

# Probability and Statistics 2

Data Science Engineering

Session 3: Structure of chains and classification of states.

---

Class structure. Irreducibility. Chain decomposition. Recurrence and transience. Positive and null recurrence. Absorption probabilities and expected absorption time. Aperiodicity.

---

## 1. CLASS STRUCTURE AND IRREDUCIBILITY

Given a Markov chain we can decompose it into smaller chains that are easier to analyse. Throughout  $\{X_n, n \geq 0\}$  denotes a Markov chain with state space  $S = \{0, 1, \dots, m\}$  if  $S$  is finite and  $S = \mathbb{N}$  or  $S = \mathbb{Z}$  if it is infinite.

**Definition 1.1.** We say that state  $i$  leads to  $j$  if there exists  $n$  such that  $p_{ij}(n) > 0$ , and we denote it by  $i \rightarrow j$ . Two states  $i, j$  communicate if  $i \rightarrow j$  and  $j \rightarrow i$ , and we denote it by  $i \leftrightarrow j$ .

These definitions are easier to understand if one looks at the graphical representation of the chain. Then  $i \rightarrow j$  if and only if there exists a directed path from  $i$  to  $j$ .

It should be apparent that  $\leftrightarrow$  is an equivalence relation; that is, it satisfies the following properties:

- i) reflexive:  $i \leftrightarrow i$ ,
- ii) symmetric: if  $i \leftrightarrow j$ , then  $j \leftrightarrow i$
- iii) transitive: if  $i \leftrightarrow j$  and  $j \leftrightarrow k$ , then  $i \leftrightarrow k$ .

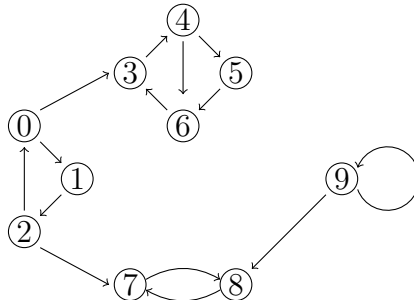
The equivalence relation  $\leftrightarrow$  partitions  $S$  into *communication classes* (or simply *classes*). Note that classes can also be composed of single states. A class  $C$  is *closed* if  $i \in C$  and  $i \rightarrow j$ , implies that  $j \in C$ . Closed classes are state subsets that the chain cannot leave once it has entered them. A state  $i$  is *absorbing* if its class is closed and only contains state  $i$ .

**Proposition 1.2.** The state space  $S$  can be univocally partitioned

$$S = C_1 \cup \dots \cup C_r \cup T,$$

where  $C_1, \dots, C_r$  is the set of closed classes.

**Example 1.3.** The picture below shows the graphical representation of a Markov chain. For instance, we have  $9 \rightarrow 7$  and  $3 \leftrightarrow 6$ . The classes are  $\{0, 1, 2\}$ ,  $\{2, 4, 5, 6\}$ ,  $\{7, 8\}$  and  $\{9\}$ . The only closed classes are the second and third ones. Thus, the desired partition is  $C_1 = \{3, 4, 5, 6\}$ ,  $C_2 = \{7, 8\}$ ,  $T = \{0, 1, 2, 9\}$ .



□

With the partition defined above it should be clear that, in the long run, the states  $T$  will be abandoned (with probability one) and the chain will get stuck in one of the closed classes  $C_i$ . This motivates the following definition.

**Definition 1.4.** A chain is irreducible if it has a single communication class. If a chain is not irreducible, we call it reducible.

In applications, it is generally desirable that the chain is irreducible as this implies nice properties of the chain that we will see in the next chapter.

**Example 1.5.** Let us revisit two of the examples we saw in the previous chapter.

- i) *Simple random walk on  $\mathbb{Z}$ :* Note that  $i \leftrightarrow j$  for every  $i, j \in \mathbb{Z}$ . Thus, the chain is irreducible.
- ii) *Gambler's ruin:* The states 0 and  $N$  are absorbing states. There is another class composed of  $\{1, 2, \dots, N-1\}$ . So the chain is reducible. We will analyse this example with more detail at the end of this chapter.

## 2. RECURRENCE AND TRANSIENCE

One classical problem in the study of Markov chains is the analysis of the first time the chain visits a given state. This leads to a study of the nature of the states.

For every  $i \in S$  define the *hitting time* to  $j$  starting at  $i$  as

$$T_{ij} = \min\{n \geq 1 : X_n = j | X_0 = i\},$$

That is, the first time we hit  $j$  if the chain starts at  $i$ . When  $i = j$ , we talk about the *return time* and we denote it by  $T_i$ . Since  $T_{ij}$  is a random variable, we denote its probability

distribution by

$$f_{ij}(n) = \Pr(T_{ij} = n).$$

**Definition 2.1.** A state  $i \in S$  is recurrent if a chain starting at  $i$  returns to  $i$  with positive probability; that is,

$$\Pr(T_i < \infty) = 1.$$

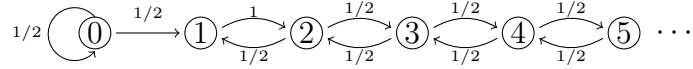
The state  $i$  is transient if it is not recurrent, namely, there is a positive probability that  $i$  is never visited again after visiting it,

$$\Pr(T_i < \infty) < 1.$$

Note that the event  $\{T_i < \infty\}$  is the union of the disjoint events  $\{T_i(n), n \geq 1\}$ . To verify if a state is recurrent/transient, it is convenient to write

$$\Pr(T_i < \infty) = \sum_{n \geq 1} f_{ii}(n)$$

**Example 2.2.** Consider the Markov chain with graphical representation depicted below. The state 0 is transient as there is a positive probability that the chain never returns to 0 once in that state.



**Theorem 2.3.** A state  $i$  is recurrent if and only if

$$\sum_{n \geq 1} p_{ii}(n) = \infty.$$

*Proof.* Starting at  $i$ , let  $N_i$  be the total number of visits at  $i$ . Let  $p$  be the probability to return to  $i$ , that is  $p = \sum_{n \geq 1} f_{ii}(n)$ . Each time we visit  $i$ , we can restart the time counter, so the probability of returning to  $i$  after visiting it is still  $p$ . So the random variable  $N_i$  is a geometric random variable with probability  $q = 1 - p$ . It has expected value  $\mathbb{E}[N_i] = 1/q = 1/(1 - p)$ . If  $i$  is recurrent,  $p = 1$  and  $\mathbb{E}[N_i] = \infty$ . If  $i$  is transient,  $p < 1$  and  $\mathbb{E}[N_i] < \infty$ .

Let  $I_n$  be the indicator random variable of the event  $\{X_n = i\}$ . We have

$$\mathbb{E}[N_i] = \sum_{n \geq 1} I_n = \sum_{n \geq 1} p_{ii}(n).$$

It follows that  $\sum_{n \geq 1} p_{ii}(n) = \infty$  if and only if  $i$  is recurrent.  $\square$

We now discuss an application of Theorem 2.3 to the analysis of the symmetric random walk on  $\mathbb{Z}^d$ . Recall the definition given in Chapter 2. The chain has state space  $\mathbb{Z}^d$  and, at every step, one of the  $2d$  possible directions is chosen, each one with probability  $1/2d$ . The following theorem determines under which conditions the random walk is recurrent.

**Theorem 2.4** (Polya 1921). *The origin in the symmetric random walk on  $\mathbb{Z}^d$  is recurrent if and only if  $d \leq 2$ .*

*Proof.* We first proof the case  $d = 1$ . The probability that the walk is at 0 after  $2n$  steps is precisely  $p_{00}(2n) = \binom{2n}{n} 2^{-2n}$ , as we need to choose the same number of right and left moves and every path of length  $2n$  has probability  $2^{-2n}$  to appear. Stirling formula is  $m! \sim \sqrt{2\pi m} (m/e)^m$  and implies that  $\binom{2n}{n} \sim 2^{2n} / \sqrt{\pi n}$ . Therefore, we have

$$\sum_{n \geq 1} p_{00}(2n) = \sum_{n \geq 1} \binom{2n}{n} 2^{-2n} \sim \sum_{n \geq 1} \frac{1}{\sqrt{\pi n}} = \infty,$$

so the origin is recurrent.

For  $d = 2$ , change the steps of the random walker. Instead of  $(1, 0), (-1, 0), (0, 1), (0, -1)$ , consider the steps  $(1, 1), (-1, 1), (-1, -1), (1, -1)$ . Therefore, one can see any path of length  $2n$  as two independent unidimensional paths of length  $n$ , one for each dimension. The probability of returning to the origin with this new set of steps, is the same as for the original set of steps. Using what we proved in the case  $d = 1$ ,

$$\sum_{n \geq 1} p_{00}(2n) = \sum_{n \geq 1} \left( \binom{2n}{n} 2^{-2n} \right)^2 \sim \sum_{n \geq 1} \frac{1}{\pi n} = \infty,$$

so the origin is still recurrent.

For  $d = 3$ , the proof is slightly more involved. One can compute  $p_{00}(2n)$  as follows:

- (1) choose what steps will be positive in each dimension, there are  $\binom{2n}{n}$  options,
- (2) choose how many steps we will take in each dimension, say  $i, j, k$  with  $i + j + k = n$ ,
- (3) choose which of the  $n$  positive and  $n$  negative steps will be moving in each dimension, there are  $\binom{n}{i, j, k}^2$  options.

Since each path of length  $2n$  has probability  $6^{-2n}$  to appear, we have

$$p_{00}(2n) = 6^{-2n} \binom{2n}{2} \sum_{i+j+k=n} \binom{n}{i, j, k}^2 \leq 6^{-2n} \binom{2n}{2} \left( \max_{i+j+k=n} \binom{n}{i, j, k} \right) \sum_{i+j+k=n} \binom{n}{i, j, k}$$

On the one hand  $\sum_{i+j+k=n} \binom{n}{i, j, k} = 3^n$ . On the other hand

$$\max_{i+j+k=n} \binom{n}{i, j, k} = \binom{n}{n/3, n/3, n/3} = \frac{n!}{(n/3)!^3} \sim \frac{C \cdot 3^n}{n}$$

for some constant  $C$ , one obtains  $p_{00}(n) \sim C n^{-3/2}$ . In the previous equation we have assumed that  $n$  is divisible by 3 so  $n/3$  is an integer. If not, one should replace  $n/3$  by  $\lfloor n/3 \rfloor$  or  $\lceil n/3 \rceil$  accordingly, but the approximation still holds.

We conclude that

$$\sum_{n \geq 1} p_{00}(n) \sim \sum_{n \geq 1} \frac{C}{n^{3/2}} < \infty,$$

so the origin is transient. For  $d \geq 4$ , the origin will still be transient, as the projected random walk on  $\mathbb{Z}^3$  induced by the first three components of  $\mathbb{Z}^d$  behaves like the case  $d = 3$  and the origin is transient.  $\square$

We now relate the classification of states with the previous section.

**Theorem 2.5.** *Given the partition  $S = C_1 \cup \dots \cup C_r \cup T$  from Proposition 1.2, a state  $i$  is recurrent if and only if  $i \in C_j$  for some  $j$ .*

In words, closed classes correspond to recurrent states.

**Definition 2.6.** *The mean return time of a state  $i$  is the expected number of steps until the chain returns to  $i$  for the first time when starting at  $i$ , that is*

$$\mu_i = \mathbb{E}(T_i) = \sum_{n \geq 1} n f_{ii}(n),$$

*. This is infinite if  $i$  is transient and can be finite or infinite if  $i$  is recurrent. A recurrent state  $i$  is positive-recurrent if  $\mu_i < \infty$  and null-recurrent if  $\mu_i = \infty$ .*

If  $S$  is finite and  $i$  is recurrent, we expect to visit the starting state  $i$  in finite time and all recurrent states are also positive-recurrent. So null-recurrence is only interesting in the case of infinite state space chains.

It is intuitively clear that, for a recurrent state  $i$  we have

$$p_{ii}(n) \rightarrow \frac{1}{\mu_i}, \quad (n \rightarrow \infty).$$

**Theorem 2.7.** *A recurrent state  $i$  is null-recurrent if and only if  $p_{ii}(n) \rightarrow 0$  as  $n \rightarrow \infty$ .*

### 3. GAMBLER'S RUIN REVISITED

If a chain is not irreducible, given an initial state one may ask

- a) **hitting probability**: what is the probability of eventually ending up in a given closed class (including absorbing states)?
- b) **expected hitting time**: what is the average time until we hit a closed class?

There are general techniques to solve the above two questions, but we will explain how to answer a) and b) using the example of the Gambler's ruin. Later in the exercises, we will see other examples.

Recall that the Gambler's ruin chain  $\{X_n, n \geq 0\}$  has state space  $S = \{0, 1, \dots, N\}$ , where  $N = a + b$  is the total capital of gamblers  $A$  and  $B$  and  $X_n$  denotes the capital of player  $A$  at time  $n$ . For the sake of simplicity, let us assume that the bet is fair; that is,  $p = 1/2$ .

According to the above classification, the states  $\{1, \dots, N-1\}$  are transient while 0 and  $N$  are recurrent, even more, they are absorbing.

In this setting, a) asks to compute the probability to end up at 0 or at  $N$  when starting at  $X_0 = i$ . Denote by  $q_i$  the probability that  $X_n$  ends up at 0 when  $X_0 = i$ . We certainly have

$$q_0 = 1 \text{ and } q_N = 0.$$

Moreover, by playing once, we see that, for  $0 < i < N$ ,

$$q_i = \frac{1}{2}(q_{i+1} + q_{i-1}).$$

As an exercise you can prove by induction that  $q_i = 1 - i/N$  satisfies the previous equations.

This argument can be adapted for all  $p \in (0, 1)$ , giving the following result.

**Theorem 3.1.** *Let  $\{X_n, n \geq 0\}$  be a random walk with barriers at 0 and  $N$  with probability of a step up  $p \in (0, 1)$ . The probability that the chain ends up at 0 if  $X_0 = i$  is*

$$q_i = \begin{cases} 1 - i/N, & p = 1/2 \\ \frac{b^i - b^N}{1 - b^N}, & p \neq 1/2. \end{cases}$$

where  $b = (1 - p)/p$ .

What happens at a casino, where  $B$  has a huge capital? Given a fix capital  $a$  for gambler  $A$ , as  $N \rightarrow \infty$  we see that the probability of ruin goes to one as long as  $p \leq 1/2$ . So even in a fair game we expect the gambler to get ruined.

Question b) on the average time till absorption can also be treated. Let  $M_i$  denote the time till the chain reaches any absorbing state if  $X_0 = i$ . By the same argument as before, we have  $\mathbb{E}(M_0) = \mathbb{E}(M_N) = 0$  and for  $p = 1/2$

$$\mathbb{E}(M_i) = 1 + \frac{1}{2}(\mathbb{E}(M_{i+1}) + \mathbb{E}(M_{i-1})),$$

By induction one can check that  $\mathbb{E}(M_i) = i(N - i)$ . For any  $p \in (0, 1)$  we have the following result.

**Theorem 3.2.** *Let  $\{X_n, n \geq 0\}$  be a random walk with barriers at 0 and  $N$  with probability of a step up  $p$ . The average time till the chain ends up at 0 or  $N$  if  $X_0 = i$  is*

$$\mathbb{E}(M_i) = \begin{cases} i(N - i), & p = 1/2 \\ \frac{1}{1 - 2p} \left( i - \frac{1 - b^i}{1 - b^N} N \right), & p \neq 1/2. \end{cases}$$

where  $b = (1 - p)/p$ .

Again, if  $N \rightarrow \infty$  and the game is fair ( $p = 1/2$ ) while the probability of ruin is 1 by Theorem 3.1, the time until ruin is infinite!

#### 4. APERIODICITY

We conclude with another notion that will be important in the next chapter

**Definition 4.1.** *The period of a recurrent state  $i$  is*

$$d(i) = \gcd\{n \geq 1 \mid p_{ii}(n) > 0\},$$

*the greatest common divisor of the times where return to  $i$  from  $i$  is possible.*

Note that the set of numbers from which we take the gcd is non-empty, as the state is recurrent.

**Example 4.2.** *Consider a random walk on a circle of length  $2m$ . If the chain starts at a state  $i$  then it can only return to  $i$  after an even number of steps. Moreover, it is possible to return on two steps. So the period of any state is 2.*

**Definition 4.3.** *A chain is aperiodic if and only if every recurrent state has period 1.*

A sufficient condition for a state to have period 1 is that  $p_{ii} > 0$ ; that is, state  $i$  has a loop. For an irreducible chain, if there exists a state with period 1, then the chain is aperiodic.

#### 5. EXERCISES AND PROBLEMS

- (1) Find the communication classes of the following chains, and identify the closed ones. Classify the transient and recurrent states and compute its periodicity.

$$\begin{pmatrix} 1/2 & 1/2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1/3 & 0 & 0 & 1/3 & 1/3 & 0 \\ 0 & 0 & 0 & 1/2 & 1/2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \quad \begin{pmatrix} 1/2 & 0 & 0 & 0 & 1/2 \\ 0 & 1/2 & 0 & 1/2 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1/4 & 1/4 & 1/4 & 1/4 \\ 1 & 0 & 0 & 0 & 0 \end{pmatrix}$$

- (2) In the following chains that appeared in previous chapters, identify the classes, classify the states (transient, positive/null-recurrent) and if recurrent, give their period:
- (a) The random walk with one and with two reflecting barriers.
  - (b) The random walk on a circle of  $m$  points.
  - (c) The Ehrenfest model and the modified Ehrenfest model.
  - (d) Gene frequencies chain.
  - (e) Tennis game chain.
  - (f) *Juego de la oca* chain.
- (3) Let  $\{X_n, n \geq 0\}$  be a Markov chain with transition matrix

$$P = \begin{pmatrix} 1-2p & 2p & 0 \\ p & 1-2p & p \\ 0 & 2p & 1-2p \end{pmatrix}$$

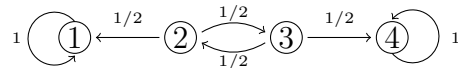
Classify the states of the chain.

- (4) An urn contains  $k$  green balls and  $k + 2$  red balls. A ball is picked at random. If it is green, a red ball is selected and both are removed. If it is red it is replaced together with an extra green ball. The process terminates when there are no green balls left. Show that the probability of termination is  $1/(k + 1)$ .
- (5) Consider the random walk on a circle of  $m$  points  $\{0, 1, \dots, m - 1\}$ . We start at 0. Use the gambler's ruin problem to show that the probability that 1 is the last point visited is  $1/m$ . The same result holds with 1 replaced by  $k$  for any  $1 \leq k \leq m - 1$ .
- (6) A Gambler has an initial capital of 100€. The Gambler plays a fair game ( $p = 1/2$ ) until it doubles its capital or gets ruined. Show that if their stake is 1€, the expected time until the end of the game is 4 times the expected time if the stake is 2€.
- (7) A Markov chain  $\{X_n, n \geq 0\}$  has state space  $\{0, 1, 2, \dots\}$  and transitions

$$p_{i,j} = \begin{cases} 2/3 & \text{if } j = i + 1 \\ 1/3 & \text{if } j = 0 \\ 0 & \text{otherwise} \end{cases}$$

Show that the chain is irreducible. Determine if 0 is transient or recurrent. Is it positive-recurrent?

- (8) \* A Markov chain  $\{X_n, n \geq 0\}$  has state space  $\{0, 1, 2, \dots\}$  and transitions from each  $i > 0$  to  $i + 1$  with probability  $1 - \frac{1}{2i^\alpha}$  and to 0 with probability  $\frac{1}{2i^\alpha}$ . It transitions from 0 to 1 with probability 1. Show that the chain is irreducible. Determine if 0 is transient or recurrent depending on  $\alpha$ .
- (9) Consider a symmetric random walk  $\{X_n, n \geq 0\}$  with  $X_0 = 0$ . For  $m > 0$ , let  $Y_m$  be the number of visits to  $m$  before the walk returns to 0 for the first time.
  - (a) By using the gambler's ruin problem show that  $\Pr(Y_m > 0 | X_1 = 1) = 1/m$ . Deduce that  $\Pr(Y_m > 0) = \frac{1}{2m}$ .
  - (b) By using the gambler's ruin problem show that  $Y_m | Y_m > 0$  follows a geometric law  $Geom(1/2m)$ .
  - (c) Deduce that  $\mathbb{E}(Y_m) = 1$  (for all  $m > 0$ , a surprising result!).
  - (d) Prove that the digraph of communication classes does not contain directed cycles.
- (10) Consider the Markov chain



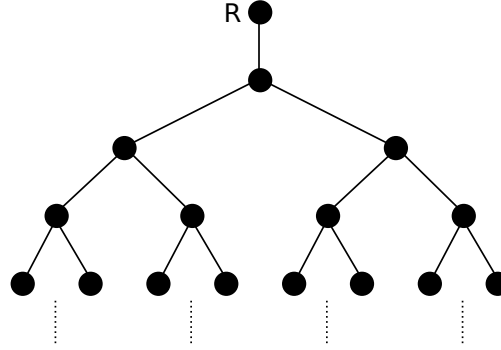
Starting at 2, compute the probability to hit 4 and the expected time until it happens.

- (11) A gambler has 1€ and wants to obtain 5€. He plays a game where a fair coin is tossed: if heads, he gets twice the amount of his stake, and if tails, he loses his stake. The gambler can choose between two strategies:
  - he stakes 1€ at every step.
  - he stakes all his money if he has at most 2€, and otherwise he only stakes just enough to get a capital of 5€.



What is the strategy that maximises the probability of getting his goal? Can you find a better strategy?

- (12) (Random walk on infinite trees) Consider an infinite binary random tree with a selected vertex  $R$  which comes with a single edge, and every other vertex of degree 3. The random walk on  $T$  starts at  $R$  and at every time it jumps to a random neighbour of the current vertex. Show that the random walk is transient. (Hint: use the Gambler's ruin problem)



- (13) A process moves on the integers  $\{1, 2, 3, 4, 5\}$ . It starts at 1 and, on each successive step, moves to an integer greater than its present position, moving with equal probability to each of the remaining larger integers. Find the expected number of steps to reach state five. Can you make a conjecture on the expected number of steps for the case  $\{1, 2, \dots, n\}$ ?
- (14) Recall the tennis game Markov chain designed in Chapter 1. Consider the game of tennis when *deuce* is reached; that is, you can restrict the chain to the states 1: A wins; 2: advantage A; 3: deuce; 4: advantage B; 5: B wins.
- Find the absorption probabilities for states 1 and 5.
  - At deuce, find the expected duration of the game.
- (15) Consider the random walk on an oriented torus. Fix  $a, b \in \mathbb{N}$ , the state space is

$$\{(x, y) \mid x \in \{0, \dots, a-1\}, y \in \{0, \dots, b-1\}\}.$$

The transitions are given by

$$\Pr(X_n = (x, y) \mid X_{n-1} = (x', y')) = \begin{cases} 1/2, & x = x' + 1 \text{ and } y = y' \\ 1/2, & x = x' \text{ and } y = y' + 1, \\ 0 & \text{otherwise} \end{cases}$$

Show that the chain is aperiodic if and only if  $\gcd(a, b) = 1$ .

- (16) \* Recall the tennis game Markov chain designed in Chapter 1. We will use it to predict the outcome of a game, where the players are Federer and Nadal. Here we have some stats of matches played between them (updated on July 2019):

Given the point probability, compute the probability that Federer wins a game. Is it a good approximation of the real one? You can use a mathematical software (e.g. R, Mathematica, Maple, Matlab, Sage ...) to solve the problem.

	points	% points	games	% games
Federer	3675	$p = 0.4935$	569	0.4806
Nadal	3772	$q = 0.5065$	615	0.5194