)$ , we call  $\{X(t)\}_{t \geq 0}$  a  $(\mu, \sigma^2)$  Brownian motion. Clearly,  $X(t) \sim N(\mu t, \sigma^2 t)$ .

## 2 The Hitting Time of a Brownian Motion

We consider the first arrival time of position  $b$  in a Brownian motion. This is called the *hitting time* of  $b$ . Let us first consider the standard Brownian motion  $\{W(t)\}$ . Define  $\tau_b \triangleq \inf\{t \geq 0 \mid W(t) > b\}$ . For any  $t > 0$ ,

$$\begin{aligned} \Pr[\tau_b < t] &= \Pr[\tau_b < t \wedge W(t) > b] + \Pr[\tau_b < t \wedge W(t) < b] \\ &= \Pr[W(t) > b] + \Pr[W(t) < b \mid \tau_b < t] \cdot \Pr[\tau_b < t]. \end{aligned}$$

Note that  $W(t) \sim \mathcal{N}(0, t)$ . Let  $\Phi$  be the cumulative distribution function of standard Gaussian distribution, that is,  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$ . Then

$$\Pr[W(t) > b] = \Pr\left[\frac{W(t)}{\sqrt{t}} > \frac{b}{\sqrt{t}}\right] = 1 - \Phi\left(\frac{b}{\sqrt{t}}\right).$$

<sup>1</sup> Let  $\Omega$  be the sample space. Then  $W$  can be viewed as a mapping from  $\mathbb{R} \times \Omega$  to  $\mathbb{R}$ . The meaning of “ $W(t)$  is continuous almost surely” is:  $\exists \Omega_0 \subseteq \Omega$  with  $\Pr[\Omega_0] = 1$  such that  $\forall \omega \in \Omega_0$ ,  $W(t, \omega)$  is continuous with regard to  $t$ .

This is called the *principle of reflection* of a standard Brownian motion.

Assuming we have known the value of  $\tau_b$  and  $\tau_b < t$ , we can regard  $\{W(t)\}_{t \geq \tau_b}$  as a Brownian motion starting from  $b$ . Thus, as Figure 1 shows,  $\Pr[W(t) < b \mid \tau_b < t] = \frac{1}{2}$ .

By direct calculation, we have  $\Pr[\tau_b < t] = 2 \left( 1 - \Phi \left( \frac{b}{\sqrt{t}} \right) \right)$ .

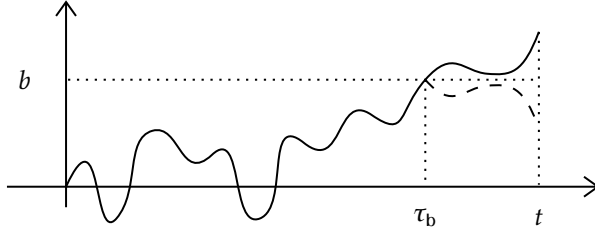


Figure 1: The hitting time and the reflection principle

It is more challenging to find the hitting time of a  $(\mu, \sigma^2)$  Brownian motion. The main difficulty is that the principle of reflection no longer holds when a nonzero drift  $\mu$  is present.

We can overcome the difficulty by leveraging the following useful lemma.

**Lemma 2.** Let  $Y_1, \dots, Y_n$  be i.i.d.  $N(\theta, \nu^2)$  random variables. Then the distribution of  $(Y_1, \dots, Y_n)$  conditioned on  $\sum_{i=1}^n Y_i = y$  is irrelevant to  $\theta$ .

*Proof.* Let  $X = \sum_{i=1}^n Y_i$ . We use  $f_{Y_1, \dots, Y_n | X}$  to denote the density of  $Y_1, \dots, Y_n$  conditioned on  $X$ . Then

$$\begin{aligned} f_{Y_1, \dots, Y_n | X}(y_1, \dots, y_n, x) &= \frac{f_{Y_1, \dots, Y_n, X}(y_1, \dots, y_n, x)}{f_X(x)} \\ &= \frac{f_{Y_1, \dots, Y_n}(y_1, \dots, y_{n-1}, x - \sum_{i=1}^{n-1} y_i)}{f_X(x)} \\ &= \frac{\exp\left(-\frac{(x - \sum_{i=1}^{n-1} y_i - n\theta)^2}{2\nu^2}\right) \prod_{i=1}^{n-1} \exp\left(-\frac{(y_i - \theta)^2}{2\nu^2}\right)}{\exp\left(-\frac{(x - n\theta)^2}{2n\nu^2}\right)}. \end{aligned}$$

A direct calculation shows that all terms on  $\theta$  cancel and therefore the lemma is proved.  $\square$

Here  $A \sim B$  means  $A = c \cdot B$  for some universal constant  $c$ .

The following corollary is immediate since all relevant random variables can be expressed as the sum of independent Gaussians.

**Corollary 3.** Let  $\{X(t)\}_{t \geq 0}$  be a  $(\mu, \sigma^2)$  Brownian motion. Conditioned on  $X(t) = x$ , for any  $t_1 \leq t_2 \leq \dots \leq t_n \leq t$ , the joint distribution of  $(X(t_1), X(t_2), \dots, X(t_n))$  is the same for all  $\mu$ .

Armed with this, we can calculate the hitting time of a  $(\mu, \sigma^2)$  Brownian motion.

**Lemma 4.** Let  $X(t)$  be a  $(\mu, \sigma^2)$  Brownian motion. For any  $y > x$ ,

$$\Pr[\tau_y \leq t \mid X(t) = x] = e^{-\frac{2y(y-x)}{t\sigma^2}}.$$

*Proof.* Applying Corollary 3, we know that  $\Pr[\tau_y \leq t \mid X(t) = x] = \Pr[\tau'_y \leq t \mid X'(t) = x]$  where  $X'(t)$  is a  $N(0, \sigma^2)$  Brownian motion and  $\tau'_y$  is the hitting time  $X'(t)$ .

Consider an *infinitesimal change*  $dx$ . It holds that

$$\Pr[\tau'_y \leq t \mid X'(t) \in [x, x+dx]] = \frac{\Pr[\tau'_y \leq t \wedge X'(t) \in [x, x+dx]]}{\Pr[X'(t) \in [x, x+dx]]}.$$

Since  $\Pr[X'(t) \in [x, x+dx]] = f_{X'(t)}(x)dx$ , we only need to calculate the numerator. Note that

$$\Pr[\tau'_y \leq t \wedge X'(t) \in [x, x+dx]] = \Pr[\tau'_y \leq t] \cdot \Pr[X'(t) \in [x, x+dx] \mid \tau'_y \leq t].$$

Applying the principle of reflection, the above is equal to

$$\begin{aligned} \Pr[\tau'_y \leq t] \cdot \Pr[X'(t) \in [2y-x-dx, 2y-x] \mid \tau'_y \leq t] &= \Pr[X'(t) \in [2y-x-dx, 2y-x] \wedge \tau'_y \leq t] \\ &= \Pr[X'(t) \in [2y-x-dx, 2y-x]] \\ &= f_{X'(t)}(2y-x)dx \end{aligned}$$

The second equality is due to that  $dx$  is infinitesimal and therefore  $x+dx < y$ . As a result, we have

$$\Pr[\tau_y \leq t \mid X'(t) = x] = \frac{f_{X'(t)}(2y-x)}{f_{X'(t)}(x)} = e^{-\frac{2y(y-x)}{t\sigma^2}}.$$

□

We are now ready to compute the hitting time  $\tau_y$ . When  $y \leq x$ , clearly  $\Pr[\tau_y \leq t \mid X(t) = y] = 1$ . Therefore,

$$\begin{aligned} \Pr[\tau_y \leq t] &= \int_{-\infty}^{\infty} \Pr[\tau_y \leq t \mid X(t) = x] \cdot f_{X(t)}(x) dx \\ &= \int_{-\infty}^y \Pr[\tau_y \leq t \mid X(t) = x] \cdot f_{X(t)}(x) dx + \Pr[X(t) \geq y] \\ &= \int_{-\infty}^y e^{-\frac{2y(y-x)}{t\sigma^2}} \cdot \frac{1}{\sqrt{2\pi t\sigma^2}} e^{-\frac{(x-\mu t)^2}{2t\sigma^2}} dx + \left(1 - \Phi\left(\frac{y-\mu t}{\sigma\sqrt{t}}\right)\right) \\ &= e^{\frac{2y\mu}{\sigma^2}} \left(1 - \Phi\left(\frac{\mu t + y}{\sigma\sqrt{t}}\right)\right) + \left(1 - \Phi\left(\frac{y-\mu t}{\sigma\sqrt{t}}\right)\right). \end{aligned}$$

It is not completely rigorous here since Corollary 3 only applies to the joint distribution of *finite* many random variables. Nevertheless, it is conceivable that the same holds for the whole process  $X(t)$ .

# [AI2613 Lecture 12] Gaussian Processes, Brownian Bridge

June 9, 2023

## 1 Gaussian Processes and Brownian Motion

In the last lecture, we defined the standard Brownian motion:

**Definition 1** (Standard Brownian Motion / Wiener Process). We say a stochastic process  $\{W(t)\}_{t \geq 0}$  is a standard Brownian motion or Wiener process if it satisfies

- $W(0) = 0$ ;
- **Independent increments:**  $\forall 0 \leq t_0 \leq t_1 \leq \dots \leq t_n, W(t_1) - W(t_0), W(t_2) - W(t_1), \dots, W(t_n) - W(t_{n-1})$  are mutually independent;
- **Stationary increments:**  $\forall s, t > 0, W(s+t) - W(s) \sim \mathcal{N}(0, t)$ ;
- $W(t)$  is continuous almost surely.<sup>1</sup>

<sup>1</sup> Let  $\Omega$  be the sample space. Then  $W$  can be viewed as a mapping from  $\mathbb{R} \times \Omega$  to  $\mathbb{R}$ . The meaning of “ $W(t)$  is continuous almost surely” is:  $\exists \Omega_0 \subseteq \Omega$  with  $\Pr[\Omega_0] = 1$  such that  $\forall \omega \in \Omega_0, W(t, \omega)$  is continuous with regard to  $t$ .

Today we will give another characterization of Brownian motions in terms of the *Gaussian process*. First recall the notion of high dimensional Gaussian distribution. A vector of random variables  $(X_1, X_2, \dots, X_n)$  is said to be Gaussian iff  $\forall a_1, a_2, \dots, a_n, \sum_{i=1}^n a_i X_i$  is a one-dimensional Gaussian. Let  $\mu = (\mu_1, \mu_2, \dots, \mu_n)$  where  $\mu_i = \mathbb{E}[X_i]$ . Let  $\Sigma = (\text{Cov}(X_i, X_j))_{i,j}$ . Then the probability density function  $f$  of  $(X_1, X_2, \dots, X_n)$  is

$$\text{for } x = (x_1, x_2, \dots, x_n), f(x) = (2\pi)^{-\frac{n}{2}} \cdot |\det \Sigma|^{-\frac{1}{2}} \cdot e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}.$$

**Definition 2** (Gaussian Process). A stochastic process  $\{X(t)\}_{t \geq 0}$  is called Gaussian process if  $\forall 0 \leq t_1 \leq t_2 \leq \dots \leq t_n, (X(t_1), X(t_2), \dots, X(t_n))$  is Gaussian.

Note that a Gaussian vector can be characterized by its mean vector and the covariance matrix. Standard Brownian motion is a special family of Gaussian processes where the covariance of  $X(s)$  and  $X(t)$  is  $s \wedge t$ .

**Definition 3** (Standard Brownian Motion/Standard Wiener Process). We say a stochastic process  $\{W(t)\}_{t \geq 0}$  is a standard Brownian motion or Wiener process if it satisfies

- $\{W(t)\}_{t \geq 0}$  is an almost surely continuous Gaussian Process;
- $\forall s \geq 0, \mathbb{E}[W(s)] = 0$ ;
- $\forall 0 \leq s \leq t, \text{Cov}(W(s), W(t)) = s$ .

We will show that it is easier to use Definition 3 to verify that a certain stochastic process is a Brownian motion. Let us first verify that the two definitions are equivalent.

**Proposition 4.** *The two definitions of standard Brownian motions are equivalent.*

*Proof.* Given Definition 1, it is easy to know that  $E[W(s)] = 0$  for all  $s \geq 0$  since  $W(s) \sim \mathcal{N}(0, s)$ . What we need is to verify that  $\{W(t)\}_{t \geq 0}$  in Definition 1 is a Gaussian process and to compute the covariance of  $W(s)$  and  $W(t)$  in Definition 1.

Note that  $\forall 0 \leq s < t$  and  $\forall a, b$ , we have

$$aW(s) + bW(t) = (a + b)W(s) + b(W(t) - W(s)).$$

Since  $W(s)$  and  $W(t) - W(s)$  are two independent Gaussian's,  $aW(s) + bW(t)$  is still a Gaussian.

By the distributive law of covariance, for any  $0 \leq s \leq t$ , we have

$$\begin{aligned} \text{Cov}(W(s), W(t)) &= \text{Cov}(W(s), W(t) - W(s) + W(s)) \\ &= \text{Cov}(W(s), W(t) - W(s)) + \text{Cov}(W(s), W(s)) \\ &= \text{Var}[W(s)] = s. \end{aligned}$$

Then we consider the counterpart. Given Definition 3, we can deduce the first and fourth property in Definition 1 directly. For any  $0 \leq t_{i-1} \leq t_i \leq t_{j-1} \leq t_j$ , we have

$$\begin{aligned} &\text{Cov}(W(t_i) - W(t_{i-1}), W(t_j) - W(t_{j-1})) \\ &= \text{Cov}(W(t_i), W(t_j)) + \text{Cov}(W(t_{i-1}), W(t_{j-1})) \\ &\quad - \text{Cov}(W(t_i), W(t_{j-1})) - \text{Cov}(W(t_{i-1}), W(t_j)) \\ &= t_i + t_{i-1} - t_i - t_{i-1} = 0, \end{aligned}$$

which yields the independence of  $W(t_i) - W(t_{i-1})$  and  $W(t_j) - W(t_{j-1})$ . Thus, the  $\{W(t)\}_{t \geq 0}$  in Definition 3 satisfies independent increments.

It is easy to verify that  $\forall s, t > 0$ ,  $W(s + t) - W(s)$  is a Gaussian with mean 0. Note that

$$\begin{aligned} \text{Var}[W(t + s) - W(s)] &= E[(W(t + s) - W(s))^2] \\ &= E[W(t + s)^2] + E[W(s)^2] - 2E[W(t + s)W(s)] \\ &= \text{Var}[W(t + s)] + \text{Var}[W(s)] - 2\text{Cov}(W(t + s), W(s)) \\ &= t + s + s - 2s = t. \end{aligned}$$

Thus, the  $\{W(t)\}_{t \geq 0}$  in Definition 3 satisfies stationary increments.  $\square$

**Example 1.** Suppose  $\{W(t)\}_{t \geq 0}$  is a standard Brownian motion. We claim that  $\{X(t)\}_{t \geq 0}$  is also a standard Brownian motion where  $X(0) = 0$  and  $X(t) = t \cdot W(\frac{1}{t})$  for  $t > 0$ .

We verify the three requirements in Definition 3.

Since  $X(t) = t \cdot W(\frac{1}{t})$  which is the composition of two (almost surely) continuous function,  $\{X(t)\}_{t \geq 0}$  is continuous almost surely as well. For any

It is worth noting that the sum of two Gaussians is not necessarily a Gaussian, unless they are joint Gaussian. Independence is just a special case of joint Gaussian (the covariance is zero).

$a_1, a_2, \dots, a_n$  and  $t_1, t_2, \dots, t_n \geq 0$ ,  $\sum_{i=1}^n a_i X(t_i) = \sum_{i=1}^n a_i t_i \cdot W(\frac{1}{t_i})$ . Since  $\{W(t)\}$  is standard Brownian motion,  $\sum_{i=1}^n a_i t_i \cdot W(\frac{1}{t_i})$  is Gaussian. Thus,  $\{X(t)\}_{t \geq 0}$  is a Gaussian process. For  $0 \leq s < t$ ,

$$\begin{aligned} \text{Cov}(X(s), X(t)) &= \text{Cov}(sW(\frac{1}{s}), tW(\frac{1}{t})) \\ &= st \cdot \text{Cov}(W(\frac{1}{s}), W(\frac{1}{t})) \\ &= st \cdot \frac{1}{t} = s. \end{aligned}$$

Thus,  $\{X(t)\}_{t \geq 0}$  is a standard Brownian motion.

## 2 Brownian Bridge

In the last lecture, we already calculated the distribution of  $W(t)$  conditioned on  $W(u) = x$  for some  $u \geq t$ . We use  $X(t)$  to denote this process, and  $X(t)$  is usually called a *Brownian bridge*.

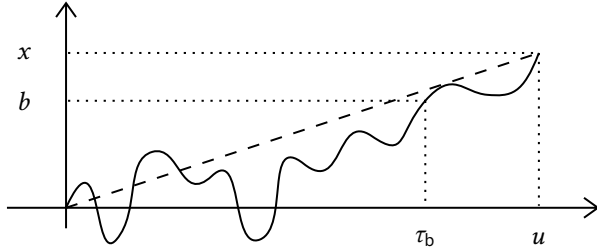


Figure 1: A Brownian bridge

We know from previous calculations that  $X(t) \sim N(\frac{t}{u}x, \frac{t(u-t)}{u})$  is a Gaussian. Since the *conditional distribution of a multidimensional Gaussian distribution is Gaussian as well*,  $X(t)$  is a Gaussian process. As a result, it is useful to compute the covariance of this process.

Recall  $W(t)$  is a standard Brownian motion. For any  $s \leq t$ , we have

$$\begin{aligned} &\text{Cov}(X(s), X(t)) \\ &= \text{Cov}(W(s), W(t) \mid W(u) = x) \\ &= \mathbb{E}[W(s) \cdot W(t) \mid W(u) = x] - \mathbb{E}[W(s) \mid W(u) = x] \cdot \mathbb{E}[W(t) \mid W(u) = x] \\ &= \int_{-\infty}^{\infty} y \mathbb{E}[W(s) \mid W(t) = y, W(u) = x] \cdot f_{W(t) \mid W(u)}(y \mid x) dy - \frac{st}{u^2} x^2. \\ &= \frac{s}{t} \mathbb{E}[W(t)^2 \mid W(u) = x] - \frac{st}{u^2} x^2 \\ &= \frac{s(u-t)}{u}. \end{aligned}$$

$$\mathbb{E}[W(t)^2 \mid W(u) = x] = \text{Var}[X(t)] + \mathbb{E}[X(t)]^2$$

**Definition 5** (Standard Brownian Bridge). *A standard Brownian motion ending at  $W(1) = 0$  is called a standard Brownian bridge.*

We can verify that  $X(t) = W(t) - tW(1)$  is a standard Brownian bridge by calculating its mean and covariances.

Again like we did in the last lecture, we can compute the hitting time of a standard Brownian bridge using the principle of reflection.

**Example 2** (Hitting Time in a Brownian Bridge). Let  $\{W(t)\}_{t \geq 0}$  be a standard Brownian motion. Let  $\tau_b \triangleq \inf \{t \geq 0 \mid W(t) > b\}$ . Then we compute  $\Pr[\tau_b < u \mid W(u) = x]$ . Note that if  $b < x$ ,  $\Pr[\tau_b < u \mid W(u) = x] = 1$ . Let  $\psi$  be the probability density function of standard Gaussian distribution, that is,  $\psi(x) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{x^2}{2}}$ . If  $b > x$ , letting  $dx$  be an infinitesimal change, we have

$$\begin{aligned} \Pr[\tau_b < u \mid W(u) = x] &= \frac{\Pr[\tau_b < u \wedge W(u) \in [x, x + dx]]}{\Pr[W(u) \in [x, x + dx]]} \\ &= \frac{\Pr[\tau_b < u] \cdot \Pr[W(u) \in [x, x + dx] \mid \tau_b < u]}{f_u(x) dx}. \end{aligned}$$

If we have known the value of  $\tau_b$  and  $\tau_b < u$ , we can regard  $\{W(u)\}_{t \geq \tau_b}$  as a Brownian motion starting from  $b$ . Then we have

$$\begin{aligned} \Pr[\tau_b < u] \cdot \Pr[W(u) \in [x, x + dx] \mid \tau_b < u] &= \Pr[\tau_b < u] \cdot \Pr[W(u) \in [2b - x - dx, 2b - x] \mid \tau_b < u] \\ &= \Pr[\tau_b < u \wedge W(u) \in [2b - x - dx, 2b - x]] \\ &= \Pr[W(u) \in [2b - x - dx, 2b - x]] \\ &= f_u(2b - x) dx \end{aligned}$$

Thus, when  $b > x$ ,  $\Pr[\tau_b < u \mid W(u) = x] = \frac{f_u(2b-x)}{f_u(x)} = e^{-\frac{2b(b-x)}{u}}$ .

When  $b = x$ , we have

$$\Pr[\tau_b < u \mid W(u) = b] = \frac{\Pr[\tau_b < u \wedge W(u) \in [b, b + db]]}{\Pr[W(u) \in [b, b + db]]}.$$

Note that

$$\Pr[\tau_b < u \wedge W(u) \in [b, b + db]] = \Pr[\tau_b < u] - \Pr[\tau_b < u \wedge W(u) > b + db] - \Pr[\tau_b < u \wedge W(u) < b]. \quad (1)$$

We know that  $\Pr[\tau_b < u] = 2 \left(1 - \Phi\left(\frac{b}{\sqrt{u}}\right)\right)$ . Note that

$$\begin{aligned} \Pr[\tau_b < u \wedge W(u) > b + db] &= \Pr[W(u) > b + db] \\ &= 1 - \Phi\left(\frac{b}{\sqrt{u}}\right) - \Pr[W(u) \in [b, b + db]]. \end{aligned}$$

And

$$\begin{aligned} \Pr[\tau_b < u \wedge W(u) < b] &= \Pr[\tau_b < u] \cdot \Pr[W(u) < b \mid \tau_b < u] \\ &= \frac{1}{2} \cdot \Pr[\tau_b < u] = 1 - \Phi\left(\frac{b}{\sqrt{u}}\right). \end{aligned}$$

Thus, Equation (1) equals to  $\Pr[W(u) \in [b, b + db]]$  and

$$\Pr[\tau_b < u \mid W(u) = b] = 1.$$



### 3 Kolmogorov-Smirnov Test

In this section, we introduce an application of Brownian Bridge, the Kolmogorov-Smirnov test.

Suppose that  $U_1, U_2, \dots, U_n$  are independently sampled from some distribution  $[0, 1]$  with CDF  $F$ . We would like to check if it is a uniform distribution, i.e., if the  $F$  satisfies  $F(t) = t$  for every  $t \in [0, 1]$ .

Let  $\widehat{F}_n$  be the empirical cumulative distribution function, that is, for  $t \in [0, 1]$ ,  $\widehat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}[U_i \leq t]$ . It then follows from the law of large numbers that

$$\widehat{F}_n(t) \xrightarrow{n \rightarrow \infty} \mathbf{E} [\widehat{F}_n(t)] = \frac{1}{n} \sum_{i=1}^n \Pr [U_i \leq t] = F(t).$$

The idea of Kolmogorov-Smirnov test is to monitor the variable  $\widehat{F}_n(t) - t$  for every  $t \in [0, 1]$  and reject the uniformity hypothesis if there exists some  $t$  that  $|\widehat{F}_n(t) - t|$  is large. Then our goal is to find a suitable rejection threshold  $b$  such that if  $F$  is indeed a uniform distribution, the failure probability  $\lim_{n \rightarrow \infty} \Pr \left[ \max_{t \in [0, 1]} |\widehat{F}_n(t) - t| \geq b \right]$  is sufficiently small (i.e.,  $\leq \frac{1}{100}$ ). If  $F$  is a uniform distribution, for a fixed  $t$ , we have

$$\begin{aligned} \mathbf{E} [\widehat{F}_n(t)] &= F(t) = t; \\ \text{Var} [\widehat{F}_n(t)] &= \frac{1}{n^2} \sum_{i=1}^n \text{Var} [\mathbf{1}[U_i \leq t]] = \frac{1}{n} \cdot t(1-t). \end{aligned}$$

Let  $X_n(t) \triangleq \sqrt{n} \cdot (\widehat{F}_n(t) - t)$  for  $t \in [0, 1]$ . By the Central Limit Theorem, we have  $X_n(t) \sim \mathcal{N}(0, t(1-t))$  when  $n \rightarrow \infty$ . For any  $0 \leq s \leq t \leq 1$ ,

$$\begin{aligned} \text{Cov}(X_n(s), X_n(t)) &= n \cdot \text{Cov}(\widehat{F}_n(s) - s, \widehat{F}_n(t) - t) \\ &= \frac{1}{n} \text{Cov} \left( \sum_{i=1}^n \mathbf{1}[U_i \leq s], \sum_{i=1}^n \mathbf{1}[U_i \leq t] \right) \\ &= \text{Cov}(\mathbf{1}[U_1 \leq s], \mathbf{1}[U_1 \leq t]) \\ &= \Pr[U_1 \leq s, U_1 \leq t] - \Pr[U_1 \leq s] \Pr[U_1 \leq t] \\ &= s(1-t). \end{aligned}$$

For any  $0 \leq t_1 \leq t_2 \leq \dots \leq t_k \leq 1$ , let  $\Sigma = (\text{Cov}(X_n(t_i), X_n(t_j)))_{i,j}$ . It follows from the high-dimensional Central Limit Theorem that

$$(X_n(t_1), X_n(t_2), \dots, X_n(t_k))^T \xrightarrow{D} \mathcal{N}(\mathbf{0}, \Sigma) \sim (X(t_1), X(t_2), \dots, X(t_k))^T,$$

where  $\{X(t)\}$  is a standard Brownian Bridge. Then using the result in Example 2, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr \left[ \max_{t \in [0, 1]} \widehat{F}_n(t) - t \geq b \right] &= \Pr \left[ \max_{t \in [0, 1]} X(t) \geq \sqrt{nb} \right] \\ &= \Pr \left[ \tau_{\sqrt{nb}} < 1 \mid W(1) = 0 \right] = \exp\{-2nb^2\}. \end{aligned}$$