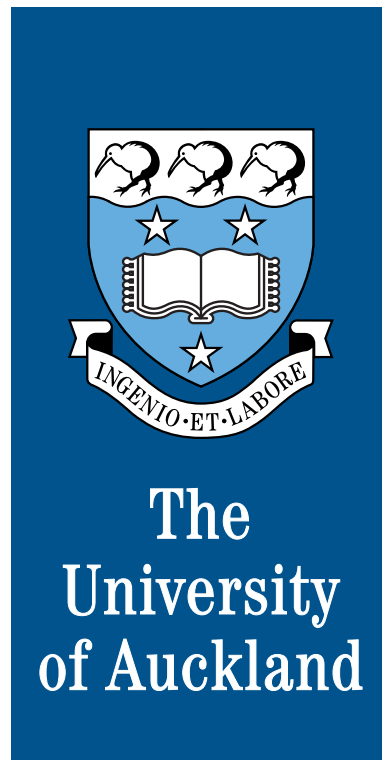


COURSE NOTES

STATS 325

Stochastic Processes



Department of Statistics
University of Auckland

Contents

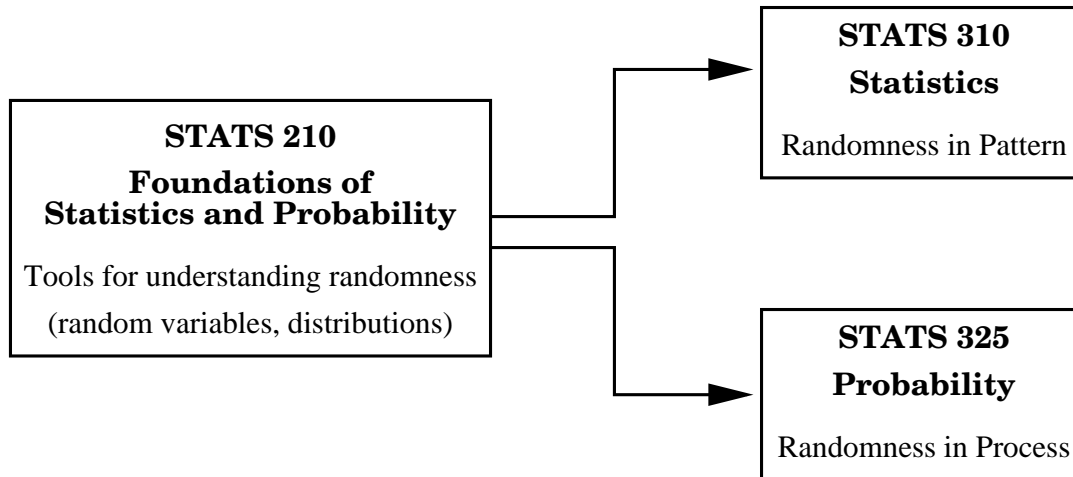
1. Stochastic Processes	4
1.1 Revision: Sample spaces and random variables	8
1.2 Stochastic Processes	9
2. Probability	16
2.1 Sample spaces and events	16
2.2 Probability Reference List	17
2.3 Conditional Probability	18
2.4 The Partition Theorem (Law of Total Probability)	23
2.5 Bayes' Theorem: inverting conditional probabilities	25
2.6 First-Step Analysis for calculating probabilities in a process	28
2.7 Special Process: the Gambler's Ruin	32
2.8 Independence	35
2.9 The Continuity Theorem	36
2.10 Random Variables	38
2.11 Continuous Random Variables	39
2.12 Discrete Random Variables	41
2.13 Independent Random Variables	43
3. Expectation and Variance	44
3.1 Expectation	45
3.2 Variance, covariance, and correlation	48
3.3 Conditional Expectation and Conditional Variance	51
3.4 Examples of Conditional Expectation and Variance	57
3.5 First-Step Analysis for calculating expected reaching times	63
3.6 Probability as a conditional expectation	67
3.7 Special process: a model for gene spread	71
4. Generating Functions	74
4.1 Common sums	74
4.2 Probability Generating Functions	76
4.3 Using the probability generating function to calculate probabilities	79
4.4 Expectation and moments from the PGF	81
4.5 Probability generating function for a sum of independent r.v.s	82
4.6 Randomly stopped sum	83

4.7	Summary: Properties of the PGF	89
4.8	Convergence of PGFs	90
4.9	Special Process: the Random Walk	95
4.10	Defective random variables	101
4.11	Random Walk: the probability we never reach our goal	105
5.	Mathematical Induction	108
5.1	Proving things in mathematics	108
5.2	Mathematical Induction by example	110
5.3	Some harder examples of mathematical induction	113
6.	Branching Processes:	116
6.1	Branching Processes	117
6.2	Questions about the Branching Process	119
6.3	Analysing the Branching Process	120
6.4	What does the distribution of \mathbf{Z}_n look like?	124
6.5	Mean and variance of \mathbf{Z}_n	128
7.	Extinction in Branching Processes	132
7.1	Extinction Probability	133
7.2	Conditions for ultimate extinction	139
7.3	Time to Extinction	143
7.4	Case Study: Geometric Branching Processes	146
8.	Markov Chains	149
8.1	Introduction	149
8.2	Definitions	151
8.3	The Transition Matrix	152
8.4	Example: setting up the transition matrix	154
8.5	Matrix Revision	154
8.6	The \mathbf{t} -step transition probabilities	155
8.7	Distribution of \mathbf{X}_t	157
8.8	Trajectory Probability	160
8.9	Worked Example: distribution of \mathbf{X}_t and trajectory probabilities	161
8.10	Class Structure	163
8.11	Hitting Probabilities	164
8.12	Expected hitting times	170
9.	Equilibrium	174
9.1	Equilibrium distribution in pictures	175
9.2	Calculating equilibrium distributions	176
9.3	Finding an equilibrium distribution	177
9.4	Long-term behaviour	179
9.5	Irreducibility	183

9.6	Periodicity	184
9.7	Convergence to Equilibrium	187
9.8	Examples	188
9.9	Special Process: the Two-Armed Bandit	192

Chapter 1: Stochastic Processes

What are Stochastic Processes, and how do they fit in?



Stats 210: laid the foundations of both Statistics and Probability: the tools for understanding randomness.

Stats 310: develops the theory for understanding *randomness in pattern*: tools for estimating parameters (maximum likelihood), testing hypotheses, modelling patterns in data (regression models).

Stats 325: develops the theory for understanding *randomness in process*. A process is a sequence of events where each step follows from the last after a random choice.

What sort of problems will we cover in Stats 325?

Here are some examples of the sorts of problems that we study in this course.

Gambler's Ruin

You start with \$30 and toss a fair coin repeatedly. Every time you throw a Head, you win \$5. Every time you throw a Tail, you lose \$5. You will stop when you reach \$100 or when you lose everything. What is the probability that you lose everything?

Answer: 70%.

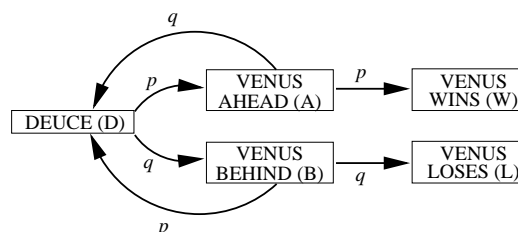


With permission, copyright Martin Ouellet

Winning at tennis

What is your probability of winning a game of tennis, starting from the even score Deuce (40-40), if your probability of winning each point is 0.3 and your opponent's is 0.7?

Answer: 15%.



Winning a lottery



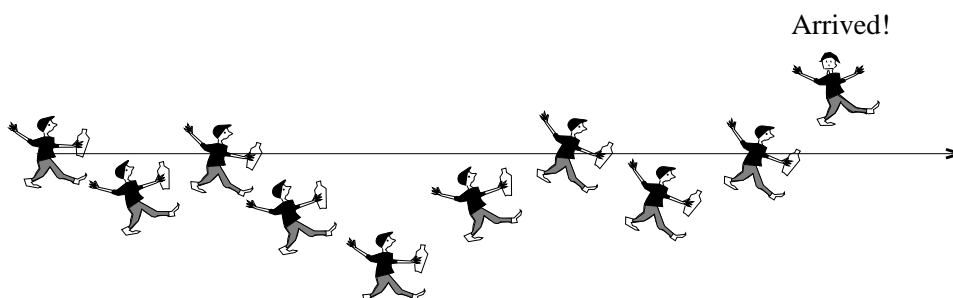
A million people have bought tickets for the weekly lottery draw. Each person has a probability of one-in-a-million of selecting the winning numbers. If more than one person selects the winning numbers, the winner will be chosen at random from all those with matching numbers.

You watch the lottery draw on TV and your numbers match the winning numbers!!! Only a one-in-a-million chance, and there were only a million players, so surely you will win the prize?

Not quite... What is the probability you will win? **Answer:** only 63%.

Drunkard's walk

A very drunk person staggers to left and right as he walks along. With each step he takes, he staggers one pace to the left with probability 0.5, and one pace to the right with probability 0.5. What is the expected number of paces he must take before he ends up one pace to the left of his starting point?



Answer: the expectation is infinite!

Pyramid selling schemes

Have you received a chain letter like this one? Just send \$10 to the person whose name comes at the top of the list, and add your own name to the bottom of the list. Send the letter to as many people as you can. Within a few months, the letter promises, you will have received \$77,000 in \$10 notes! Will you?

I WAS AMAZED WHEN I SAW HOW MUCH MONEY CAME FLOODING THROUGH MY LETTER BOX...I TURNED \$218 INTO \$78190 WITHIN THE FIRST 80 DAYS OF OPERATING THIS BUSINESS PLAN

**DO NOT BIN THIS IMMEDIATELY
THINK ABOUT IT FOR A FEW DAYS
FILE IN PENDING**

My name is David Rhodes and in September 1997 I lost my job. At the time I was living at the edge of my means and in debt. Consequently, this started a chain reaction that ended with the repossession of my home and car. If that wasn't enough several debt collectors were constantly hounding me. I imagine life looked bleak.

In January 1998 I received a letter telling me how to make my money. I ignored it because I was sceptical. However by March I was in debt. I finally realised that I had absolutely nothing to lose apart from that, I couldn't stop myself from thinking what if...

In the summer of 1999 my family and I went on a cruise and now Mercedes with cash and we are currently building our \$ home and I don't owe a single cent.

To date I have made over \$1,100,000. Even now as I write this it is hard to come to terms with the fact that like most people, I

**THIS IS HOW THE SYSTEM WORKS
WITHIN 60 DAYS**

You have sent off your \$10 note then mailed 200 letters (minimum) your details are printed at No5 on each of them. Your tasks are now complete. Sit back and relax- you deserve it.

If only 3% of 200 people respond to your letter, 6 people will mail 200 letters each = 1200 letters with your name at No4.

If only 3% of 1200 people respond to your letter, 36 people will mail 200 letters each = 7,200 letters with your name at No3.

If only 3% of 7,200 people respond to your letter, 216 people will mail 200 letters each = 43,200 letters with your name at No2.

If only 3% of 43,200 people respond to your letter, 1296 people will mail 200 letters each = 259,200 letters with your name at No1.

If only 3% of 259,200 people respond to their letters 7,776 people will send you \$10 each because your name is at No1 position therefore you will receive
\$77,760.00 in \$10 notes

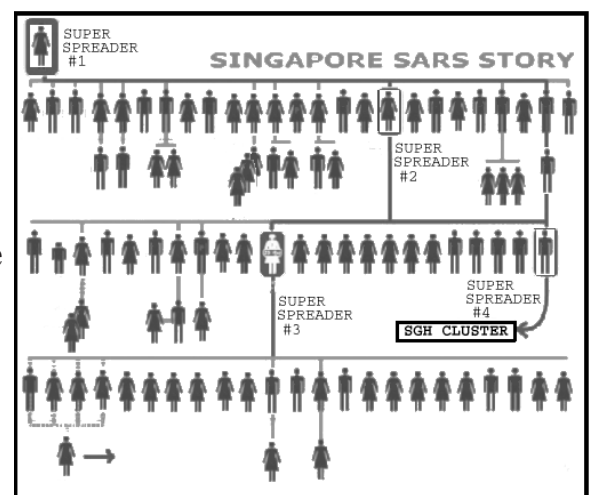
Answer: it depends upon the response rate. However, with a fairly realistic assumption about response rate, we can calculate an expected return of \$76 with a 64% chance of getting nothing!

Note: Pyramid selling schemes like this are prohibited under the Fair Trading Act, and it is illegal to participate in them.

Spread of SARS

The figure to the right shows the spread of the disease SARS (Severe Acute Respiratory Syndrome) through Singapore in 2003. With this pattern of infections, what is the probability that the disease eventually dies out of its own accord?

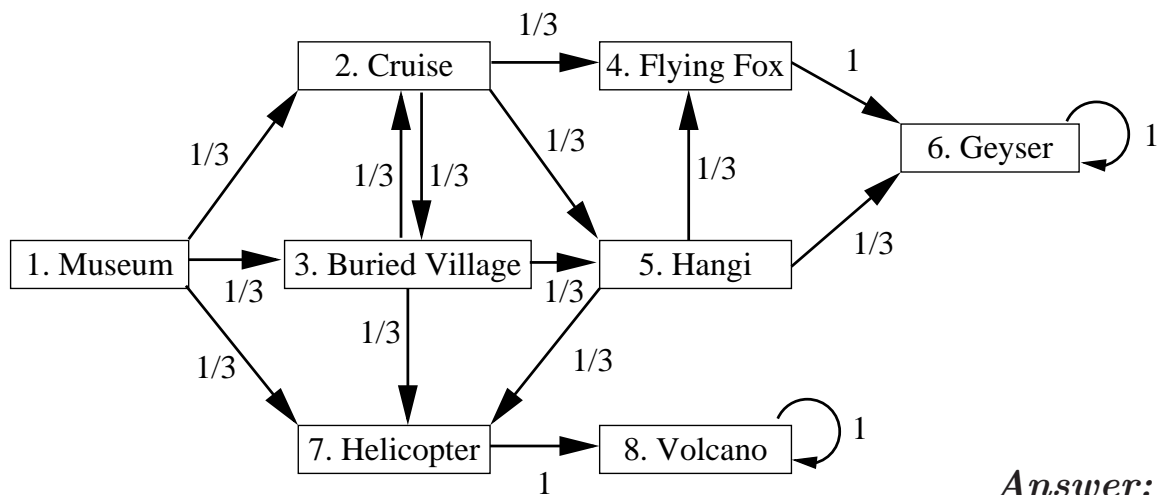
Answer: 0.997.



Markov's Marvellous Mystery Tours

Mr Markov's Marvellous Mystery Tours promises an All-Stochastic Tourist Experience for the town of Rotorua. Mr Markov has eight tourist attractions, to which he will take his clients completely at random with the probabilities shown below. He promises at least three exciting attractions per tour, ending at either the Lady Knox Geyser or the Tarawera Volcano. (Unfortunately he makes no mention of how the hapless tourist might get home from these places.)

What is the expected number of activities for a tour starting from the museum?



Answer: 4.2.

Structure of the course

- **Probability.** Probability and random variables, with special focus on conditional probability. Finding hitting probabilities for stochastic processes.
- **Expectation.** Expectation and variance. Introduction to conditional expectation, and its application in finding expected reaching times in stochastic processes.
- **Generating functions.** Introduction to probability generating functions, and their applications to stochastic processes, especially the Random Walk.
- **Branching process.** This process is a simple model for reproduction. Examples are the pyramid selling scheme and the spread of SARS above.

- **Markov chains.** Almost all the examples we look at throughout the course can be formulated as Markov chains. By developing a single unifying theory, we can easily tackle complex problems with many states and transitions like Markov’s Marvellous Mystery Tours above.

The rest of this chapter covers:

- quick revision of sample spaces and random variables;
- formal definition of stochastic processes.

1.1 Revision: Sample spaces and random variables

Definition: A **random experiment** is a physical situation whose outcome cannot be predicted until it is observed.

Definition: A **sample space**, Ω , is a set of possible outcomes of a random experiment.

Example:

Random experiment: Toss a coin once.

Sample space: $\Omega = \{\text{head}, \text{tail}\}$

Definition: A **random variable**, X , is defined as a function from the sample space to the real numbers: $X : \Omega \rightarrow \mathbb{R}$.

That is, *a random variable assigns a real number to every possible outcome of a random experiment.*

Example:

Random experiment: Toss a coin once.

Sample space: $\Omega = \{\text{head}, \text{tail}\}$.

An example of a random variable: $X : \Omega \rightarrow \mathbb{R}$ maps “head” $\rightarrow 1$, “tail” $\rightarrow 0$.

Essential point: A random variable is a way of producing random real numbers.

1.2 Stochastic Processes

Definition: A stochastic process is a family of random variables, $\{X(t) : t \in T\}$, where t usually denotes time. That is, at every time t in the set T , a random number $X(t)$ is observed.

Definition: $\{X(t) : t \in T\}$ is a discrete-time process if the set T is finite or countable.

In practice, this generally means $T = \{0, 1, 2, 3, \dots\}$

Thus a discrete-time process is $\{X(0), X(1), X(2), X(3), \dots\}$: a random number associated with every time $0, 1, 2, 3, \dots$

Definition: $\{X(t) : t \in T\}$ is a continuous-time process if T is not finite or countable.

In practice, this generally means $T = [0, \infty)$, or $T = [0, K]$ for some K .

Thus a continuous-time process $\{X(t) : t \in T\}$ has a random number $X(t)$ associated with every instant in time.

(Note that $X(t)$ need not *change* at every instant in time, but it is *allowed* to change at any time; i.e. not just at $t = 0, 1, 2, \dots$, like a discrete-time process.)

Definition: The state space, S , is *the set of real values that $X(t)$ can take*.

Every $X(t)$ takes a value in \mathbb{R} , but S will often be a smaller set: $S \subseteq \mathbb{R}$. For example, if $X(t)$ is the outcome of a coin tossed at time t , then the state space is $S = \{0, 1\}$.

Definition: The state space S is discrete if it is *finite or countable*. Otherwise it is *continuous*.

The state space S is the set of states that the stochastic process can be in.

For Reference: Discrete Random Variables

1. Binomial distribution

Notation: $X \sim \text{Binomial}(n, p)$.

Description: number of successes in n independent trials, each with probability p of success.

Probability function:

$$f_X(x) = \mathbb{P}(X = x) = \binom{n}{x} p^x (1 - p)^{n-x} \quad \text{for } x = 0, 1, \dots, n.$$

Mean: $\mathbb{E}(X) = np$.

Variance: $\text{Var}(X) = np(1 - p) = npq$, where $q = 1 - p$.

Sum: If $X \sim \text{Binomial}(n, p)$, $Y \sim \text{Binomial}(m, p)$, and X and Y are independent, then

$$X + Y \sim \text{Bin}(n + m, p).$$

2. Poisson distribution

Notation: $X \sim \text{Poisson}(\lambda)$.

Description: arises out of the Poisson process as the number of events in a fixed time or space, when events occur at a constant average rate. Also used in many other situations.

Probability function: $f_X(x) = \mathbb{P}(X = x) = \frac{\lambda^x}{x!} e^{-\lambda} \quad \text{for } x = 0, 1, 2, \dots$

Mean: $\mathbb{E}(X) = \lambda$.

Variance: $\text{Var}(X) = \lambda$.

Sum: If $X \sim \text{Poisson}(\lambda)$, $Y \sim \text{Poisson}(\mu)$, and X and Y are independent, then

$$X + Y \sim \text{Poisson}(\lambda + \mu).$$

3. Geometric distribution

Notation: $X \sim \text{Geometric}(p)$.

Description: number of failures before the **first** success in a sequence of independent trials, each with $\mathbb{P}(\text{success}) = p$.

Probability function: $f_X(x) = \mathbb{P}(X = x) = (1 - p)^x p$ for $x = 0, 1, 2, \dots$

Mean: $\mathbb{E}(X) = \frac{1 - p}{p} = \frac{q}{p}$, where $q = 1 - p$.

Variance: $\text{Var}(X) = \frac{1 - p}{p^2} = \frac{q}{p^2}$, where $q = 1 - p$.

Sum: if X_1, \dots, X_k are **independent**, and each $X_i \sim \text{Geometric}(p)$, then

$$X_1 + \dots + X_k \sim \text{Negative Binomial}(k, p).$$

4. Negative Binomial distribution

Notation: $X \sim \text{NegBin}(k, p)$.

Description: number of failures before the **kth** success in a sequence of independent trials, each with $\mathbb{P}(\text{success}) = p$.

Probability function:

$$f_X(x) = \mathbb{P}(X = x) = \binom{k + x - 1}{x} p^k (1 - p)^x \quad \text{for } x = 0, 1, 2, \dots$$

Mean: $\mathbb{E}(X) = \frac{k(1 - p)}{p} = \frac{kq}{p}$, where $q = 1 - p$.

Variance: $\text{Var}(X) = \frac{k(1 - p)}{p^2} = \frac{kq}{p^2}$, where $q = 1 - p$.

Sum: If $X \sim \text{NegBin}(k, p)$, $Y \sim \text{NegBin}(m, p)$, and X and Y are **independent**, then

$$X + Y \sim \text{NegBin}(k + m, p).$$

5. Hypergeometric distribution

Notation: $X \sim \text{Hypergeometric}(N, M, n)$.

Description: Sampling without replacement from a finite population. Given N objects, of which M are ‘special’. Draw n objects without replacement. X is the number of the n objects that are ‘special’.

Probability function:

$$f_X(x) = \mathbb{P}(X = x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}} \quad \text{for } \begin{cases} x = \max(0, n + M - N) \\ \text{to } x = \min(n, M). \end{cases}$$

Mean: $\mathbb{E}(X) = np$, where $p = \frac{M}{N}$.

Variance: $\text{Var}(X) = np(1-p) \left(\frac{N-n}{N-1} \right)$, where $p = \frac{M}{N}$.

6. Multinomial distribution

Notation: $\mathbf{X} = (X_1, \dots, X_k) \sim \text{Multinomial}(n; p_1, p_2, \dots, p_k)$.

Description: there are n independent trials, each with k possible outcomes. Let $p_i = \mathbb{P}(\text{outcome } i)$ for $i = 1, \dots, k$. Then $\mathbf{X} = (X_1, \dots, X_k)$, where X_i is the number of trials with outcome i , for $i = 1, \dots, k$.

Probability function:

$$f_{\mathbf{X}}(\mathbf{x}) = \mathbb{P}(X_1 = x_1, \dots, X_k = x_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

for $x_i \in \{0, \dots, n\} \forall_i$ with $\sum_{i=1}^k x_i = n$, and where $p_i \geq 0 \forall_i$, $\sum_{i=1}^k p_i = 1$.

Marginal distributions: $X_i \sim \text{Binomial}(n, p_i)$ for $i = 1, \dots, k$.

Mean: $\mathbb{E}(X_i) = np_i$ for $i = 1, \dots, k$.

Variance: $\text{Var}(X_i) = np_i(1 - p_i)$, for $i = 1, \dots, k$.

Covariance: $\text{cov}(X_i, X_j) = -np_i p_j$, for all $i \neq j$.

Continuous Random Variables

1. Uniform distribution

Notation: $X \sim \text{Uniform}(a, b)$.

Probability density function (pdf): $f_X(x) = \frac{1}{b-a}$ for $a < x < b$.

Cumulative distribution function:

$$F_X(x) = \mathbb{P}(X \leq x) = \frac{x-a}{b-a} \quad \text{for } a < x < b.$$

$$F_X(x) = 0 \text{ for } x \leq a, \text{ and } F_X(x) = 1 \text{ for } x \geq b.$$

Mean: $\mathbb{E}(X) = \frac{a+b}{2}$.

Variance: $\text{Var}(X) = \frac{(b-a)^2}{12}$.

2. Exponential distribution

Notation: $X \sim \text{Exponential}(\lambda)$.

Probability density function (pdf): $f_X(x) = \lambda e^{-\lambda x}$ for $0 < x < \infty$.

Cumulative distribution function:

$$F_X(x) = \mathbb{P}(X \leq x) = 1 - e^{-\lambda x} \quad \text{for } 0 < x < \infty.$$

$$F_X(x) = 0 \text{ for } x \leq 0.$$

Mean: $\mathbb{E}(X) = \frac{1}{\lambda}$.

Variance: $\text{Var}(X) = \frac{1}{\lambda^2}$.

Sum: if X_1, \dots, X_k are independent, and each $X_i \sim \text{Exponential}(\lambda)$, then

$$X_1 + \dots + X_k \sim \text{Gamma}(k, \lambda).$$

3. Gamma distribution

Notation: $X \sim \text{Gamma}(k, \lambda)$.

Probability density function (pdf):

$$f_X(x) = \frac{\lambda^k}{\Gamma(k)} x^{k-1} e^{-\lambda x} \quad \text{for } 0 < x < \infty,$$

where $\Gamma(k) = \int_0^\infty y^{k-1} e^{-y} dy$ (the Gamma function).

Cumulative distribution function: no closed form.

Mean: $\mathbb{E}(X) = \frac{k}{\lambda}$.

Variance: $\text{Var}(X) = \frac{k}{\lambda^2}$.

Sum: if X_1, \dots, X_n are independent, and $X_i \sim \text{Gamma}(k_i, \lambda)$, then

$$X_1 + \dots + X_n \sim \text{Gamma}(k_1 + \dots + k_n, \lambda).$$

4. Normal distribution

Notation: $X \sim \text{Normal}(\mu, \sigma^2)$.

Probability density function (pdf):

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\{-(x-\mu)^2/2\sigma^2\}} \quad \text{for } -\infty < x < \infty.$$

Cumulative distribution function: no closed form.

Mean: $\mathbb{E}(X) = \mu$.

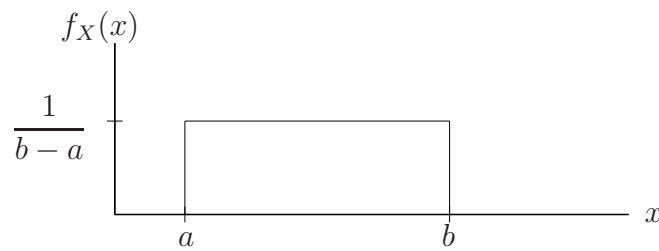
Variance: $\text{Var}(X) = \sigma^2$.

Sum: if X_1, \dots, X_n are independent, and $X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$, then

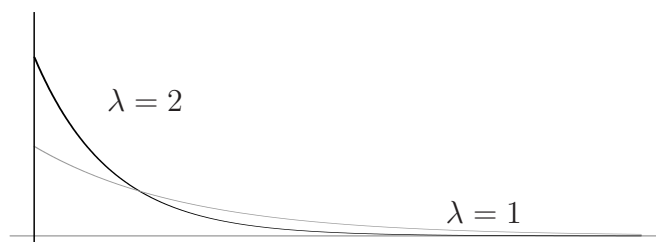
$$X_1 + \dots + X_n \sim \text{Normal}(\mu_1 + \dots + \mu_n, \sigma_1^2 + \dots + \sigma_n^2).$$

Probability Density Functions

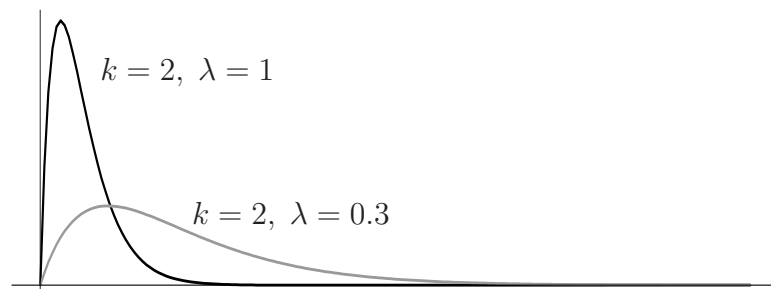
Uniform(a, b)



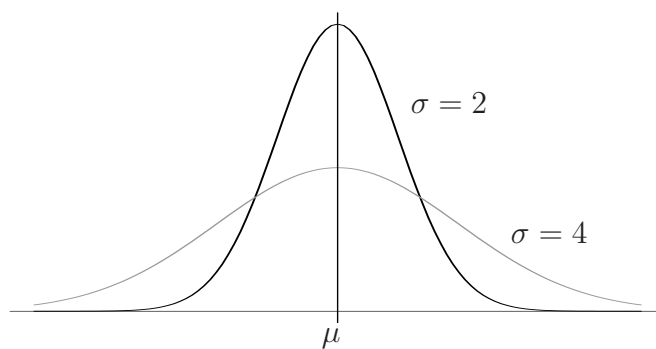
Exponential(λ)



Gamma(k, λ)



Normal(μ, σ^2)



Chapter 2: Probability

The aim of this chapter is to revise the basic rules of probability. By the end of this chapter, you should be comfortable with:

- conditional probability, and what you can and can't do with conditional expressions;
 - the Partition Theorem and Bayes' Theorem;
 - First-Step Analysis for finding the probability that a process reaches some state, by conditioning on the outcome of the first step;
 - calculating probabilities for continuous and discrete random variables.
-

2.1 Sample spaces and events

Definition: A sample space, Ω , is a *set of possible outcomes of a random experiment*.

Definition: An event, A , is a *subset of the sample space*.

This means that event A is simply *a collection of outcomes*.

Example:

Random experiment: Pick a person in this class at random.

Sample space: $\Omega = \{\text{all people in class}\}$

Event A : $A = \{\text{all males in class}\}$.

Definition: Event A occurs if *the outcome of the random experiment is a member of the set A* .

In the example above, event A occurs if *the person we pick is male*.

2.2 Probability Reference List

The following properties hold for all events A, B .

- $\mathbb{P}(\emptyset) = 0$.
- $0 \leq \mathbb{P}(A) \leq 1$.
- **Complement:** $\mathbb{P}(\overline{A}) = 1 - \mathbb{P}(A)$.
- **Probability of a union:** $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$.

For three events A, B, C :

$$\mathbb{P}(A \cup B \cup C) = \mathbb{P}(A) + \mathbb{P}(B) + \mathbb{P}(C) - \mathbb{P}(A \cap B) - \mathbb{P}(A \cap C) - \mathbb{P}(B \cap C) + \mathbb{P}(A \cap B \cap C).$$

If A and B are **mutually exclusive**, then $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$.

- **Conditional probability:** $\mathbb{P}(A | B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$.
- **Multiplication rule:** $\mathbb{P}(A \cap B) = \mathbb{P}(A | B)\mathbb{P}(B) = \mathbb{P}(B | A)\mathbb{P}(A)$.
- **The Partition Theorem:** if B_1, B_2, \dots, B_m form a partition of Ω , then

$$\mathbb{P}(A) = \sum_{i=1}^m \mathbb{P}(A \cap B_i) = \sum_{i=1}^m \mathbb{P}(A | B_i)\mathbb{P}(B_i) \quad \text{for any event } A.$$

As a special case, B and \overline{B} partition Ω , so:

$$\begin{aligned} \mathbb{P}(A) &= \mathbb{P}(A \cap B) + \mathbb{P}(A \cap \overline{B}) \\ &= \mathbb{P}(A | B)\mathbb{P}(B) + \mathbb{P}(A | \overline{B})\mathbb{P}(\overline{B}) \quad \text{for any } A, B. \end{aligned}$$

- **Bayes' Theorem:** $\mathbb{P}(B | A) = \frac{\mathbb{P}(A | B)\mathbb{P}(B)}{\mathbb{P}(A)}$.

More generally, if B_1, B_2, \dots, B_m form a partition of Ω , then

$$\mathbb{P}(B_j | A) = \frac{\mathbb{P}(A | B_j)\mathbb{P}(B_j)}{\sum_{i=1}^m \mathbb{P}(A | B_i)\mathbb{P}(B_i)} \quad \text{for any } j.$$

- **Chains of events:** for any events A_1, A_2, \dots, A_n ,

$$\mathbb{P}(A_1 \cap A_2 \cap \dots \cap A_n) = \mathbb{P}(A_1)\mathbb{P}(A_2 | A_1)\mathbb{P}(A_3 | A_2 \cap A_1) \dots \mathbb{P}(A_n | A_{n-1} \cap \dots \cap A_1).$$



2.3 Conditional Probability

Suppose we are working with sample space $\Omega = \{\text{people in class}\}$. I want to find the proportion of people in the class who ski. What do I do?

Count up the number of people in the class who ski, and divide by the total number of people in the class.

$$\mathbb{P}(\text{person skis}) = \frac{\text{number of skiers in class}}{\text{total number of people in class}}.$$

Now suppose I want to find the proportion of *females* in the class who ski. What do I do?

Count up the number of females in the class who ski, and divide by the total number of females in the class.

$$\mathbb{P}(\text{female skis}) = \frac{\text{number of female skiers in class}}{\text{total number of females in class}}.$$

By changing from asking about everyone to asking about females only, we have:

- *restricted attention to the set of females only,*

or: reduced the sample space from the set of everyone to the set of females,

or: conditioned on the event $\{\text{females}\}$.

We could write the above as:

$$\mathbb{P}(\text{skis} \mid \text{female}) = \frac{\text{number of female skiers in class}}{\text{total number of females in class}}.$$

Conditioning is like changing the sample space: we are now working in a new sample space of females in class.

In the above example, we could replace ‘skiing’ with *any* attribute B . We have:

$$\mathbb{P}(\text{skis}) = \frac{\# \text{ skiers in class}}{\# \text{ class}}; \quad \mathbb{P}(\text{skis} \mid \text{female}) = \frac{\# \text{ female skiers in class}}{\# \text{ females in class}};$$

so:

$$\mathbb{P}(B) = \frac{\# B\text{'s in class}}{\text{total } \# \text{ people in class}},$$

and:

$$\begin{aligned} \mathbb{P}(B \mid \text{female}) &= \frac{\# \text{ female } B\text{'s in class}}{\text{total } \# \text{ females in class}} \\ &= \frac{\# \text{ in class who are } B \text{ and female}}{\# \text{ in class who are female}}. \end{aligned}$$

Likewise, we could replace ‘female’ with any attribute A :

$$\mathbb{P}(B \mid A) = \frac{\text{number in class who are } B \text{ and } A}{\text{number in class who are } A}.$$

This is how we get the definition of conditional probability:

$$\mathbb{P}(B \mid A) = \frac{\mathbb{P}(B \text{ and } A)}{\mathbb{P}(A)} = \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)}.$$

By conditioning on event A , we have *changed the sample space to the set of A ’s only*.

Definition: Let A and B be events on the same sample space: so $A \subseteq \Omega$ and $B \subseteq \Omega$.
The conditional probability of event B , given event A , is

$$\mathbb{P}(B \mid A) = \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)}.$$

Multiplication Rule: (Immediate from above). For any events A and B ,

$$\mathbb{P}(A \cap B) = \mathbb{P}(A | B)\mathbb{P}(B) = \mathbb{P}(B | A)\mathbb{P}(A) = \mathbb{P}(B \cap A).$$

Conditioning as ‘changing the sample space’

The idea that “*conditioning*” = “*changing the sample space*” can be very helpful in understanding how to manipulate conditional probabilities.

Any ‘unconditional’ probability can be written as a conditional probability:

$$\mathbb{P}(B) = \mathbb{P}(B | \Omega).$$

Writing $\mathbb{P}(B) = \mathbb{P}(B | \Omega)$ just means that we are looking for the probability of event B , out of all possible outcomes in the set Ω .

In fact, the symbol \mathbb{P} *belongs* to the set Ω : it has *no meaning without* Ω . To remind ourselves of this, we can write

$$\mathbb{P} = \mathbb{P}_\Omega.$$

Then $\mathbb{P}(B) = \mathbb{P}(B | \Omega) = \mathbb{P}_\Omega(B)$.

Similarly, $\mathbb{P}(B | A)$ means that we are looking for the probability of event B , out of all possible outcomes in the set A .

So A is just another sample space. Thus *we can manipulate conditional probabilities $\mathbb{P}(\cdot | A)$ just like any other probabilities, as long as we always stay inside the same sample space A .*

The trick: Because we can think of A as just another sample space, let’s write

$$\mathbb{P}(\cdot | A) = \mathbb{P}_A(\cdot)$$

***Note: NOT
standard notation!***

Then *we can use \mathbb{P}_A just like \mathbb{P} , as long as we remember to keep the A subscript on **EVERY** \mathbb{P} that we write.*

This helps us to make quite complex manipulations of conditional probabilities without thinking too hard or making mistakes. There is only one rule you need to learn to use this tool effectively:

$$\mathbb{P}_A(B | C) = \mathbb{P}(B | C \cap A) \text{ for any } A, B, C.$$

(Proof: Exercise).

The rules:

$$\begin{aligned} \mathbb{P}(\cdot | A) &= \mathbb{P}_A(\cdot) \\ \mathbb{P}_A(B | C) &= \mathbb{P}(B | C \cap A) \text{ for any } A, B, C. \end{aligned}$$

Examples:

1. Probability of a union. In general,

$$\mathbb{P}(B \cup C) = \mathbb{P}(B) + \mathbb{P}(C) - \mathbb{P}(B \cap C).$$

So, $\mathbb{P}_A(B \cup C) = \mathbb{P}_A(B) + \mathbb{P}_A(C) - \mathbb{P}_A(B \cap C).$

Thus, $\mathbb{P}(B \cup C | A) = \mathbb{P}(B | A) + \mathbb{P}(C | A) - \mathbb{P}(B \cap C | A).$

2. Which of the following is equal to $\mathbb{P}(B \cap C | A)$?

(a) $\mathbb{P}(B | C \cap A).$ (c) $\mathbb{P}(B | C \cap A)\mathbb{P}(C | A).$

(b) $\frac{\mathbb{P}(B | C)}{\mathbb{P}(A)}.$ (d) $\mathbb{P}(B | C)\mathbb{P}(C | A).$

Solution:

$$\begin{aligned} \mathbb{P}(B \cap C | A) &= \mathbb{P}_A(B \cap C) \\ &= \mathbb{P}_A(B | C) \mathbb{P}_A(C) \\ &= \mathbb{P}(B | C \cap A) \mathbb{P}(C | A). \end{aligned}$$

Thus the correct answer is (c).

3. Which of the following is true?

(a) $\mathbb{P}(\overline{B} | A) = 1 - \mathbb{P}(B | A)$.

(b) $\mathbb{P}(\overline{B} | A) = \mathbb{P}(B) - \mathbb{P}(B | A)$.

Solution:

$$\mathbb{P}(\overline{B} | A) = \mathbb{P}_A(\overline{B}) = 1 - \mathbb{P}_A(B) = 1 - \mathbb{P}(B | A).$$

Thus the correct answer is (a).

4. Which of the following is true?

(a) $\mathbb{P}(\overline{B} \cap A) = \mathbb{P}(A) - \mathbb{P}(B \cap A)$.

(b) $\mathbb{P}(\overline{B} \cap A) = \mathbb{P}(B) - \mathbb{P}(B \cap A)$.

Solution:

$$\begin{aligned} \mathbb{P}(\overline{B} \cap A) &= \mathbb{P}(\overline{B} | A) \mathbb{P}(A) = \mathbb{P}_A(\overline{B}) \mathbb{P}(A) \\ &= (1 - \mathbb{P}_A(B)) \mathbb{P}(A) \\ &= \mathbb{P}(A) - \mathbb{P}(B | A) \mathbb{P}(A) \\ &= \mathbb{P}(A) - \mathbb{P}(B \cap A). \end{aligned}$$

Thus the correct answer is (a).

5. True or false: $\mathbb{P}(B | A) = 1 - \mathbb{P}(B | \overline{A})$?

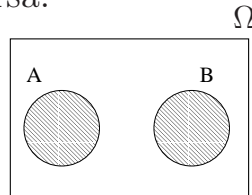
Answer: False. $\mathbb{P}(B | A) = \mathbb{P}_A(B)$. Once we have \mathbb{P}_A , we are stuck with it! There is no easy way of converting from \mathbb{P}_A to $\mathbb{P}_{\overline{A}}$: or anything else. Probabilities in one sample space (\mathbb{P}_A) cannot tell us anything about probabilities in a different sample space ($\mathbb{P}_{\overline{A}}$).

Exercise: if we wish to express $\mathbb{P}(B | A)$ in terms of only B and \overline{A} , show that
$$\mathbb{P}(B | A) = \frac{\mathbb{P}(B) - \mathbb{P}(B | \overline{A}) \mathbb{P}(\overline{A})}{1 - \mathbb{P}(\overline{A})}.$$
 Note that this does not simplify nicely!

2.4 The Partition Theorem (Law of Total Probability)

Definition: Events A and B are mutually exclusive, or disjoint, if $A \cap B = \emptyset$.

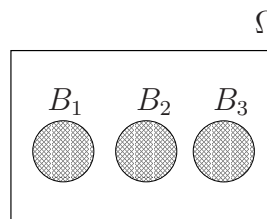
This means events A and B cannot happen together. If A happens, it excludes B from happening, and vice-versa.



If A and B are mutually exclusive, $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$.

For all other A and B , $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$.

Definition: Any number of events B_1, B_2, \dots, B_k are mutually exclusive if every pair of the events is mutually exclusive: ie. $B_i \cap B_j = \emptyset$ for all i, j with $i \neq j$.



Definition: A partition of Ω is a *collection of mutually exclusive events whose union is Ω* .

That is, sets B_1, B_2, \dots, B_k form a partition of Ω if

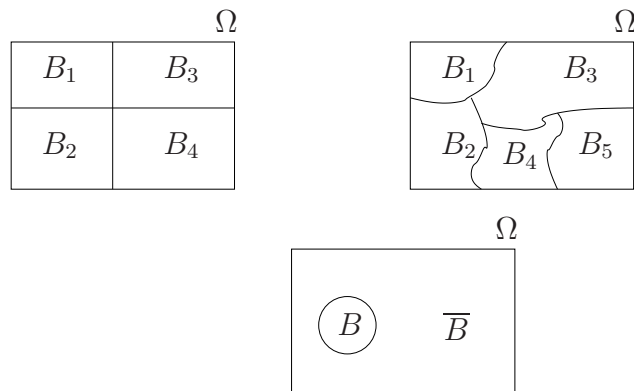
$$B_i \cap B_j = \emptyset \text{ for all } i, j \text{ with } i \neq j,$$

and

$$\bigcup_{i=1}^k B_i = B_1 \cup B_2 \cup \dots \cup B_k = \Omega.$$

B_1, \dots, B_k form a partition of Ω if they *have no overlap*
and collectively cover all possible outcomes.

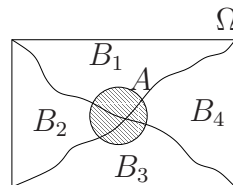
Examples:



Partitioning an event A

Any set A can be partitioned: it doesn't have to be Ω .

In particular, if B_1, \dots, B_k form a partition of Ω , then $(A \cap B_1), \dots, (A \cap B_k)$ form a partition of A .



Theorem 2.4: The Partition Theorem (Law of Total Probability)

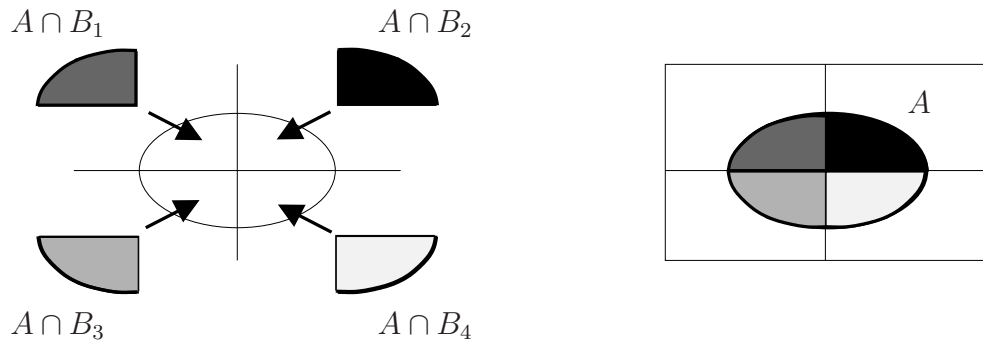
Let B_1, \dots, B_m form a partition of Ω . Then for any event A ,

$$\mathbb{P}(A) = \sum_{i=1}^m \mathbb{P}(A \cap B_i) = \sum_{i=1}^m \mathbb{P}(A | B_i) \mathbb{P}(B_i)$$

Both formulations of the Partition Theorem are very widely used, but especially the conditional formulation $\sum_{i=1}^m \mathbb{P}(A | B_i) \mathbb{P}(B_i)$.

Intuition behind the Partition Theorem:

The Partition Theorem is easy to understand because it simply states that “the whole is the sum of its parts.”



$$\mathbb{P}(A) = \mathbb{P}(A \cap B_1) + \mathbb{P}(A \cap B_2) + \mathbb{P}(A \cap B_3) + \mathbb{P}(A \cap B_4).$$

2.5 Bayes' Theorem: inverting conditional probabilities

Bayes' Theorem allows us to “invert” a conditional statement, ie. *to express* $\mathbb{P}(B | A)$ *in terms of* $\mathbb{P}(A | B)$.

Theorem 2.5: Bayes' Theorem

For any events A and B:

$$\mathbb{P}(B | A) = \frac{\mathbb{P}(A | B)\mathbb{P}(B)}{\mathbb{P}(A)}.$$

Proof:

$$\mathbb{P}(B \cap A) = \mathbb{P}(A \cap B)$$

$$\mathbb{P}(B | A)\mathbb{P}(A) = \mathbb{P}(A | B)\mathbb{P}(B) \quad (\text{multiplication rule})$$

$$\therefore \mathbb{P}(B | A) = \frac{\mathbb{P}(A | B)\mathbb{P}(B)}{\mathbb{P}(A)}. \quad \square$$

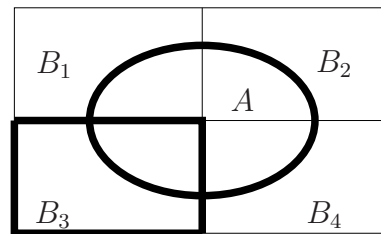
Extension of Bayes' Theorem

Suppose that B_1, B_2, \dots, B_m form a partition of Ω . By the Partition Theorem,

$$\mathbb{P}(A) = \sum_{i=1}^m \mathbb{P}(A | B_i) \mathbb{P}(B_i).$$

Thus, for *any single partition member* B_j , put $B = B_j$ in Bayes' Theorem to obtain:

$$\mathbb{P}(B_j | A) = \frac{\mathbb{P}(A | B_j) \mathbb{P}(B_j)}{\mathbb{P}(A)} = \frac{\mathbb{P}(A | B_j) \mathbb{P}(B_j)}{\sum_{i=1}^m \mathbb{P}(A | B_i) \mathbb{P}(B_i)}.$$



Special case: $m = 2$

Given any event B , the events B and \overline{B} form a partition of Ω . Thus:

$$\mathbb{P}(B | A) = \frac{\mathbb{P}(A | B) \mathbb{P}(B)}{\mathbb{P}(A | B) \mathbb{P}(B) + \mathbb{P}(A | \overline{B}) \mathbb{P}(\overline{B})}.$$

Example: In screening for a certain disease, the probability that a healthy person wrongly gets a positive result is 0.05. The probability that a diseased person wrongly gets a negative result is 0.002. The overall rate of the disease in the population being screened is 1%. If my test gives a positive result, what is the probability I actually have the disease?

1. Define events:

$$D = \{\text{have disease}\} \quad \overline{D} = \{\text{do not have the disease}\}$$

$$P = \{\text{positive test}\} \quad N = \overline{P} = \{\text{negative test}\}$$

2. Information given:

$$\text{False positive rate is } 0.05 \Rightarrow \mathbb{P}(P | \overline{D}) = 0.05$$

$$\text{False negative rate is } 0.002 \Rightarrow \mathbb{P}(N | D) = 0.002$$

$$\text{Disease rate is } 1\% \Rightarrow \mathbb{P}(D) = 0.01.$$

3. Looking for $\mathbb{P}(D | P)$:

$$\text{We have} \quad \mathbb{P}(D | P) = \frac{\mathbb{P}(P | D)\mathbb{P}(D)}{\mathbb{P}(P)}.$$

$$\begin{aligned} \text{Now} \quad \mathbb{P}(P | D) &= 1 - \mathbb{P}(\overline{P} | D) \\ &= 1 - \mathbb{P}(N | D) \\ &= 1 - 0.002 \\ &= 0.998. \end{aligned}$$

$$\begin{aligned} \text{Also} \quad \mathbb{P}(P) &= \mathbb{P}(P | D)\mathbb{P}(D) + \mathbb{P}(P | \overline{D})\mathbb{P}(\overline{D}) \\ &= 0.998 \times 0.01 + 0.05 \times (1 - 0.01) \\ &= 0.05948. \end{aligned}$$

Thus

$$\mathbb{P}(D | P) = \frac{0.998 \times 0.01}{0.05948} = 0.168.$$

Given a positive test, my chance of having the disease is only 16.8%.

2.6 First-Step Analysis for calculating probabilities in a process

In a stochastic process, what happens at the next step depends upon the current state of the process. We often wish to know the probability of *eventually* reaching some particular state, given our current position.

Throughout this course, we will tackle this sort of problem using a technique called **First-Step Analysis**.

The idea is to consider all possible first steps away from the current state. We derive a system of equations that specify the probability of the eventual outcome given each of the possible first steps. We then try to solve these equations for the probability of interest.

First-Step Analysis depends upon **conditional probability** and the **Partition Theorem**. Let S_1, \dots, S_k be the k possible first steps we can take away from our current state. We wish to find the probability that event E happens eventually. First-Step Analysis calculates $\mathbb{P}(E)$ as follows:

$$\mathbb{P}(E) = \mathbb{P}(E|S_1)\mathbb{P}(S_1) + \dots + \mathbb{P}(E|S_k)\mathbb{P}(S_k).$$

Here, $\mathbb{P}(S_1), \dots, \mathbb{P}(S_k)$ give the probabilities of taking the different first steps $1, 2, \dots, k$.

Example: Tennis game at Deuce.

Venus and Serena are playing tennis, and have reached the score Deuce (40-40). (*Deuce* comes from the French word *Deux* for ‘two’, meaning that each player needs to win two consecutive points to win the game.)

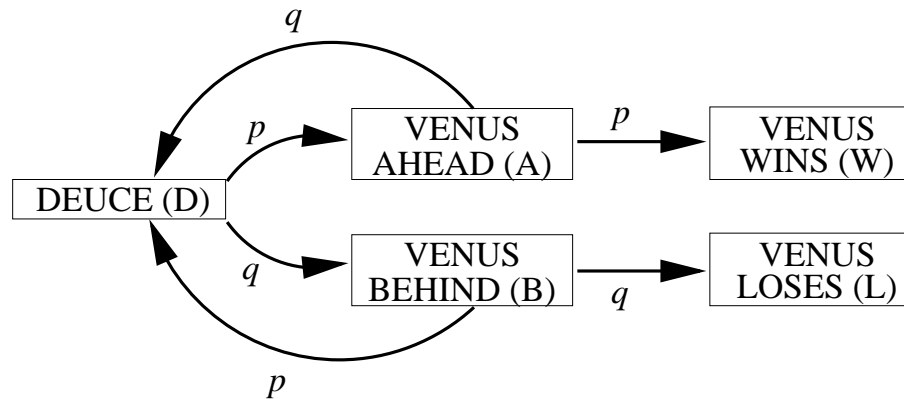


For each point, let:

$$p = \mathbb{P}(\text{Venus wins point}), \quad q = 1 - p = \mathbb{P}(\text{Serena wins point}).$$

Assume that all points are independent.

Let v be the probability that Venus wins the game eventually, starting from Deuce. Find v .



Use First-step analysis. The possible steps starting from Deuce are:

1. Venus wins the next point (probability p): move to state A;
2. Venus loses the next point (probability q): move to state B.

Let V be the event that Venus wins *EVENTUALLY* starting from Deuce, so $v = \mathbb{P}(V \mid D)$. Starting from Deuce (D), the possible steps are to states A and B. So:

$$\begin{aligned}
 v &= \mathbb{P}(\text{Venus wins} \mid D) = \mathbb{P}(V \mid D) \\
 &= \mathbb{P}_D(V) \\
 &= \mathbb{P}_D(V \mid A)\mathbb{P}_D(A) + \mathbb{P}_D(V \mid B)\mathbb{P}_D(B) \\
 &= \mathbb{P}(V \mid A)p + \mathbb{P}(V \mid B)q. \quad (\star)
 \end{aligned}$$

Now we need to find $\mathbb{P}(V \mid A)$, and $\mathbb{P}(V \mid B)$, again using First-step analysis:

$$\begin{aligned}
 \mathbb{P}(V \mid A) &= \mathbb{P}(V \mid W)p + \mathbb{P}(V \mid D)q \\
 &= 1 \times p + v \times q \\
 &= p + qv. \quad (a)
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \mathbb{P}(V \mid B) &= \mathbb{P}(V \mid L)q + \mathbb{P}(V \mid D)p \\
 &= 0 \times q + v \times p \\
 &= pv. \quad (b)
 \end{aligned}$$

Substituting (a) and (b) into (★),

$$\begin{aligned} v &= (p + qv)p + (pv)q \\ v &= p^2 + 2pqv \\ v(1 - 2pq) &= p^2 \\ v &= \frac{p^2}{1 - 2pq}. \end{aligned}$$

Note: Because $p + q = 1$, we have:

$$1 = (p + q)^2 = p^2 + q^2 + 2pq.$$

So the final probability that Venus wins the game is:

$$v = \frac{p^2}{1 - 2pq} = \frac{p^2}{p^2 + q^2}.$$

Note how this result makes intuitive sense. For the game to finish from Deuce, either Venus has to win two points in a row (probability p^2), or Serena does (probability q^2). The ratio $p^2/(p^2 + q^2)$ describes Venus's 'share' of the winning probability.

First-step analysis as the Partition Theorem:

Our approach to finding $v = \mathbb{P}(\text{Venus wins})$ can be summarized as:

$$\mathbb{P}(\text{Venus wins}) = v = \sum_{\text{first steps}} \mathbb{P}(V \mid \text{first step})\mathbb{P}(\text{first step}).$$

First-step analysis is just the **Partition Theorem**:

The sample space is $\Omega = \{ \text{all possible routes from Deuce to the end} \}$.

An example of a sample point is: $D \rightarrow A \rightarrow D \rightarrow B \rightarrow D \rightarrow B \rightarrow L$.

Another example is: $D \rightarrow B \rightarrow D \rightarrow A \rightarrow W$.

The **partition** of the sample space that we use in first-step analysis is:

$$R_1 = \{ \text{all possible routes from Deuce to the end that start with } D \rightarrow A \}$$

$$R_2 = \{ \text{all possible routes from Deuce to the end that start with } D \rightarrow B \}$$

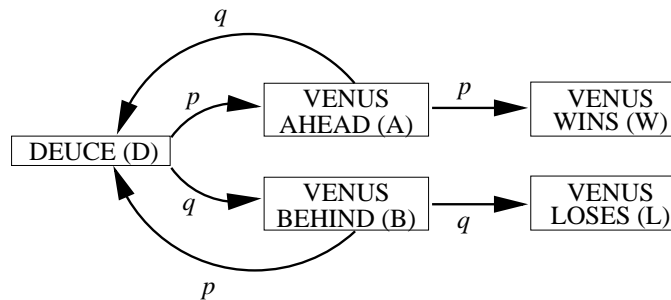
Then first-step analysis simply states:

$$\begin{aligned}\mathbb{P}(V) &= \mathbb{P}(V | R_1)\mathbb{P}(R_1) + \mathbb{P}(V | R_2)\mathbb{P}(R_2) \\ &= \mathbb{P}_D(V | A)\mathbb{P}_D(A) + \mathbb{P}_D(V | B)\mathbb{P}_D(B).\end{aligned}$$

Notation for quick solutions of first-step analysis problems

Defining a helpful notation is central to modelling with stochastic processes. Setting up well-defined notation helps you to solve problems quickly and easily. Defining your notation is one of the most important steps in modelling, because it provides the conversion from words (which is how your problem starts) to mathematics (which is how your problem is solved).

Several marks are allotted on first-step analysis questions for setting up a well-defined and helpful notation.



Here is the correct way to formulate and solve this first-step analysis problem.

Need the probability that Venus wins eventually, starting from Deuce.

1. Define notation: let

$$v_D = \mathbb{P}(\text{Venus wins eventually} \mid \text{start at state } D)$$

$$v_A = \mathbb{P}(\text{Venus wins eventually} \mid \text{start at state } A)$$

$$v_B = \mathbb{P}(\text{Venus wins eventually} \mid \text{start at state } B)$$

2. First-step analysis:

$$v_D = pv_A + qv_B \quad (a)$$

$$v_A = p \times 1 + qv_D \quad (b)$$

$$v_B = pv_D + q \times 0 \quad (c)$$

3. Substitute (b) and (c) in (a):

$$\begin{aligned}\Rightarrow v_D &= p(p + qv_D) + q(pv_D) \\ v_D(1 - pq - pq) &= p^2 \\ \therefore v_D &= \frac{p^2}{1 - 2pq}\end{aligned}$$

as before.

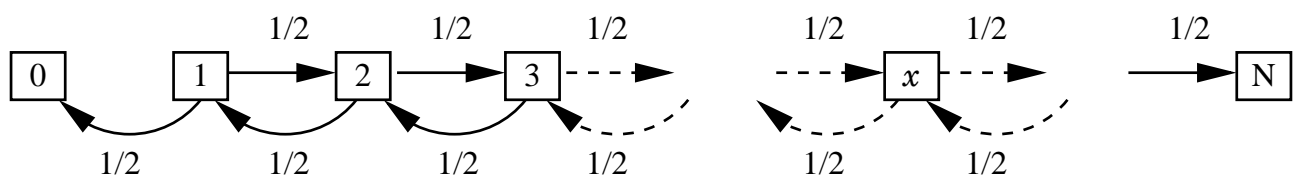
2.7 Special Process: the Gambler's Ruin

This is a famous problem in probability. A gambler starts with \$ x . She tosses a fair coin repeatedly.

If she gets a Head, she wins \$1. If she gets a Tail, she loses \$1.



The coin tossing is repeated until the gambler has either \$0 or \$ N , when she stops. What is the probability of the Gambler's Ruin, i.e. that the gambler ends up with \$0?



Wish to find

$$\mathbb{P}(\text{ends with } \$0 \mid \text{starts with } \$x).$$

Define event

$$R = \{\text{eventual Ruin}\} = \{\text{ends with } \$0\}.$$

We wish to find $\mathbb{P}(R \mid \text{starts with } \$x)$.

Define notation:

$$p_x = \mathbb{P}(R \mid \text{currently has } \$x) \quad \text{for } x = 0, 1, \dots, N.$$

Information given:

$$\begin{aligned} p_0 &= \mathbb{P}(R \mid \text{currently has \$0}) = 1, \\ p_N &= \mathbb{P}(R \mid \text{currently has \$N}) = 0. \end{aligned}$$

First-step analysis:

$$\begin{aligned} p_x &= \mathbb{P}(R \mid \text{has \$}x) \\ &= \frac{1}{2}\mathbb{P}(R \mid \text{has \$}(x+1)) + \frac{1}{2}\mathbb{P}(R \mid \text{has \$}(x-1)) \\ &= \frac{1}{2}p_{x+1} + \frac{1}{2}p_{x-1} \quad (\star) \end{aligned}$$

True for $x = 1, 2, \dots, N-1$, with boundary conditions $p_0 = 1, p_N = 0$.

Solution of difference equation (\star) :

$$\begin{aligned} p_x &= \frac{1}{2}p_{x+1} + \frac{1}{2}p_{x-1} \quad \text{for } x = 1, 2, \dots, N-1; \\ p_0 &= 1 \\ p_N &= 0. \end{aligned} \quad (\star)$$

We usually solve equations like this using the theory of 2nd-order difference equations. For this special case we will also verify the answer by two other methods.

1. Theory of linear 2nd order difference equations

Theory tells us that the general solution of (\star) is $p_x = A + Bx$ for some constants A, B and for $x = 0, 1, \dots, N$. Our job is to find A and B using the boundary conditions:

$$p_x = A + Bx \text{ for constants } A \text{ and } B \text{ and for } x = 0, 1, \dots, N.$$

So

$$\begin{aligned} p_0 &= A + B \times 0 = 1 \quad \Rightarrow \quad A = 1; \\ p_N &= A + B \times N = 1 + BN = 0 \quad \Rightarrow \quad B = -\frac{1}{N}. \end{aligned}$$

So our solution is:

$$p_x = A + Bx = 1 - \frac{x}{N} \text{ for } x = 0, 1, \dots, N.$$

For Stats 325, you will be told the general solution of the 2nd-order difference equation and expected to solve it using the boundary conditions.

For Stats 721, we will study the theory of 2nd-order difference equations. You will be able to derive the general solution for yourself before solving it.

Question: What is the probability that the gambler wins (ends with \$N), starting with \$x?

$$\mathbb{P}(\text{ends with } \$N) = 1 - \mathbb{P}(\text{ends with } \$0) = 1 - p_x = \frac{x}{N} \text{ for } x = 0, 1, \dots, N.$$

2. Solution by inspection

The problem shown in this section is the *symmetric* Gambler's Ruin, where the probability is $\frac{1}{2}$ of moving up or down on any step. For this special case, we can solve the difference equation by inspection.

We have:

$$\begin{aligned} p_x &= \frac{1}{2}p_{x+1} + \frac{1}{2}p_{x-1} \\ \frac{1}{2}p_x + \frac{1}{2}p_x &= \frac{1}{2}p_{x+1} + \frac{1}{2}p_{x-1} \end{aligned}$$

Rearranging: $p_{x-1} - p_x = p_x - p_{x+1}.$

Boundaries: $p_0 = 1, p_N = 0.$

There are N steps to go down from $p_0 = 1$ to $p_N = 0$.

Each step is the same size, because

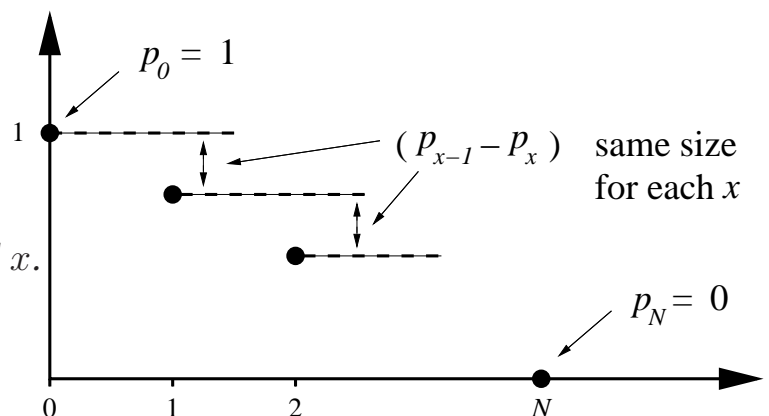
$$(p_{x-1} - p_x) = (p_x - p_{x+1}) \text{ for all } x.$$

So each step has size $1/N$,

$$\Rightarrow p_0 = 1, p_1 = 1 - 1/N, p_2 = 1 - 2/N, \text{ etc.}$$

So

$$p_x = 1 - \frac{x}{N} \text{ as before.}$$



3. Solution by repeated substitution.

In principle, all systems could be solved by this method, but it is usually too tedious to apply in practice.

Rearrange (★) to give:

$$\begin{aligned}
 p_{x+1} &= 2p_x - p_{x-1} \\
 \Rightarrow (x=1) \quad p_2 &= 2p_1 - 1 \quad (\text{recall } p_0 = 1) \\
 (x=2) \quad p_3 &= 2p_2 - p_1 = 2(2p_1 - 1) - p_1 = 3p_1 - 2 \\
 (x=3) \quad p_4 &= 2p_3 - p_2 = 2(3p_1 - 2) - (2p_1 - 1) = 4p_1 - 3 \quad \text{etc} \\
 &\vdots \\
 \text{giving} \quad p_x &= xp_1 - (x-1) \quad \text{in general,} \quad (**) \\
 \text{likewise} \quad p_N &= Np_1 - (N-1) \quad \text{at endpoint.}
 \end{aligned}$$

Boundary condition: $p_N = 0 \Rightarrow Np_1 - (N-1) = 0 \Rightarrow p_1 = 1 - 1/N.$

*Substitute in (**):*

$$\begin{aligned}
 p_x &= xp_1 - (x-1) \\
 &= x \left(1 - \frac{1}{N}\right) - (x-1) \\
 &= x - \frac{x}{N} - x + 1 \\
 p_x &= 1 - \frac{x}{N} \quad \text{as before.} \quad \square
 \end{aligned}$$

2.8 Independence

Definition: Events A and B are statistically independent if and only if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B).$$

This implies that A and B are statistically independent if and only if $\mathbb{P}(A | B) = \mathbb{P}(A).$

Note: If events are *physically* independent, they will also be statistically indept.

For interest: more than two events

Definition: For more than two events, A_1, A_2, \dots, A_n , we say that A_1, A_2, \dots, A_n are mutually independent if

$$\mathbb{P}\left(\bigcap_{i \in J} A_i\right) = \prod_{i \in J} \mathbb{P}(A_i) \quad \text{for ALL finite subsets } J \subseteq \{1, 2, \dots, n\}.$$

Example: events A_1, A_2, A_3, A_4 are mutually independent if

- i) $\mathbb{P}(A_i \cap A_j) = \mathbb{P}(A_i)\mathbb{P}(A_j)$ for all i, j with $i \neq j$; AND
- ii) $\mathbb{P}(A_i \cap A_j \cap A_k) = \mathbb{P}(A_i)\mathbb{P}(A_j)\mathbb{P}(A_k)$ for all i, j, k that are all different; AND
- iii) $\mathbb{P}(A_1 \cap A_2 \cap A_3 \cap A_4) = \mathbb{P}(A_1)\mathbb{P}(A_2)\mathbb{P}(A_3)\mathbb{P}(A_4)$.

Note: For mutual independence, it is **not** enough to check that $\mathbb{P}(A_i \cap A_j) = \mathbb{P}(A_i)\mathbb{P}(A_j)$ for all $i \neq j$. Pairwise independence does not imply mutual independence.

2.9 The Continuity Theorem

The Continuity Theorem states that probability is a *continuous set function*:

Theorem 2.9: The Continuity Theorem

- a) Let A_1, A_2, \dots be an *increasing sequence of events*: i.e.

$$A_1 \subseteq A_2 \subseteq \dots \subseteq A_n \subseteq A_{n+1} \subseteq \dots$$

Then

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} A_n\right) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n).$$

Note: because $A_1 \subseteq A_2 \subseteq \dots$, we have: $\lim_{n \rightarrow \infty} A_n = \bigcup_{n=1}^{\infty} A_n$.

b) Let B_1, B_2, \dots be a *decreasing sequence of events*: i.e.

$$B_1 \supseteq B_2 \supseteq \dots \supseteq B_n \supseteq B_{n+1} \supseteq \dots$$

Then

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} B_n\right) = \lim_{n \rightarrow \infty} \mathbb{P}(B_n).$$

Note: because $B_1 \supseteq B_2 \supseteq \dots$, we have: $\lim_{n \rightarrow \infty} B_n = \bigcap_{n=1}^{\infty} B_n$.

Proof (a) only: for (b), take complements and use (a).

Define $C_1 = A_1$, and $C_i = A_i \setminus A_{i-1}$ for $i = 2, 3, \dots$. Then C_1, C_2, \dots are mutually exclusive, and $\bigcup_{i=1}^n C_i = \bigcup_{i=1}^n A_i$, and likewise, $\bigcup_{i=1}^{\infty} C_i = \bigcup_{i=1}^{\infty} A_i$.

Thus

$$\begin{aligned} \mathbb{P}\left(\lim_{n \rightarrow \infty} A_n\right) &= \mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \mathbb{P}\left(\bigcup_{i=1}^{\infty} C_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(C_i) \quad (C_i \text{ mutually exclusive}) \\ &= \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{P}(C_i) \\ &= \lim_{n \rightarrow \infty} \mathbb{P}\left(\bigcup_{i=1}^n C_i\right) \\ &= \lim_{n \rightarrow \infty} \mathbb{P}\left(\bigcup_{i=1}^n A_i\right) = \lim_{n \rightarrow \infty} \mathbb{P}(A_n). \quad \square \end{aligned}$$

2.10 Random Variables

Definition: A **random variable**, X , is defined as a *function from the sample space to the real numbers*: $X : \Omega \rightarrow \mathbb{R}$.

A random variable therefore *assigns a real number to every possible outcome of a random experiment*.

A random variable is essentially *a rule or mechanism for generating random real numbers*.

The Distribution Function

Definition: The **cumulative distribution function** of a random variable X is given by

$$F_X(x) = \mathbb{P}(X \leq x)$$

$F_X(x)$ is often referred to as simply the **distribution function**.

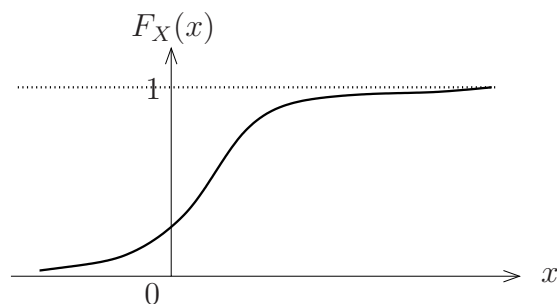
Properties of the distribution function

- 1) $F_X(-\infty) = \mathbb{P}(X \leq -\infty) = 0$.
 $F_X(+\infty) = \mathbb{P}(X \leq \infty) = 1$.
 - 2) $F_X(x)$ is a non-decreasing function of x :
if $x_1 < x_2$, then $F_X(x_1) \leq F_X(x_2)$.
 - 3) If $b > a$, then $\mathbb{P}(a < X \leq b) = F_X(b) - F_X(a)$.
 - 4) F_X is right-continuous: i.e. $\lim_{h \downarrow 0} F_X(x + h) = F_X(x)$.
-

2.11 Continuous Random Variables

Definition: The random variable X is continuous if *the distribution function $F_X(x)$ is a continuous function.*

In practice, this means that a continuous random variable *takes values in a continuous subset of \mathbb{R} : e.g. $X : \Omega \rightarrow [0, 1]$ or $X : \Omega \rightarrow [0, \infty)$.*

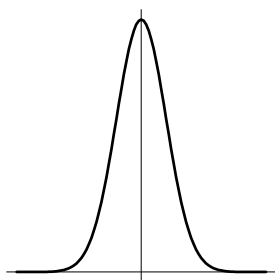


Probability Density Function for continuous random variables

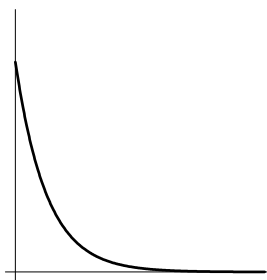
Definition: Let X be a continuous random variable with continuous distribution function $F_X(x)$. The probability density function (p.d.f.) of X is defined as

$$f_X(x) = F'_X(x) = \frac{d}{dx}(F_X(x))$$

The pdf, $f_X(x)$, gives the *shape* of the distribution of X .



Normal distribution



Exponential distribution



Gamma distribution

By the Fundamental Theorem of Calculus, the distribution function $F_X(x)$ can be written in terms of the probability density function, $f_X(x)$, as follows:

$$F_X(x) = \int_{-\infty}^x f_X(u) du$$

Endpoints of intervals

For continuous random variables, every point x has $\mathbb{P}(X = x) = 0$. This means that the endpoints of intervals are not important for continuous random variables.

Thus, $\mathbb{P}(a \leq X \leq b) = \mathbb{P}(a < X \leq b) = \mathbb{P}(a \leq X < b) = \mathbb{P}(a < X < b)$.

This is *only* true for *continuous* random variables.

Calculating probabilities for continuous random variables

To calculate $\mathbb{P}(a \leq X \leq b)$, use *either*

$$\mathbb{P}(a \leq X \leq b) = F_X(b) - F_X(a)$$

or

$$\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(x) dx$$

Example: Let X be a continuous random variable with p.d.f.

$$f_X(x) = \begin{cases} 2x^{-2} & \text{for } 1 < x < 2, \\ 0 & \text{otherwise.} \end{cases}$$

- Find the cumulative distribution function, $F_X(x)$.
- Find $\mathbb{P}(X \leq 1.5)$.

$$a) \quad F_X(x) = \int_{-\infty}^x f_X(u) du = \int_1^x 2u^{-2} du = \left[\frac{2u^{-1}}{-1} \right]_1^x = 2 - \frac{2}{x} \quad \text{for } 1 < x < 2.$$

Thus

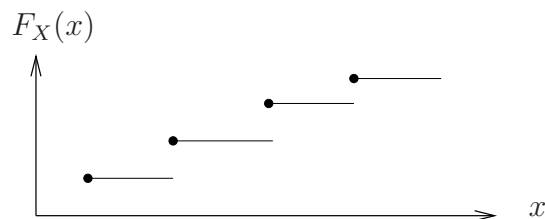
$$F_X(x) = \begin{cases} 0 & \text{for } x \leq 1, \\ 2 - \frac{2}{x} & \text{for } 1 < x < 2, \\ 1 & \text{for } x \geq 2. \end{cases}$$

$$b) \quad \mathbb{P}(X \leq 1.5) = F_X(1.5) = 2 - \frac{2}{1.5} = \frac{2}{3}.$$

2.12 Discrete Random Variables

Definition: The random variable X is **discrete** if X takes values in a finite or countable subset of \mathbb{R} : thus, $X : \Omega \rightarrow \{x_1, x_2, \dots\}$.

When X is a discrete random variable, the distribution function $F_X(x)$ is a *step function*.



Probability function

Definition: Let X be a discrete random variable with distribution function $F_X(x)$.

The **probability function** of X is defined as

$$f_X(x) = \mathbb{P}(X = x).$$

Endpoints of intervals

For discrete random variables, *individual points can have* $\mathbb{P}(X = x) > 0$.

This means that *the endpoints of intervals ARE important for discrete random variables*.

For example, if X takes values $0, 1, 2, \dots$, and a, b are integers with $b > a$, then

$$\mathbb{P}(a \leq X \leq b) = \mathbb{P}(a - 1 < X \leq b) = \mathbb{P}(a \leq X < b + 1) = \mathbb{P}(a - 1 < X < b + 1).$$

Calculating probabilities for discrete random variables

To calculate $\mathbb{P}(X \in A)$ for any countable set A , use

$$\mathbb{P}(X \in A) = \sum_{x \in A} \mathbb{P}(X = x).$$

Partition Theorem for probabilities of discrete random variables

Recall the Partition Theorem: for any event A , and for events B_1, B_2, \dots that form a *partition* of Ω ,

$$\mathbb{P}(A) = \sum_y \mathbb{P}(A | B_y) \mathbb{P}(B_y).$$

We can use the Partition Theorem to find probabilities for random variables. Let X and Y be discrete random variables.

- Define event A as $A = \{X = x\}$.
- Define event B_y as $B_y = \{Y = y\}$ for $y = 0, 1, 2, \dots$ (or whatever values Y takes).
- Then, by the Partition Theorem,

$$\mathbb{P}(X = x) = \sum_y \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y).$$

2.13 Independent Random Variables

Random variables X and Y are independent if they have no effect on each other. This means that the probability that they both take specified values simultaneously is the product of the individual probabilities.

Definition: Let X and Y be random variables. The joint distribution function of X and Y is given by

$$F_{X,Y}(x, y) = \mathbb{P}(X \leq x \text{ and } Y \leq y) = \mathbb{P}(X \leq x, Y \leq y).$$

Definition: Let X and Y be any random variables (continuous or discrete). X and Y are independent if

$$F_{X,Y}(x, y) = F_X(x)F_Y(y) \text{ for ALL } x, y \in \mathbb{R}.$$

If X and Y are **discrete**, they are independent if and only if their joint probability function is the product of their individual probability functions:

$$\begin{aligned} \text{Discrete } X, Y \text{ are indept} \quad &\Longleftrightarrow \quad \mathbb{P}(X = x \text{ AND } Y = y) = \mathbb{P}(X = x)\mathbb{P}(Y = y) \\ &\text{for ALL } x, y \\ &\Longleftrightarrow \quad f_{X,Y}(x, y) = f_X(x)f_Y(y) \quad \text{for ALL } x, y. \end{aligned}$$

Chapter 3: Expectation and Variance

In the previous chapter we looked at probability, with three major themes:

1. Conditional probability: $\mathbb{P}(A | B)$.
2. First-step analysis for calculating eventual probabilities in a stochastic process.
3. Calculating probabilities for continuous and discrete random variables.

In this chapter, we look at the same themes for **expectation** and **variance**. The expectation of a random variable is the *long-term average of the random variable*.

Imagine observing many thousands of independent random values from the random variable of interest. Take the average of these random values. The expectation is the value of this average as the sample size tends to infinity.

We will repeat the three themes of the previous chapter, but in a different order.

1. Calculating expectations for continuous and discrete random variables.
2. Conditional expectation: the expectation of a random variable X , *conditional* on the value taken by another random variable Y . If the value of Y affects the value of X (i.e. X and Y are *dependent*), the conditional expectation of X given the value of Y will be different from the overall expectation of X .
3. First-step analysis for calculating the expected amount of time needed to reach a particular state in a process (e.g. the expected number of shots before we win a game of tennis).

We will also study similar themes for variance.

3.1 Expectation

The mean, expected value, or expectation of a random variable X is written as $\mathbb{E}(X)$ or μ_X . If we observe N random values of X , then the mean of the N values will be approximately equal to $\mathbb{E}(X)$ for large N . The expectation is defined differently for continuous and discrete random variables.

Definition: Let X be a continuous random variable with p.d.f. $f_X(x)$. The expected value of X is

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x f_X(x) dx.$$

Definition: Let X be a discrete random variable with probability function $f_X(x)$. The expected value of X is

$$\mathbb{E}(X) = \sum_x x f_X(x) = \sum_x x \mathbb{P}(X = x).$$

Expectation of $g(X)$

Let $g(X)$ be a function of X . We can imagine a long-term average of $g(X)$ just as we can imagine a long-term average of X . This average is written as $\mathbb{E}(g(X))$. Imagine observing X many times (N times) to give results x_1, x_2, \dots, x_N . Apply the function g to each of these observations, to give $g(x_1), \dots, g(x_N)$. The mean of $g(x_1), g(x_2), \dots, g(x_N)$ approaches $\mathbb{E}(g(X))$ as the number of observations N tends to infinity.

Definition: Let X be a continuous random variable, and let g be a function. The expected value of $g(X)$ is

$$\mathbb{E}(g(X)) = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

Definition: Let X be a discrete random variable, and let g be a function. The expected value of $g(X)$ is

$$\mathbb{E}(g(X)) = \sum_x g(x) f_X(x) = \sum_x g(x) \mathbb{P}(X = x).$$

Expectation of XY : the definition of $\mathbb{E}(XY)$

Suppose we have two random variables, X and Y . These might be independent, in which case the value of X has no effect on the value of Y . Alternatively, X and Y might be *dependent*: when we observe a random value for X , it might influence the random values of Y that we are most likely to observe. For example, X might be the height of a randomly selected person, and Y might be the weight. On the whole, larger values of X will be associated with larger values of Y .

To understand what $\mathbb{E}(XY)$ means, think of observing a large number of *pairs* $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$. If X and Y are dependent, the value x_i might affect the value y_i , and vice versa, so we have to keep the observations together in their pairings. As the number of pairs N tends to infinity, the average $\frac{1}{N} \sum_{i=1}^N x_i \times y_i$ approaches the expectation $\mathbb{E}(XY)$.

For example, if X is height and Y is weight, $\mathbb{E}(XY)$ is the average of (height \times weight). We are interested in $\mathbb{E}(XY)$ because it is used for calculating the *covariance* and *correlation*, which are measures of how closely related X and Y are (see Section 3.2).

Properties of Expectation

- i) Let g and h be functions, and let a and b be constants. For any random variable X (discrete or continuous),

$$\mathbb{E}\{ag(X) + bh(X)\} = a\mathbb{E}\{g(X)\} + b\mathbb{E}\{h(X)\}.$$

In particular,

$$\mathbb{E}(aX + b) = a\mathbb{E}(X) + b.$$

- ii) Let X and Y be ANY random variables (discrete, continuous, independent, or non-independent). Then

$$\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y).$$

More generally, for ANY random variables X_1, \dots, X_n ,

$$\mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}(X_1) + \dots + \mathbb{E}(X_n).$$

iii) Let X and Y be **independent** random variables, and g, h be functions. Then

$$\begin{aligned}\mathbb{E}(XY) &= \mathbb{E}(X)\mathbb{E}(Y) \\ \mathbb{E}(g(X)h(Y)) &= \mathbb{E}(g(X))\mathbb{E}(h(Y)).\end{aligned}$$

Notes: 1. $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ is ONLY generally true if X and Y are **INDEPENDENT**.

2. If X and Y are independent, then $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$. However, the converse is not generally true: it is possible for $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ even though X and Y are dependent.

Probability as an Expectation

Let A be any event. We can write $\mathbb{P}(A)$ as an expectation, as follows. Define the **indicator function**:

$$I_A = \begin{cases} 1 & \text{if event } A \text{ occurs,} \\ 0 & \text{otherwise.} \end{cases}$$

Then I_A is a **random variable**, and

$$\begin{aligned}\mathbb{E}(I_A) &= \sum_{r=0}^1 r\mathbb{P}(I_A = r) \\ &= 0 \times \mathbb{P}(I_A = 0) + 1 \times \mathbb{P}(I_A = 1) \\ &= \mathbb{P}(I_A = 1) \\ &= \mathbb{P}(A).\end{aligned}$$

Thus

$$\mathbb{P}(A) = \mathbb{E}(I_A) \text{ for any event } A.$$

3.2 Variance, covariance, and correlation

The variance of a random variable X is a measure of how *spread out* it is. Are the values of X clustered tightly around their mean, or can we commonly observe values of X a long way from the mean value? The *variance* measures how far the values of X are from their mean, on average.

Definition: Let X be any random variable. The **variance** of X is

$$\text{Var}(X) = \mathbb{E}\left((X - \mu_X)^2\right) = \mathbb{E}(X^2) - \left(\mathbb{E}(X)\right)^2.$$

The variance is the *mean squared deviation* of a random variable from its own mean.

If X has *high variance*, we can observe values of X a long way from the mean.

If X has *low variance*, the values of X tend to be clustered tightly around the mean value.

Example: Let X be a continuous random variable with p.d.f.

$$f_X(x) = \begin{cases} 2x^{-2} & \text{for } 1 < x < 2, \\ 0 & \text{otherwise.} \end{cases}$$

Find $\mathbb{E}(X)$ and $\text{Var}(X)$.

$$\begin{aligned} E(X) &= \int_{-\infty}^{\infty} x f_X(x) dx = \int_1^2 x \times 2x^{-2} dx = \int_1^2 2x^{-1} dx \\ &= \left[2 \log(x) \right]_1^2 \\ &= 2 \log(2) - 2 \log(1) \\ &= 2 \log(2). \end{aligned}$$

For $\text{Var}(X)$, we use

$$\text{Var}(X) = \mathbb{E}(X^2) - \{\mathbb{E}(X)\}^2.$$

Now

$$\begin{aligned} E(X^2) &= \int_{-\infty}^{\infty} x^2 f_X(x) dx = \int_1^2 x^2 \times 2x^{-2} dx = \int_1^2 2 dx \\ &= \left[2x \right]_1^2 \\ &= 2 \times 2 - 2 \times 1 \\ &= 2. \end{aligned}$$

Thus

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}(X^2) - \{\mathbb{E}(X)\}^2 \\ &= 2 - \{2 \log(2)\}^2 \\ &= 0.0782. \end{aligned}$$

Covariance

Covariance is a measure of the association or dependence between two random variables X and Y . Covariance can be either positive or negative. (*Variance* is always positive.)

Definition: Let X and Y be any random variables. The covariance between X and Y is given by

$$\text{cov}(X, Y) = \mathbb{E}\{(X - \mu_X)(Y - \mu_Y)\} = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y),$$

where $\mu_X = \mathbb{E}(X)$, $\mu_Y = \mathbb{E}(Y)$.

1. $\text{cov}(X, Y)$ will be **positive** if large values of X tend to occur with large values of Y , and small values of X tend to occur with small values of Y . For example, if X is height and Y is weight of a randomly selected person, we would expect $\text{cov}(X, Y)$ to be positive.

2. $\text{cov}(X, Y)$ will be **negative** if large values of X tend to occur with small values of Y , and small values of X tend to occur with large values of Y . For example, if X is age of a randomly selected person, and Y is heart rate, we would expect X and Y to be negatively correlated (older people have slower heart rates).
3. If X and Y are independent, then there is no pattern between large values of X and large values of Y , so $\text{cov}(X, Y) = 0$. However, $\text{cov}(X, Y) = 0$ does NOT imply that X and Y are independent, unless X and Y are Normally distributed.

Properties of Variance

- i) Let g be a function, and let a and b be constants. For any random variable X (discrete or continuous),

$$\text{Var}\{ag(X) + b\} = a^2 \text{Var}\{g(X)\}.$$

In particular, $\text{Var}(aX + b) = a^2 \text{Var}(X)$.

- ii) Let X and Y be **independent** random variables. Then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

- iii) If X and Y are **NOT independent**, then

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y).$$

Correlation (non-examinable)

The correlation coefficient of X and Y is a measure of the linear association between X and Y . It is given by the covariance, scaled by the overall variability in X and Y . As a result, the correlation coefficient is always between -1 and $+1$, so it is easily compared for different quantities.

Definition: The **correlation** between X and Y , also called the **correlation coefficient**, is given by

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

The correlation measures linear association between X and Y . It takes values only between -1 and $+1$, and has the same sign as the covariance.

The correlation is ± 1 if and only if there is a perfect linear relationship between X and Y , i.e. $\text{corr}(X, Y) = 1 \iff Y = aX + b$ for some constants a and b .

The correlation is 0 if X and Y are independent, but a correlation of 0 does not *imply* that X and Y are independent.

3.3 Conditional Expectation and Conditional Variance

Throughout this section, we will assume for simplicity that X and Y are discrete random variables. However, exactly the same results hold for continuous random variables too.

Suppose that X and Y are discrete random variables, possibly dependent on each other. Suppose that we fix Y at the value y . This gives us a set of conditional probabilities $\mathbb{P}(X = x | Y = y)$ for all possible values x of X . This is called the **conditional distribution of X , given that $Y = y$** .

Definition: Let X and Y be discrete random variables. The **conditional probability function** of X , given that $Y = y$, is:

$$\mathbb{P}(X = x | Y = y) = \frac{\mathbb{P}(X = x \text{ AND } Y = y)}{\mathbb{P}(Y = y)}.$$

We write the conditional probability function as:

$$f_{X|Y}(x | y) = \mathbb{P}(X = x | Y = y).$$

Note: The conditional probabilities $f_{X|Y}(x | y)$ sum to one, just like any other probability function:

$$\sum_x \mathbb{P}(X = x | Y = y) = \sum_x \mathbb{P}_{\{Y=y\}}(X = x) = 1,$$

using the subscript notation $\mathbb{P}_{\{Y=y\}}$ of Section 2.3.

We can also find the expectation and variance of X with respect to this conditional distribution. That is, if we know that the value of Y is fixed at y , then we can find the mean value of X *given that* Y takes the value y , and also the variance of X given that $Y = y$.

Definition: Let X and Y be discrete random variables. The **conditional expectation of X , given that $Y = y$** , is

$$\mu_{X|Y=y} = \mathbb{E}(X | Y = y) = \sum_x x f_{X|Y}(x | y).$$

$\mathbb{E}(X | Y = y)$ is the *mean value of X , when Y is fixed at y* .

Conditional expectation as a random variable

The unconditional expectation of X , $\mathbb{E}(X)$, is just *a number*:
e.g. $\mathbb{E}X = 2$ or $\mathbb{E}X = 5.8$.

The conditional expectation, $\mathbb{E}(X | Y = y)$, is *a number depending on y* .

If Y has an influence on the value of X , then Y will have an influence on the *average* value of X . So, for example, we would expect $\mathbb{E}(X | Y = 2)$ to be different from $\mathbb{E}(X | Y = 3)$.

We can therefore view $\mathbb{E}(X | Y = y)$ as a *function of y* , say $\mathbb{E}(X | Y=y) = h(y)$.

To evaluate this function, $h(y) = \mathbb{E}(X | Y = y)$, we:

- i) *fix Y at the chosen value y* ;
- ii) *find the expectation of X when Y is fixed at this value.*

However, we could also evaluate the function at a *random value* of Y :

- i) *observe a random value of Y* ;
- ii) *fix Y at that observed random value*;
- iii) *evaluate $\mathbb{E}(X | Y = \text{observed random value})$.*

We obtain a random variable: $\mathbb{E}(X | Y) = h(Y)$.

The randomness comes from the randomness in Y , not in X .

Conditional expectation, $\mathbb{E}(X | Y)$, is a random variable with randomness inherited from Y , not X .

Example: Suppose $Y = \begin{cases} 1 & \text{with probability } 1/8, \\ 2 & \text{with probability } 7/8, \end{cases}$

and $X | Y = \begin{cases} 2Y & \text{with probability } 3/4, \\ 3Y & \text{with probability } 1/4. \end{cases}$

Conditional expectation of X given $Y = y$ is a number depending on y :

If $Y = 1$, then: $X | (Y = 1) = \begin{cases} 2 & \text{with probability } 3/4 \\ 3 & \text{with probability } 1/4 \end{cases}$

$$\text{so } \mathbb{E}(X | Y = 1) = 2 \times \frac{3}{4} + 3 \times \frac{1}{4} = \frac{9}{4}.$$

If $Y = 2$, then: $X | (Y = 2) = \begin{cases} 4 & \text{with probability } 3/4 \\ 6 & \text{with probability } 1/4 \end{cases}$

$$\text{so } \mathbb{E}(X | Y = 2) = 4 \times \frac{3}{4} + 6 \times \frac{1}{4} = \frac{18}{4}.$$

Thus $\mathbb{E}(X | Y = y) = \begin{cases} 9/4 & \text{if } y = 1 \\ 18/4 & \text{if } y = 2. \end{cases}$

So $\mathbb{E}(X | Y = y)$ is a number depending on y , or a function of y .

Conditional expectation of X given random Y is a random variable:

From above, $\mathbb{E}(X | Y) = \begin{cases} 9/4 & \text{if } Y = 1 \text{ (probability } 1/8), \\ 18/4 & \text{if } Y = 2 \text{ (probability } 7/8). \end{cases}$

So $\mathbb{E}(X | Y) = \begin{cases} 9/4 & \text{with probability } 1/8, \\ 18/4 & \text{with probability } 7/8. \end{cases}$

Thus $\mathbb{E}(X | Y)$ is a random variable.

The randomness in $\mathbb{E}(X | Y)$ is inherited from Y , not from X .

Conditional expectation is a very useful tool for finding the **unconditional** expectation of X (see below). Just like the Partition Theorem, it is useful because it is often easier to specify conditional probabilities than to specify overall probabilities.

Conditional variance

The conditional variance is similar to the conditional expectation.

- $\text{Var}(X | Y = y)$ is the variance of X , when Y is fixed at the value $Y = y$.
- $\text{Var}(X | Y)$ is a random variable, giving the variance of X when Y is fixed at a value to be selected randomly.

Definition: Let X and Y be random variables. The conditional variance of X , given Y , is given by

$$\text{Var}(X | Y) = \mathbb{E}(X^2 | Y) - \left\{ \mathbb{E}(X | Y) \right\}^2 = \mathbb{E} \left\{ (X - \mu_{X|Y})^2 | Y \right\}$$

Like expectation, $\text{Var}(X | Y = y)$ is a number depending on y (a function of y), while $\text{Var}(X | Y)$ is a random variable with randomness inherited from Y .

Laws of Total Expectation and Variance

If all the expectations below are finite, then for ANY random variables X and Y , we have:

i) $\boxed{\mathbb{E}(X) = \mathbb{E}_Y(\mathbb{E}(X | Y))}$ *Law of Total Expectation.*

Note that we can pick any r.v. Y , to make the expectation as easy as we can.

ii) $\mathbb{E}(g(X)) = \mathbb{E}_Y(\mathbb{E}(g(X) | Y))$ *for any function g .*

iii) $\boxed{\text{Var}(X) = \mathbb{E}_Y(\text{Var}(X | Y)) + \text{Var}_Y(\mathbb{E}(X | Y))}$

Law of Total Variance.

Note: \mathbb{E}_Y and Var_Y denote expectation over Y and variance over Y ,

i.e. the expectation or variance is computed over the distribution of the random variable Y .

The Law of Total Expectation says that *the total average is the average of case-by-case averages*.

- The total average is $\mathbb{E}(X)$;
- The case-by-case averages are $\mathbb{E}(X | Y)$ *for the different values of Y* ;
- The average of case-by-case averages is *the average over Y of the Y -case averages*: $\mathbb{E}_Y(\mathbb{E}(X | Y))$.

Example: In the example above, we had: $\mathbb{E}(X | Y) = \begin{cases} 9/4 & \text{with probability } 1/8, \\ 18/4 & \text{with probability } 7/8. \end{cases}$

The total average is:

$$\mathbb{E}(X) = \mathbb{E}_Y \left\{ \mathbb{E}(X | Y) \right\} = \frac{9}{4} \times \frac{1}{8} + \frac{18}{4} \times \frac{7}{8} = 4.22.$$

Proof of (i), (ii), (iii):

(i) is a special case of (ii), so we just need to prove (ii). Begin at RHS:

$$\begin{aligned} \text{RHS} &= \mathbb{E}_Y \left[\mathbb{E}(g(X) | Y) \right] = \mathbb{E}_Y \left[\sum_x g(x) \mathbb{P}(X = x | Y) \right] \\ &= \sum_y \left[\sum_x g(x) \mathbb{P}(X = x | Y = y) \right] \mathbb{P}(Y = y) \\ &= \sum_y \sum_x g(x) \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y) \\ &= \sum_x g(x) \sum_y \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y) \\ &= \sum_x g(x) \mathbb{P}(X = x) \quad (\text{partition rule}) \\ &= \mathbb{E}(g(X)) = \text{LHS}. \end{aligned}$$

(iii) Wish to prove $\text{Var}(X) = \mathbb{E}_Y[\text{Var}(X | Y)] + \text{Var}_Y[\mathbb{E}(X | Y)]$. Begin at RHS:

$$\begin{aligned} &\mathbb{E}_Y[\text{Var}(X | Y)] + \text{Var}_Y[\mathbb{E}(X | Y)] \\ &= \mathbb{E}_Y \left\{ \mathbb{E}(X^2 | Y) - (\mathbb{E}(X | Y))^2 \right\} + \left\{ \mathbb{E}_Y \left\{ [\mathbb{E}(X | Y)]^2 \right\} - \left[\underbrace{\mathbb{E}_Y(\mathbb{E}(X | Y))}_{\mathbb{E}(X) \text{ by part (i)}} \right]^2 \right\} \\ &= \underbrace{\mathbb{E}_Y \{ \mathbb{E}(X^2 | Y) \}}_{\mathbb{E}(X^2) \text{ by part (i)}} - \mathbb{E}_Y \{ [\mathbb{E}(X | Y)]^2 \} + \mathbb{E}_Y \{ [\mathbb{E}(X | Y)]^2 \} - (\mathbb{E}X)^2 \\ &= \mathbb{E}(X^2) - (\mathbb{E}X)^2 \\ &= \text{Var}(X) = \text{LHS}. \quad \square \end{aligned}$$

3.4 Examples of Conditional Expectation and Variance

1. Swimming with dolphins

Fraser runs a dolphin-watch business. Every day, he is unable to run the trip due to bad weather with probability p , independently of all other days. Fraser works every day except the bad-weather days, which he takes as holiday.



Let Y be the number of consecutive days Fraser has to work between bad-weather days. Let X be the total number of customers who go on Fraser's trip in this period of Y days. Conditional on Y , the distribution of X is

$$(X | Y) \sim \text{Poisson}(\mu Y).$$

- (a) Name the distribution of Y , and state $\mathbb{E}(Y)$ and $\text{Var}(Y)$.
- (b) Find the expectation and the variance of the number of customers Fraser sees between bad-weather days, $\mathbb{E}(X)$ and $\text{Var}(X)$.

- (a) Let 'success' be 'bad-weather day' and 'failure' be 'work-day'.

Then $\mathbb{P}(\text{success}) = \mathbb{P}(\text{bad-weather}) = p$.

Y is the number of failures before the first success.

So

$$Y \sim \text{Geometric}(p).$$

Thus

$$\mathbb{E}(Y) = \frac{1-p}{p},$$

$$\text{Var}(Y) = \frac{1-p}{p^2}.$$

- (b) We know $(X | Y) \sim \text{Poisson}(\mu Y)$: so

$$\mathbb{E}(X | Y) = \text{Var}(X | Y) = \mu Y.$$

By the Law of Total Expectation:

$$\begin{aligned}\mathbb{E}(X) &= \mathbb{E}_Y\left\{\mathbb{E}(X | Y)\right\} \\ &= \mathbb{E}_Y(\mu Y) \\ &= \mu \mathbb{E}_Y(Y) \\ \therefore \mathbb{E}(X) &= \frac{\mu(1-p)}{p}.\end{aligned}$$

By the Law of Total Variance:

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}_Y\left(\text{Var}(X | Y)\right) + \text{Var}_Y\left(\mathbb{E}(X | Y)\right) \\ &= \mathbb{E}_Y(\mu Y) + \text{Var}_Y(\mu Y) \\ &= \mu \mathbb{E}_Y(Y) + \mu^2 \text{Var}_Y(Y) \\ &= \mu \left(\frac{1-p}{p}\right) + \mu^2 \left(\frac{1-p}{p^2}\right) \\ &= \frac{\mu(1-p)(p+\mu)}{p^2}.\end{aligned}$$

Checking your answer in R:

If you know how to use a statistical package like *R*, you can check your answer to the question above as follows.

```
> # Pick a value for p, e.g. p = 0.2.
> # Pick a value for mu, e.g. mu = 25
>
> # Generate 10,000 random values of Y ~ Geometric(p = 0.2):
> y <- rgeom(10000, prob=0.2)
>
> # Generate 10,000 random values of X conditional on Y:
> # use (X | Y) ~ Poisson(mu * Y) ~ Poisson(25 * Y)
> x <- rpois(10000, lambda = 25*y)
```

```
> # Find the sample mean of X (should be close to E(X)):
> mean(x)
[1] 100.6606
>
> # Find the sample variance of X (should be close to var(X)):
> var(x)
[1] 12624.47
>
> # Check the formula for E(X):
> 25 * (1 - 0.2) / 0.2
[1] 100
>
> # Check the formula for var(X):
> 25 * (1 - 0.2) * (0.2 + 25) / 0.2^2
[1] 12600
```

The formulas we obtained by working give $\mathbb{E}(X) = 100$ and $\text{Var}(X) = 12600$. The sample mean was $\bar{x} = 100.6606$ (close to 100), and the sample variance was 12624.47 (close to 12600). Thus our working seems to have been correct.

2. Randomly stopped sum

This model arises very commonly in stochastic processes. A random number N of events occur, and each event i has associated with it some cost, penalty, or reward X_i . The question is to find the mean and variance of the total cost / reward:

$$T_N = X_1 + X_2 + \dots + X_N.$$

The difficulty is that the number N of terms in the sum is itself random.

T_N is called a *randomly stopped sum*: it is a sum of X_i 's, randomly stopped at the random number of N terms.



Example: Think of a cash machine, which has to be loaded with enough money to cover the day's business. The number of customers per day is a random number N . Customer i withdraws a random amount X_i . The total amount withdrawn during the day is a randomly stopped sum: $T_N = X_1 + \dots + X_N$.

Cash machine example

The citizens of Remuera withdraw money from a cash machine according to the following probability function (X):

Amount, x (\$)	50	100	200
$\mathbb{P}(X = x)$	0.3	0.5	0.2

The number of customers per day has the distribution $N \sim \text{Poisson}(\lambda)$.

Let $T_N = X_1 + X_2 + \dots + X_N$ be the total amount of money withdrawn in a day, where each X_i has the probability function above, and X_1, X_2, \dots are independent of each other and of N .

T_N is a randomly stopped sum, stopped by the random number of N customers.

- Show that $\mathbb{E}(X) = 105$, and $\text{Var}(X) = 2725$.
- Find $\mathbb{E}(T_N)$ and $\text{Var}(T_N)$: the mean and variance of the amount of money withdrawn each day.

Solution

- Exercise.
- Let $T_N = \sum_{i=1}^N X_i$. If we knew how many terms were in the sum, we could easily find $\mathbb{E}(T_N)$ and $\text{Var}(T_N)$ as the mean and variance of a sum of independent r.v.s. So ‘pretend’ we know how many terms are in the sum: i.e. condition on N .

$$\begin{aligned}
 \mathbb{E}(T_N | N) &= \mathbb{E}(X_1 + X_2 + \dots + X_N | N) \\
 &= \mathbb{E}(X_1 + X_2 + \dots + X_N) \\
 &\quad \text{(because all } X_i\text{'s are independent of } N\text{)} \\
 &= \mathbb{E}(X_1) + \mathbb{E}(X_2) + \dots + \mathbb{E}(X_N) \\
 &\quad \text{where } N \text{ is now considered constant;} \\
 &\quad \text{(we do NOT need independence of } X_i\text{'s for this)} \\
 &= N \times \mathbb{E}(X) \quad \text{(because all } X_i\text{'s have same mean, } \mathbb{E}(X)\text{)} \\
 &= 105N.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \text{Var}(T_N | N) &= \text{Var}(X_1 + X_2 + \dots + X_N | N) \\
 &= \text{Var}(X_1 + X_2 + \dots + X_N) \\
 &\quad \text{where } N \text{ is now considered constant;} \\
 &\quad \text{(because all } X_i \text{'s are independent of } N\text{)} \\
 &= \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_N) \\
 &\quad \text{(we DO need independence of } X_i \text{'s for this)} \\
 &= N \times \text{Var}(X) \quad \text{(because all } X_i \text{'s have same variance, } \text{Var}(X)\text{)} \\
 &= 2725N.
 \end{aligned}$$

So

$$\begin{aligned}
 \mathbb{E}(T_N) &= \mathbb{E}_N \left\{ \mathbb{E}(T_N | N) \right\} \\
 &= \mathbb{E}_N(105N) \\
 &= 105\mathbb{E}_N(N) \\
 &= 105\lambda,
 \end{aligned}$$

because $N \sim \text{Poisson}(\lambda)$ so $\mathbb{E}(N) = \lambda$.

Similarly,

$$\begin{aligned}
 \text{Var}(T_N) &= \mathbb{E}_N \left\{ \text{Var}(T_N | N) \right\} + \text{Var}_N \left\{ \mathbb{E}(T_N | N) \right\} \\
 &= \mathbb{E}_N \{2725N\} + \text{Var}_N \{105N\} \\
 &= 2725\mathbb{E}_N(N) + 105^2 \text{Var}_N(N) \\
 &= 2725\lambda + 11025\lambda \\
 &= 13750\lambda,
 \end{aligned}$$

because $N \sim \text{Poisson}(\lambda)$ so $\mathbb{E}(N) = \text{Var}(N) = \lambda$.

Check in R (advanced)

```
> # Create a function tn.func to calculate a single value of T_N
> # for a given value N=n:
> tn.func <- function(n){
  sum(sample(c(50, 100, 200), n, replace=T,
    prob=c(0.3, 0.5, 0.2)))
}

> # Generate 10,000 random values of N, using lambda=50:
> N <- rpois(10000, lambda=50)
> # Generate 10,000 random values of T_N, conditional on N:
> TN <- sapply(N, tn.func)
> # Find the sample mean of T_N values, which should be close to
> # 105 * 50 = 5250:
> mean(TN)
[1] 5253.255
> # Find the sample variance of T_N values, which should be close
> # to 13750 * 50 = 687500:
> var(TN)
[1] 682469.4
```

All seems well. Note that the sample variance is often some distance from the true variance, even when the sample size is 10,000.

General result for randomly stopped sums:

Suppose X_1, X_2, \dots each have the same mean μ and variance σ^2 , and X_1, X_2, \dots , and N are mutually independent. Let $T_N = X_1 + \dots + X_N$ be the randomly stopped sum. By following similar working to that above:

$$\mathbb{E}(T_N) = \mathbb{E} \left\{ \sum_{i=1}^N X_i \right\} = \mu \mathbb{E}(N)$$

$$\text{Var}(T_N) = \text{Var} \left\{ \sum_{i=1}^N X_i \right\} = \sigma^2 \mathbb{E}(N) + \mu^2 \text{Var}(N).$$

3.5 First-Step Analysis for calculating expected reaching times

Remember from Section 2.6 that we use First-Step Analysis for finding the probability of eventually reaching a particular state in a stochastic process. First-step analysis for probabilities uses *conditional probability and the Partition Theorem (Law of Total Probability)*.

In the same way, we can use first-step analysis for finding the *expected reaching time for a state*.

This is the expected number of steps that will be needed to reach a particular state from a specified start-point, or the expected length of time it will take to get there if we have a continuous time process.

Just as first-step analysis for probabilities uses conditional probability and the law of total probability (Partition Theorem), first-step analysis for expectations uses *conditional expectation and the law of total expectation*.

First-step analysis for probabilities:

The first-step analysis procedure for probabilities can be summarized as follows:

$$\mathbb{P}(\text{eventual goal}) = \sum_{\substack{\text{first-step} \\ \text{options}}} \mathbb{P}(\text{eventual goal} \mid \text{option}) \mathbb{P}(\text{option}) .$$

This is because the first-step options form a *partition of the sample space*.

First-step analysis for expected reaching times:

The expression for expected reaching times is very similar:

$$\mathbb{E}(\text{reaching time}) = \sum_{\substack{\text{first-step} \\ \text{options}}} \mathbb{E}(\text{reaching time} \mid \text{option}) \mathbb{P}(\text{option}) .$$

This follows immediately from the law of total expectation:

$$\mathbb{E}(X) = \mathbb{E}_Y \left\{ \mathbb{E}(X | Y) \right\} = \sum_y \mathbb{E}(X | Y = y) \mathbb{P}(Y = y).$$

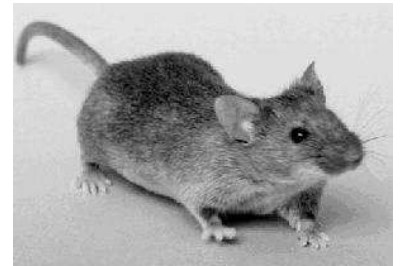
Let X be the reaching time, and let Y be the label for possible options:
i.e. $Y = 1, 2, 3, \dots$ for options 1, 2, 3, \dots

We then obtain:

$$\begin{aligned} \mathbb{E}(X) &= \sum_y \mathbb{E}(X | Y = y) \mathbb{P}(Y = y) \\ \text{i.e. } \mathbb{E}(\text{reaching time}) &= \sum_{\substack{\text{first-step} \\ \text{options}}} \mathbb{E}(\text{reaching time} | \text{option}) \mathbb{P}(\text{option}). \end{aligned}$$

Example 1: Mouse in a Maze

A mouse is trapped in a room with three exits at the centre of a maze.

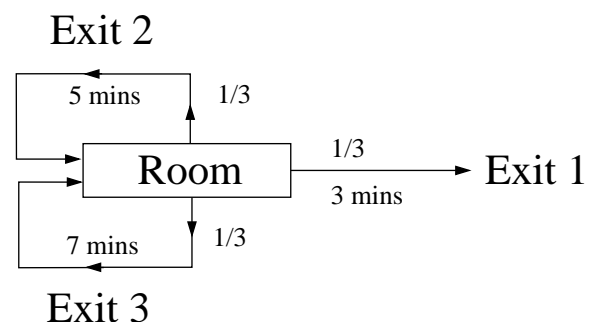


- Exit 1 leads outside the maze after 3 minutes.
- Exit 2 leads back to the room after 5 minutes.
- Exit 3 leads back to the room after 7 minutes.

Every time the mouse makes a choice, it is equally likely to choose any of the three exits. What is the expected time taken for the mouse to leave the maze?

Let X = time taken for mouse to leave maze, starting from room R .

Let Y = exit the mouse chooses first (1, 2, or 3).



Then:

$$\begin{aligned}\mathbb{E}(X) &= \mathbb{E}_Y(\mathbb{E}(X | Y)) \\ &= \sum_{y=1}^3 \mathbb{E}(X | Y = y) \mathbb{P}(Y = y) \\ &= \mathbb{E}(X | Y = 1) \times \frac{1}{3} + \mathbb{E}(X | Y = 2) \times \frac{1}{3} + \mathbb{E}(X | Y = 3) \times \frac{1}{3}.\end{aligned}$$

But:

$$\mathbb{E}(X | Y = 1) = 3 \text{ minutes}$$

$$\mathbb{E}(X | Y = 2) = 5 + \mathbb{E}(X) \text{ (after 5 mins back in Room, time } \mathbb{E}(X) \text{ to get out)}$$

$$\mathbb{E}(X | Y = 3) = 7 + \mathbb{E}(X) \text{ (after 7 mins, back in Room)}$$

So

$$\begin{aligned}\mathbb{E}(X) &= 3 \times \frac{1}{3} + (5 + \mathbb{E}X) \times \frac{1}{3} + (7 + \mathbb{E}X) \times \frac{1}{3} \\ &= 15 \times \frac{1}{3} + 2(\mathbb{E}X) \times \frac{1}{3} \\ \frac{1}{3} \mathbb{E}(X) &= 15 \times \frac{1}{3} \\ \Rightarrow \mathbb{E}(X) &= 15 \text{ minutes.}\end{aligned}$$

Notation for quick solutions of first-step analysis problems

As for probabilities, first-step analysis for expectations relies on a good notation. The best way to tackle the problem above is as follows.

Define $m_R = \mathbb{E}(\text{time to leave maze} \mid \text{start in Room})$.

First-step analysis:

$$\begin{aligned}m_R &= \frac{1}{3} \times 3 + \frac{1}{3} \times (5 + m_R) + \frac{1}{3} \times (7 + m_R) \\ \Rightarrow 3m_R &= (3 + 5 + 7) + 2m_R \\ \Rightarrow m_R &= 15 \text{ minutes} \quad (\text{as before}).\end{aligned}$$

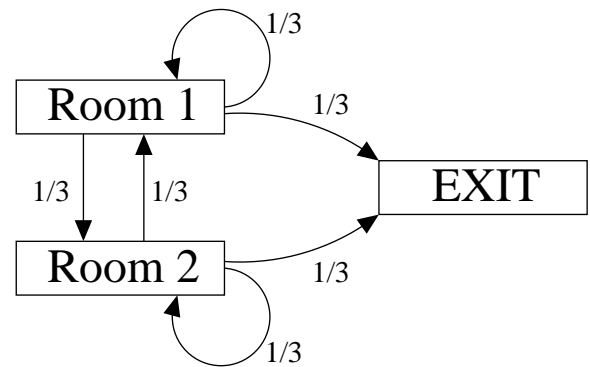
Example 2: Counting the steps

The most common questions involving first-step analysis for expectations ask for the *expected number of steps before finishing*. The number of steps is usually equal to the *number of arrows traversed from the current state to the end*.

The key point to remember is that when we take expectations, we are usually *counting something*.

You must remember to *add on whatever we are counting, to every step taken*.

The mouse is put in a new maze with two rooms, pictured here. Starting from Room 1, what is the expected number of steps the mouse takes before it reaches the exit?



1. Define notation: let

$$m_1 = \mathbb{E}(\text{number of steps to finish} \mid \text{start in Room 1})$$

$$m_2 = \mathbb{E}(\text{number of steps to finish} \mid \text{start in Room 2}).$$

2. First-step analysis:

$$m_1 = \frac{1}{3} \times 1 + \frac{1}{3} (1 + m_1) + \frac{1}{3} (1 + m_2) \quad (a)$$

$$m_2 = \frac{1}{3} \times 1 + \frac{1}{3} (1 + m_1) + \frac{1}{3} (1 + m_2) \quad (b)$$

We could solve as simultaneous equations, as usual, but in this case inspection of (a) and (b) shows immediately that $m_1 = m_2$. Thus:

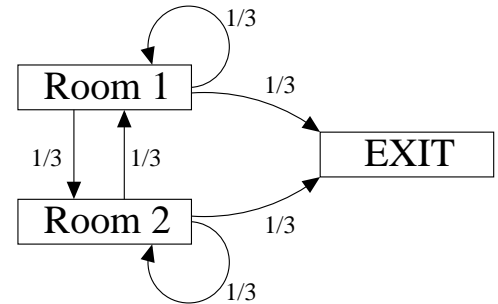
$$(a) \Rightarrow 3m_1 = 3 + 2m_1$$

$$\Rightarrow m_1 = 3 \text{ steps.}$$

Further, $m_2 = m_1 = 3$ steps also.

Incrementing before partitioning

In many problems, all possible first-step options incur the same initial penalty. The last example is such a case, because *every possible step adds 1 to the total number of steps taken*.



In a case where all steps incur the same penalty, there are two ways of proceeding:

1. *Add the penalty onto each option separately: e.g.*

$$m_1 = \frac{1}{3} \times 1 + \frac{1}{3} (1 + m_1) + \frac{1}{3} (1 + m_2).$$

2. *(Usually quicker) Add the penalty once only, at the beginning:*

$$m_1 = 1 + \frac{1}{3} \times 0 + \frac{1}{3} m_1 + \frac{1}{3} m_2.$$

In each case, we will get the same answer (check). This is because the option probabilities sum to 1, *so in Method 1 we are adding* $(\frac{1}{3} + \frac{1}{3} + \frac{1}{3}) \times 1 = 1 \times 1 = 1$, *just as we are in Method 2*.

3.6 Probability as a conditional expectation

Recall from Section 3.1 that for any event A , we can write $\mathbb{P}(A)$ as an expectation as follows.

Define the indicator random variable: $I_A = \begin{cases} 1 & \text{if event } A \text{ occurs,} \\ 0 & \text{otherwise.} \end{cases}$

Then $\mathbb{E}(I_A) = \mathbb{P}(I_A = 1) = \mathbb{P}(A)$.

We can refine this expression further, using the idea of conditional expectation. Let Y be any random variable. Then

$$\mathbb{P}(A) = \mathbb{E}(I_A) = \mathbb{E}_Y(\mathbb{E}(I_A | Y)).$$

But

$$\begin{aligned}\mathbb{E}(I_A | Y) &= \sum_{r=0}^1 r \mathbb{P}(I_A = r | Y) \\ &= 0 \times \mathbb{P}(I_A = 0 | Y) + 1 \times \mathbb{P}(I_A = 1 | Y) \\ &= \mathbb{P}(I_A = 1 | Y) \\ &= \mathbb{P}(A | Y).\end{aligned}$$

Thus

$$\mathbb{P}(A) = \mathbb{E}_Y(\mathbb{E}(I_A | Y)) = \mathbb{E}_Y(\mathbb{P}(A | Y)).$$

This means that for **any** random variable X (discrete or continuous), and for any set of values S (a discrete set or a continuous set), we can write:

- for any **discrete** random variable Y ,

$$\mathbb{P}(X \in S) = \sum_y \mathbb{P}(X \in S | Y = y) \mathbb{P}(Y = y).$$

- for any **continuous** random variable Y ,

$$\mathbb{P}(X \in S) = \int_y \mathbb{P}(X \in S | Y = y) f_Y(y) dy.$$

Example of probability as a conditional expectation: winning a lottery



Suppose that a million people have bought tickets for the weekly lottery draw. Each person has a probability of one-in-a-million of selecting the winning numbers. If more than one person selects the winning numbers, the winner will be chosen at random from all those with matching numbers.

You watch the lottery draw on TV and your numbers match the winners!! You had a one-in-a-million chance, and there were a million players, so it must be YOU, right?

Not so fast. Before you rush to claim your prize, let's calculate the probability that you really will win. You definitely win if you are the only person with matching numbers, but you can also win if there are multiple matching tickets and yours is the one selected at random from the matches.

Define Y to be the number of OTHER matching tickets out of the OTHER 1 million tickets sold. (If you are lucky, $Y = 0$ so you have definitely won.)

If there are 1 million tickets and each ticket has a one-in-a-million chance of having the winning numbers, then

$$Y \sim \text{Poisson}(1) \text{ approximately.}$$

The relationship $Y \sim \text{Poisson}(1)$ arises because of the Poisson approximation to the Binomial distribution.

(a) What is the probability function of Y , $f_Y(y)$?

$$f_Y(y) = \mathbb{P}(Y = y) = \frac{1^y}{y!} e^{-1} = \frac{1}{e \times y!} \quad \text{for } y = 0, 1, 2, \dots$$

(b) What is the probability that yours is the only matching ticket?

$$\mathbb{P}(\text{only one matching ticket}) = \mathbb{P}(Y = 0) = \frac{1}{e} = 0.368.$$

(c) The prize is chosen at random from all those who have matching tickets. What is the probability that you win if there are $Y = y$ OTHER matching tickets?

Let W be the event that I win.

$$\mathbb{P}(W \mid Y = y) = \frac{1}{y + 1}.$$

- (d) Overall, what is the probability that you win, given that you have a matching ticket?

$$\begin{aligned}
 \mathbb{P}(W) &= \mathbb{E}_Y \left\{ \mathbb{P}(W \mid Y = y) \right\} \\
 &= \sum_{y=0}^{\infty} \mathbb{P}(W \mid Y = y) \mathbb{P}(Y = y) \\
 &= \sum_{y=0}^{\infty} \left(\frac{1}{y+1} \right) \left(\frac{1}{e \times y!} \right) \\
 &= \frac{1}{e} \sum_{y=0}^{\infty} \frac{1}{(y+1)y!} \\
 &= \frac{1}{e} \sum_{y=0}^{\infty} \frac{1}{(y+1)!} \\
 &= \frac{1}{e} \left\{ \sum_{y=0}^{\infty} \frac{1}{y!} - \frac{1}{0!} \right\} \\
 &= \frac{1}{e} \{e - 1\} \\
 &= 1 - \frac{1}{e} \\
 &= 0.632.
 \end{aligned}$$

Disappointing?

3.7 Special process: a model for gene spread

Suppose that a particular gene comes in two variants (alleles): A and B. We might be interested in the case where one of the alleles, say A, is harmful — for example it causes a disease. All animals in the population must have either allele A or allele B. We want to know how long it will take before all animals have the same allele, and whether this allele will be the harmful allele A or the safe allele B. This simple model assumes asexual reproduction. It is very similar to the famous Wright-Fisher model, which is a fundamental model of population genetics.

Assumptions:

1. The population stays at constant size N for all generations.
2. At the end of each generation, the N animals create N offspring and then they immediately die.
3. If there are x parents with allele A, and $N - x$ with allele B, then each offspring gets allele A with probability x/N and allele B with $1 - x/N$.
4. All offspring are independent.

Stochastic process:

The **state** of the process at time t is $X_t =$ *the number of animals with allele A at generation t .*

Each X_t could be $0, 1, 2, \dots, N$. The state space is $\{0, 1, 2, \dots, N\}$.

Distribution of $[X_{t+1} | X_t]$

Suppose that $X_t = x$, so x of the animals at generation t have allele A.

Each of the N offspring will get A with probability $\frac{x}{N}$ and B with probability $1 - \frac{x}{N}$.

Thus the number of offspring at time $t+1$ with allele A is: $X_{t+1} \sim \text{Binomial}\left(N, \frac{x}{N}\right)$.

We write this as follows:

$$[X_{t+1} | X_t = x] \sim \text{Binomial}\left(N, \frac{x}{N}\right).$$

If

$$[X_{t+1} | X_t = x] \sim \text{Binomial}\left(N, \frac{x}{N}\right),$$

then

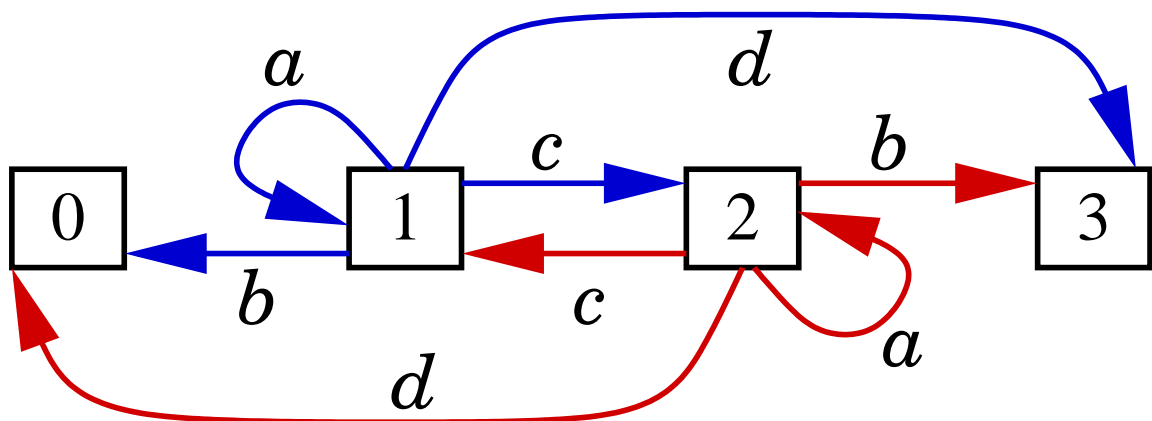
$$\mathbb{P}(X_{t+1} = y | X_t = x) = \binom{N}{y} \left(\frac{x}{N}\right)^y \left(1 - \frac{x}{N}\right)^{N-y} \quad (\text{Binomial formula})$$

Example with $N = 3$

This process becomes complicated to do by hand when N is large. We can use small N to see how to use first-step analysis to answer our questions.

Transition diagram:

Exercise: find the missing probabilities a , b , c , and d when $N = 3$. Express them all as fractions over the same denominator.



Probability the harmful allele A dies out

Suppose the process starts at generation 0. One of the three animals has the harmful allele A. Define a suitable notation, and find the probability that the harmful allele A eventually dies out.

Exercise: answer = $2/3$.

Expected number of generations to fixation

Suppose again that the process starts at generation 0, and one of the three animals has the harmful allele A. Eventually all animals will have the same allele, whether it is allele A or B. When this happens, the population is said to have reached *fixation*: it is fixed for a single allele and no further changes are possible.

Define a suitable notation, and find the expected number of generations to fixation.

Exercise: answer = 3 generations on average.

Things get more interesting for large N . When $N = 100$, and $x = 10$ animals have the harmful allele at generation 0, there is a 90% chance that the harmful allele will die out and a 10% chance that the harmful allele will take over the whole population. The expected number of generations taken to reach fixation is 63.5. If the process starts with just $x = 1$ animal with the harmful allele, there is a 99% chance the harmful allele will die out, but the expected number of generations to fixation is 10.5. Despite the allele being rare, the *average* number of generations for it to either die out or saturate the population is quite large.

Note: The model above is also an example of a process called the ***Voter Process***. The N individuals correspond to N people who each support one of two political candidates, A or B. Every day they make a new decision about whom to support, based on the amount of current support for each candidate. Fixation in the genetic model corresponds to consensus in the Voter Process.

Chapter 4: Generating Functions

This chapter looks at Probability Generating Functions (PGFs) for **discrete** random variables. PGFs are useful tools for dealing with sums and limits of random variables. For some stochastic processes, they also have a special role in telling us whether a process will *ever* reach a particular state.

By the end of this chapter, you should be able to:

- find the sum of Geometric, Binomial, and Exponential series;
 - know the definition of the PGF, and use it to calculate the mean, variance, and probabilities;
 - calculate the PGF for Geometric, Binomial, and Poisson distributions;
 - calculate the PGF for a randomly stopped sum;
 - calculate the PGF for first reaching times in the random walk;
 - use the PGF to determine whether a process will *ever* reach a given state.
-

4.1 Common sums

1. Geometric Series

$$1 + r + r^2 + r^3 + \dots = \sum_{x=0}^{\infty} r^x = \frac{1}{1-r}, \quad \text{when } |r| < 1.$$

This formula proves that $\sum_{x=0}^{\infty} \mathbb{P}(X = x) = 1$ when $X \sim \text{Geometric}(p)$:

$$\begin{aligned} \mathbb{P}(X = x) = p(1-p)^x &\Rightarrow \sum_{x=0}^{\infty} \mathbb{P}(X = x) = \sum_{x=0}^{\infty} p(1-p)^x \\ &= p \sum_{x=0}^{\infty} (1-p)^x \\ &= \frac{p}{1-(1-p)} \quad (\text{because } |1-p| < 1) \\ &= 1. \end{aligned}$$

2. Binomial Theorem For any $p, q \in \mathbb{R}$, and integer n ,

$$(p + q)^n = \sum_{x=0}^n \binom{n}{x} p^x q^{n-x}.$$

Note that $\binom{n}{x} = \frac{n!}{(n-x)!x!}$ (nC_r button on calculator.)

The Binomial Theorem proves that $\sum_{x=0}^n \mathbb{P}(X = x) = 1$ when $X \sim \text{Binomial}(n, p)$:
 $\mathbb{P}(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$ for $x = 0, 1, \dots, n$, so

$$\begin{aligned} \sum_{x=0}^n \mathbb{P}(X = x) &= \sum_{x=0}^n \binom{n}{x} p^x (1-p)^{n-x} \\ &= (p + (1-p))^n \\ &= 1^n \\ &= 1. \end{aligned}$$

3. Exponential Power Series

$$\text{For any } \lambda \in \mathbb{R}, \quad \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^\lambda.$$

This proves that $\sum_{x=0}^{\infty} \mathbb{P}(X = x) = 1$ when $X \sim \text{Poisson}(\lambda)$:

$$\mathbb{P}(X = x) = \frac{\lambda^x}{x!} e^{-\lambda} \text{ for } x = 0, 1, 2, \dots, \text{ so}$$

$$\begin{aligned} \sum_{x=0}^{\infty} \mathbb{P}(X = x) &= \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} e^{-\lambda} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} \\ &= e^{-\lambda} e^\lambda \\ &= 1. \end{aligned}$$

Note: Another useful identity is: $e^\lambda = \lim_{n \rightarrow \infty} \left(1 + \frac{\lambda}{n}\right)^n$ for $\lambda \in \mathbb{R}$.

4.2 Probability Generating Functions

The **probability generating function (PGF)** is a useful tool for dealing with **discrete** random variables taking values $0, 1, 2, \dots$. Its particular strength is that it gives us an easy way of characterizing the distribution of $X + Y$ when X and Y are independent. In general it is difficult to find the distribution of a sum using the traditional probability function. The PGF transforms a sum into a product and enables it to be handled much more easily.

Sums of random variables are particularly important in the study of stochastic processes, because many stochastic processes are formed from the sum of a sequence of repeating steps: for example, the Gambler's Ruin from Section 2.7.

The name *probability generating function* also gives us another clue to the role of the PGF. The PGF can be used to generate all the probabilities of the distribution. This is generally tedious and is not often an efficient way of calculating probabilities. However, the fact that it *can* be done demonstrates that *the PGF tells us everything there is to know about the distribution*.

Definition: Let X be a discrete random variable taking values in the non-negative integers $\{0, 1, 2, \dots\}$. The **probability generating function (PGF)** of X is $G_X(s) = \mathbb{E}(s^X)$, for all $s \in \mathbb{R}$ for which the sum converges.

Calculating the probability generating function

$$G_X(s) = \mathbb{E}(s^X) = \sum_{x=0}^{\infty} s^x \mathbb{P}(X = x).$$

Properties of the PGF:

1. $G_X(0) = \mathbb{P}(X = 0)$:

$$\begin{aligned} G_X(0) &= 0^0 \times \mathbb{P}(X = 0) + 0^1 \times \mathbb{P}(X = 1) + 0^2 \times \mathbb{P}(X = 2) + \dots \\ \therefore G_X(0) &= \mathbb{P}(X = 0). \end{aligned}$$

$$\underline{2. G_X(1) = 1 :} \quad G_X(1) = \sum_{x=0}^{\infty} 1^x \mathbb{P}(X = x) = \sum_{x=0}^{\infty} \mathbb{P}(X = x) = 1.$$

Example 1: Binomial Distribution

Let $X \sim \text{Binomial}(n, p)$, so $\mathbb{P}(X = x) = \binom{n}{x} p^x q^{n-x}$ for $x = 0, 1, \dots, n$.

$$\begin{aligned} G_X(s) &= \sum_{x=0}^n s^x \binom{n}{x} p^x q^{n-x} \\ &= \sum_{x=0}^n \binom{n}{x} (ps)^x q^{n-x} \\ &= (ps + q)^n \quad \text{by the Binomial Theorem: true for all } s. \end{aligned}$$

Thus $G_X(s) = (ps + q)^n$ for all $s \in \mathbb{R}$.

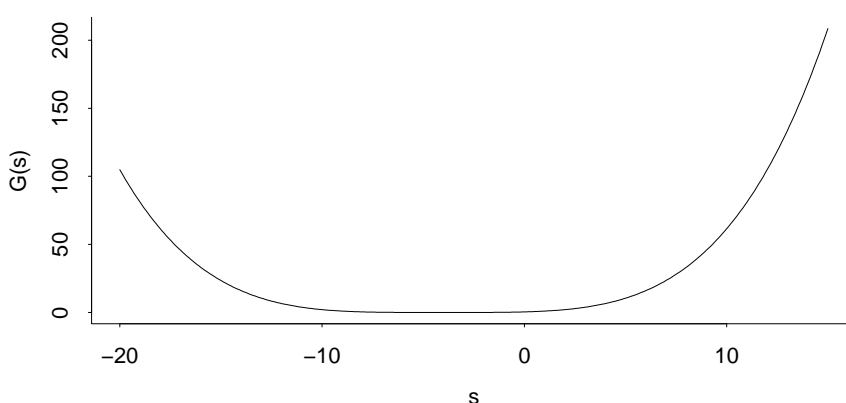
$X \sim \text{Bin}(n=4, p=0.2)$

Check $G_X(0)$:

$$\begin{aligned} G_X(0) &= (p \times 0 + q)^n \\ &= q^n \\ &= \mathbb{P}(X = 0). \end{aligned}$$

Check $G_X(1)$:

$$\begin{aligned} G_X(1) &= (p \times 1 + q)^n \\ &= (1)^n \\ &= 1. \end{aligned}$$



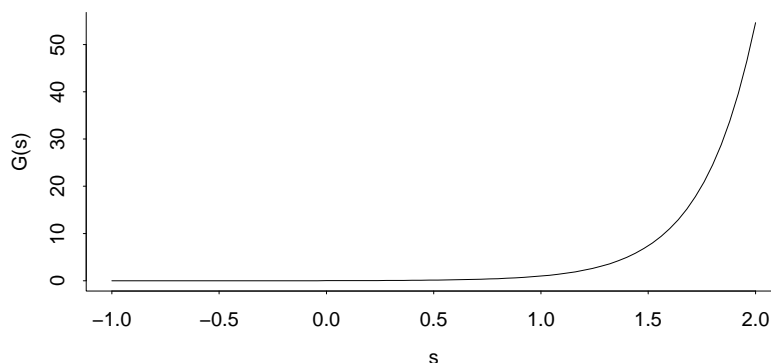
Example 2: Poisson Distribution

Let $X \sim \text{Poisson}(\lambda)$, so $\mathbb{P}(X = x) = \frac{\lambda^x}{x!} e^{-\lambda}$ for $x = 0, 1, 2, \dots$

$$\begin{aligned} G_X(s) &= \sum_{x=0}^{\infty} s^x \frac{\lambda^x}{x!} e^{-\lambda} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{(\lambda s)^x}{x!} \\ &= e^{-\lambda} e^{(\lambda s)} \quad \text{for all } s \in \mathbb{R}. \end{aligned}$$

Thus $G_X(s) = e^{\lambda(s-1)}$ for all $s \in \mathbb{R}$.

$X \sim \text{Poisson}(4)$



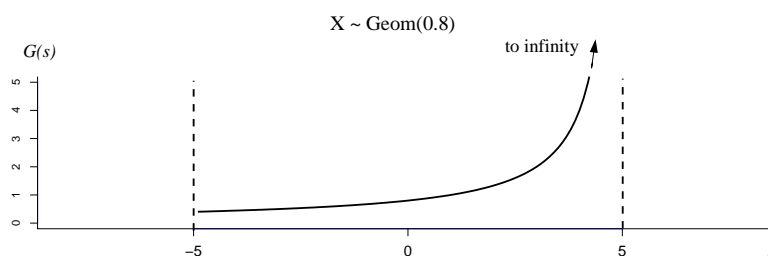
Example 3: Geometric Distribution

Let $X \sim \text{Geometric}(p)$, so $\mathbb{P}(X = x) = p(1 - p)^x = pq^x$ for $x = 0, 1, 2, \dots$, where $q = 1 - p$.

$$\begin{aligned} G_X(s) &= \sum_{x=0}^{\infty} s^x pq^x \\ &= p \sum_{x=0}^{\infty} (qs)^x \end{aligned}$$

$$= \frac{p}{1 - qs} \quad \text{for all } s \text{ such that } |qs| < 1.$$

Thus $G_X(s) = \frac{p}{1 - qs}$ for $|s| < \frac{1}{q}$.



4.3 Using the probability generating function to calculate probabilities

The probability generating function gets its name because the power series can be expanded and differentiated to reveal the individual probabilities. Thus, *given only the PGF* $G_X(s) = \mathbb{E}(s^X)$, *we can recover all probabilities* $\mathbb{P}(X = x)$.

For shorthand, write $p_x = \mathbb{P}(X = x)$. Then

$$G_X(s) = \mathbb{E}(s^X) = \sum_{x=0}^{\infty} p_x s^x = p_0 + p_1 s + p_2 s^2 + p_3 s^3 + p_4 s^4 + \dots$$

Thus $p_0 = \mathbb{P}(X = 0) = G_X(0)$.

First derivative: $G'_X(s) = p_1 + 2p_2 s + 3p_3 s^2 + 4p_4 s^3 + \dots$

Thus $p_1 = \mathbb{P}(X = 1) = G'_X(0)$.

Second derivative: $G''_X(s) = 2p_2 + (3 \times 2)p_3 s + (4 \times 3)p_4 s^2 + \dots$

Thus $p_2 = \mathbb{P}(X = 2) = \frac{1}{2} G''_X(0)$.

Third derivative: $G'''_X(s) = (3 \times 2 \times 1)p_3 + (4 \times 3 \times 2)p_4 s + \dots$

Thus $p_3 = \mathbb{P}(X = 3) = \frac{1}{3!} G'''_X(0)$.

In general:

$$p_n = \mathbb{P}(X = n) = \left(\frac{1}{n!} \right) G_X^{(n)}(0) = \left(\frac{1}{n!} \right) \frac{d^n}{ds^n} (G_X(s)) \Big|_{s=0}.$$

Example: Let X be a discrete random variable with PGF $G_X(s) = \frac{s}{5}(2 + 3s^2)$. Find the distribution of X .

$$G_X(s) = \frac{2}{5}s + \frac{3}{5}s^3 : \quad G_X(0) = \mathbb{P}(X = 0) = 0.$$

$$G'_X(s) = \frac{2}{5} + \frac{9}{5}s^2 : \quad G'_X(0) = \mathbb{P}(X = 1) = \frac{2}{5}.$$

$$G''_X(s) = \frac{18}{5}s : \quad \frac{1}{2}G''_X(0) = \mathbb{P}(X = 2) = 0.$$

$$G'''_X(s) = \frac{18}{5} : \quad \frac{1}{3!}G'''_X(0) = \mathbb{P}(X = 3) = \frac{3}{5}.$$

$$G_X^{(r)}(s) = 0 \quad \forall r \geq 4 : \quad \frac{1}{r!}G_X^{(r)}(s) = \mathbb{P}(X = r) = 0 \quad \forall r \geq 4.$$

Thus

$$X = \begin{cases} 1 & \text{with probability } 2/5, \\ 3 & \text{with probability } 3/5. \end{cases}$$

Uniqueness of the PGF

The formula $p_n = \mathbb{P}(X = n) = \left(\frac{1}{n!}\right) G_X^{(n)}(0)$ shows that the whole sequence of probabilities p_0, p_1, p_2, \dots is determined by the values of the PGF and its derivatives at $s = 0$. It follows that the PGF specifies a **unique** set of probabilities.

Fact: If two power series agree on any interval containing 0, however small, then all terms of the two series are equal.

Formally: let $A(s)$ and $B(s)$ be PGFs with $A(s) = \sum_{n=0}^{\infty} a_n s^n$, $B(s) = \sum_{n=0}^{\infty} b_n s^n$. If there exists some $R' > 0$ such that $A(s) = B(s)$ for all $-R' < s < R'$, then $a_n = b_n$ for all n .

Practical use: If we can show that two random variables have the same PGF in some interval containing 0, then we have shown that **the two random variables have the same distribution**.

Another way of expressing this is to say that **the PGF of X tells us everything there is to know about the distribution of X** .

4.4 Expectation and moments from the PGF

As well as calculating probabilities, we can also use the PGF to calculate the moments of the distribution of X . The moments of a distribution are *the mean, variance, etc.*

Theorem 4.4: Let X be a discrete random variable with PGF $G_X(s)$. Then:

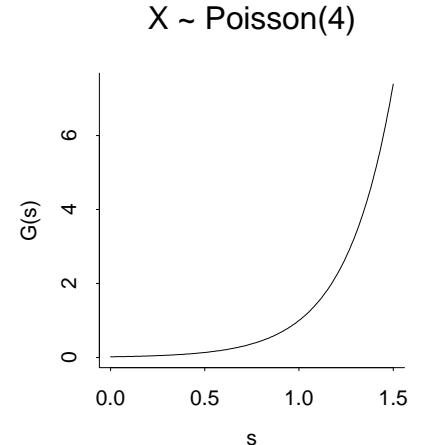
$$1. \mathbb{E}(X) = G'_X(1).$$

$$2. \mathbb{E}\left\{X(X-1)(X-2)\dots(X-k+1)\right\} = G_X^{(k)}(1) = \left.\frac{d^k G_X(s)}{ds^k}\right|_{s=1}.$$

(This is the k th factorial moment of X .)

Proof: (Sketch: see Section 4.8 for more details)

$$\begin{aligned} 1. \quad G_X(s) &= \sum_{x=0}^{\infty} s^x p_x, \\ \text{so} \quad G'_X(s) &= \sum_{x=0}^{\infty} x s^{x-1} p_x \\ \Rightarrow \quad G'_X(1) &= \sum_{x=0}^{\infty} x p_x = \mathbb{E}(X) \end{aligned}$$



$$\begin{aligned} 2. \quad G_X^{(k)}(s) &= \frac{d^k G_X(s)}{ds^k} = \sum_{x=k}^{\infty} x(x-1)(x-2)\dots(x-k+1)s^{x-k} p_x \\ \text{so} \quad G_X^{(k)}(1) &= \sum_{x=k}^{\infty} x(x-1)(x-2)\dots(x-k+1)p_x \\ &= \mathbb{E}\left\{X(X-1)(X-2)\dots(X-k+1)\right\}. \quad \square \end{aligned}$$

Example: Let $X \sim \text{Poisson}(\lambda)$. The PGF of X is $G_X(s) = e^{\lambda(s-1)}$. Find $\mathbb{E}(X)$ and $\text{Var}(X)$.

$X \sim \text{Poisson}(4)$

Solution:

$$G'_X(s) = \lambda e^{\lambda(s-1)}$$

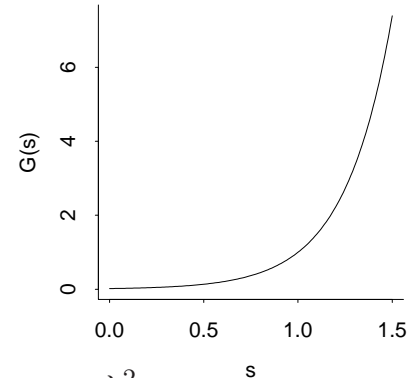
$$\Rightarrow \mathbb{E}(X) = G'_X(1) = \lambda.$$

For the variance, consider

$$\mathbb{E}\{X(X-1)\} = G''_X(1) = \lambda^2 e^{\lambda(s-1)}|_{s=1} = \lambda^2.$$

So

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}(X^2) - (\mathbb{E}X)^2 \\ &= \mathbb{E}\{X(X-1)\} + \mathbb{E}X - (\mathbb{E}X)^2 \\ &= \lambda^2 + \lambda - \lambda^2 \\ &= \lambda. \end{aligned}$$



4.5 Probability generating function for a sum of independent r.v.s

One of the PGF's greatest strengths is that it turns a sum into a product:

$$\mathbb{E}\left(s^{(X_1+X_2)}\right) = \mathbb{E}\left(s^{X_1}s^{X_2}\right).$$

This makes the PGF useful for finding the probabilities and moments of **a sum of independent random variables**.

Theorem 4.5: Suppose that X_1, \dots, X_n are ***independent*** random variables, and let $Y = X_1 + \dots + X_n$. Then

$$G_Y(s) = \prod_{i=1}^n G_{X_i}(s).$$

Proof:

$$\begin{aligned}
 G_Y(s) &= \mathbb{E}(s^{X_1+\dots+X_n}) \\
 &= \mathbb{E}(s^{X_1} s^{X_2} \dots s^{X_n}) \\
 &= \mathbb{E}(s^{X_1}) \mathbb{E}(s^{X_2}) \dots \mathbb{E}(s^{X_n}) \\
 &\quad \text{(because } X_1, \dots, X_n \text{ are independent)} \\
 &= \prod_{i=1}^n G_{X_i}(s). \quad \text{as required.} \quad \square
 \end{aligned}$$

Example: Suppose that X and Y are independent with $X \sim \text{Poisson}(\lambda)$ and $Y \sim \text{Poisson}(\mu)$. Find the distribution of $X + Y$.

Solution:

$$\begin{aligned}
 G_{X+Y}(s) &= G_X(s) \cdot G_Y(s) \\
 &= e^{\lambda(s-1)} e^{\mu(s-1)} \\
 &= e^{(\lambda+\mu)(s-1)}.
 \end{aligned}$$

But this is the PGF of the $\text{Poisson}(\lambda + \mu)$ distribution. So, by the uniqueness of PGFs, $X + Y \sim \text{Poisson}(\lambda + \mu)$.

4.6 Randomly stopped sum

Remember the randomly stopped sum model from Section 3.4. A random number N of events occur, and each event i has associated with it a cost or reward X_i . The question is to find the distribution of the total cost or reward: $T_N = X_1 + X_2 + \dots + X_N$.

T_N is called a *randomly stopped sum* because it has a random number of terms.



Example: Cash machine model. N customers arrive during the day. Customer i withdraws amount X_i . The total amount withdrawn during the day is $T_N = X_1 + \dots + X_N$.

In Chapter 3, we used the Laws of Total Expectation and Variance to show that $\mathbb{E}(T_N) = \mu \mathbb{E}(N)$ and $\text{Var}(T_N) = \sigma^2 \mathbb{E}(N) + \mu^2 \text{Var}(N)$, where $\mu = \mathbb{E}(X_i)$ and $\sigma^2 = \text{Var}(X_i)$.

In this chapter we will now use probability generating functions to investigate the *whole distribution of* T_N .

Theorem 4.6: Let X_1, X_2, \dots be a sequence of independent and identically distributed random variables with common PGF G_X . Let N be a random variable, independent of the X_i 's, with PGF G_N , and let $T_N = X_1 + \dots + X_N = \sum_{i=1}^N X_i$. Then the PGF of T_N is:

$$G_{T_N}(s) = G_N(G_X(s)).$$

Proof:

$$\begin{aligned} G_{T_N}(s) &= \mathbb{E}(s^{T_N}) = \mathbb{E}(s^{X_1 + \dots + X_N}) \\ &= \mathbb{E}_N \left\{ \mathbb{E}(s^{X_1 + \dots + X_N} \mid N) \right\} \quad (\text{conditional expectation}) \\ &= \mathbb{E}_N \left\{ \mathbb{E}(s^{X_1} \dots s^{X_N} \mid N) \right\} \\ &= \mathbb{E}_N \left\{ \mathbb{E}(s^{X_1} \dots s^{X_N}) \right\} \quad (X_i \text{'s are indept of } N) \\ &= \mathbb{E}_N \left\{ \mathbb{E}(s^{X_1}) \dots \mathbb{E}(s^{X_N}) \right\} \quad (X_i \text{'s are indept of each other}) \\ &= \mathbb{E}_N \left\{ (G_X(s))^N \right\} \\ &= G_N(G_X(s)) \quad (\text{by definition of } G_N). \quad \square \end{aligned}$$

Example: Let X_1, X_2, \dots and N be as above. Find the mean of T_N .

$$\begin{aligned}
 \mathbb{E}(T_N) = G'_{T_N}(1) &= \left. \frac{d}{ds} G_N(G_X(s)) \right|_{s=1} \\
 &= \left. G'_N(G_X(s)) \cdot G'_X(s) \right|_{s=1} \\
 &= G'_N(1) \cdot G'_X(1) \quad \text{Note: } G_X(1) = 1 \text{ for any r.v. } X \\
 &= \mathbb{E}(N) \cdot \mathbb{E}(X_1), \quad \text{— same answer as in Chapter 3.}
 \end{aligned}$$

Example: Heron goes fishing

My aunt was asked by her neighbours to feed the prize goldfish in their garden pond while they were on holiday. Although my aunt dutifully went and fed them every day, she never saw a single fish for the whole three weeks. It turned out that all the fish had been eaten by a heron when she wasn't looking!



Let N be the number of times the heron visits the pond during the neighbours' absence. Suppose that $N \sim \text{Geometric}(1 - \theta)$, so $\mathbb{P}(N = n) = (1 - \theta)\theta^n$, for $n = 0, 1, 2, \dots$. When the heron visits the pond it has probability p of catching a prize goldfish, independently of what happens on any other visit. (This assumes that there are infinitely many goldfish to be caught!) Find the distribution of

T = total number of goldfish caught.

Solution:

$$\text{Let } X_i = \begin{cases} 1 & \text{if heron catches a fish on visit } i, \\ 0 & \text{otherwise.} \end{cases}$$

Then $T = X_1 + X_2 + \dots + X_N$ (randomly stopped sum), so

$$G_T(s) = G_N(G_X(s)).$$

Now

$$G_X(s) = \mathbb{E}(s^X) = s^0 \times \mathbb{P}(X = 0) + s^1 \times \mathbb{P}(X = 1) = 1 - p + ps.$$

Also,

$$\begin{aligned} G_N(r) &= \sum_{n=0}^{\infty} r^n \mathbb{P}(N = n) = \sum_{n=0}^{\infty} r^n (1 - \theta) \theta^n \\ &= (1 - \theta) \sum_{n=0}^{\infty} (\theta r)^n \\ &= \frac{1 - \theta}{1 - \theta r}. \quad (r < 1/\theta). \end{aligned}$$

So

$$G_T(s) = \frac{1 - \theta}{1 - \theta G_X(s)} \quad (\text{putting } r = G_X(s)),$$

giving:

$$\begin{aligned} G_T(s) &= \frac{1 - \theta}{1 - \theta(1 - p + ps)} \\ &= \frac{1 - \theta}{1 - \theta + \theta p - \theta ps} \end{aligned}$$

$$[\text{could this be Geometric? } G_T(s) = \frac{1 - \pi}{1 - \pi s} \text{ for some } \pi?]$$

$$\begin{aligned} &= \frac{1 - \theta}{(1 - \theta + \theta p) - \theta ps} \\ &= \frac{\left(\frac{1 - \theta}{1 - \theta + \theta p} \right)}{\left(\frac{(1 - \theta + \theta p) - \theta ps}{1 - \theta + \theta p} \right)} \end{aligned}$$

$$\begin{aligned}
 &= \frac{\left(\frac{1 - \theta + \theta p - \theta p}{1 - \theta + \theta p}\right)}{1 - \left(\frac{\theta p}{1 - \theta + \theta p}\right)^s} \\
 &= \frac{1 - \left(\frac{\theta p}{1 - \theta + \theta p}\right)}{1 - \left(\frac{\theta p}{1 - \theta + \theta p}\right)^s}.
 \end{aligned}$$

This is the PGF of the Geometric $\left(1 - \frac{\theta p}{1 - \theta + \theta p}\right)$ distribution, so by uniqueness of PGFs, we have:

$$T \sim \text{Geometric}\left(\frac{1 - \theta}{1 - \theta + \theta p}\right).$$

Why did we need to use the PGF?

We could have solved the heron problem without using the PGF, but it is much more difficult. PGFs are very useful for dealing with sums of random variables, which are difficult to tackle using the standard probability function.

Here are the first few steps of solving the heron problem without the PGF. Recall the problem:

- Let $N \sim \text{Geometric}(1 - \theta)$, so $\mathbb{P}(N = n) = (1 - \theta)\theta^n$;
- Let X_1, X_2, \dots be independent of each other and of N , with $X_i \sim \text{Binomial}(1, p)$ (remember $X_i = 1$ with probability p , and 0 otherwise);
- Let $T = X_1 + \dots + X_N$ be the randomly stopped sum;
- Find the distribution of T .

Without using the PGF, we would tackle this by looking for an expression for $\mathbb{P}(T = t)$ for any t . Once we have obtained that expression, we might be able to see that T has a distribution we recognise (e.g. Geometric), or otherwise we would just state that T is defined by the probability function we have obtained.

To find $\mathbb{P}(T = t)$, we have to *partition over different values of N* :

$$\mathbb{P}(T = t) = \sum_{n=0}^{\infty} \mathbb{P}(T = t \mid N = n) \mathbb{P}(N = n). \quad (\star)$$

Here, we are *lucky* that we can write down the distribution of $T \mid N = n$:

- if $N = n$ is fixed, then $T = X_1 + \dots + X_n$ is a sum of n independent $\text{Binomial}(1, p)$ random variables, so $(T \mid N = n) \sim \text{Binomial}(n, p)$.

For most distributions of X , *it would be difficult or impossible to write down the distribution of $X_1 + \dots + X_n$* :

we would have to use an expression like

$$\begin{aligned} \mathbb{P}(X_1 + \dots + X_N = t \mid N = n) &= \sum_{x_1=0}^t \sum_{x_2=0}^{t-x_1} \dots \sum_{x_{n-1}=0}^{t-(x_1+\dots+x_{n-2})} \left\{ \mathbb{P}(X_1 = x_1) \times \right. \\ &\quad \left. \mathbb{P}(X_2 = x_2) \times \dots \times \mathbb{P}(X_{n-1} = x_{n-1}) \times \mathbb{P}[X_n = t - (x_1 + \dots + x_{n-1})] \right\}. \end{aligned}$$

Back to the heron problem: we are lucky in this case that we know the distribution of $(T \mid N = n)$ is $\text{Binomial}(N = n, p)$, so

$$\mathbb{P}(T = t \mid N = n) = \binom{n}{t} p^t (1-p)^{n-t} \quad \text{for } t = 0, 1, \dots, n.$$

Continuing from (\star) :

$$\mathbb{P}(T = t) = \sum_{n=0}^{\infty} \mathbb{P}(T = t \mid N = n) \mathbb{P}(N = n)$$

$$\begin{aligned}
 &= \sum_{n=t}^{\infty} \binom{n}{t} p^t (1-p)^{n-t} (1-\theta) \theta^n \\
 &= (1-\theta) \left(\frac{p}{1-p} \right)^t \sum_{n=t}^{\infty} \binom{n}{t} [\theta(1-p)]^n \quad (\star\star) \\
 &= \dots?
 \end{aligned}$$

As it happens, we can evaluate the sum in $(\star\star)$ using the fact that Negative Binomial probabilities sum to 1. You can try this if you like, but it is quite tricky. [Hint: use the Negative Binomial $(t+1, 1-\theta(1-p))$ distribution.]

Overall, we obtain the same answer that $T \sim \text{Geometric} \left(\frac{1-\theta}{1-\theta+\theta p} \right)$, but hopefully you can see why the PGF is so useful.

Without the PGF, we have two major difficulties:

1. *Writing down* $\mathbb{P}(T = t \mid N = n)$;
2. *Evaluating the sum over* n *in* $(\star\star)$.

For a general problem, both of these steps might be too difficult to do without a computer. The PGF has none of these difficulties, and even if $G_T(s)$ does not simplify readily, it still tells us everything there is to know about the distribution of T .

4.7 Summary: Properties of the PGF

Definition:	$G_X(s) = \mathbb{E}(s^X)$	
Used for:	Discrete r.v.s with values $0, 1, 2, \dots$	
Moments:	$\mathbb{E}(X) = G'_X(1)$	$\mathbb{E}\{X(X-1)\dots(X-k+1)\} = G_X^{(k)}(1)$
Probabilities:	$\mathbb{P}(X = n) = \frac{1}{n!} G_X^{(n)}(0)$	
Sums:	$G_{X+Y}(s) = G_X(s)G_Y(s)$ for independent X, Y	

4.8 Convergence of PGFs

We have been using PGFs throughout this chapter without paying much attention to their mathematical properties. For example, are we sure that the power series $G_X(s) = \sum_{x=0}^{\infty} s^x \mathbb{P}(X = x)$ converges? Can we differentiate and integrate the infinite power series term by term as we did in Section 4.4? When we said in Section 4.4 that $\mathbb{E}(X) = G'_X(1)$, can we be sure that $G_X(1)$ and its derivative $G'_X(1)$ even exist?

This technical section introduces the **radius of convergence** of the PGF. Although it isn't obvious, it is always safe to assume convergence of $G_X(s)$ at least for $|s| < 1$. Also, there are results that assure us that $\mathbb{E}(X) = G'_X(1)$ will work for all non-defective random variables X .

Definition: The **radius of convergence** of a probability generating function is a number $R > 0$, such that the sum $G_X(s) = \sum_{x=0}^{\infty} s^x \mathbb{P}(X = x)$ converges if $|s| < R$ and diverges ($\rightarrow \infty$) if $|s| > R$.

(No general statement is made about what happens when $|s| = R$.)

Fact: For any PGF, the radius of convergence exists.

It is always ≥ 1 : every PGF converges for at least $s \in (-1, 1)$.

The radius of convergence could be anything from $R = 1$ to $R = \infty$.

Note: This gives us the surprising result that the set of s for which the PGF $G_X(s)$ converges is symmetric about 0: the PGF converges for all $s \in (-R, R)$, and for no $s < -R$ or $s > R$.

This is surprising because the PGF itself is not usually symmetric about 0: i.e. $G_X(-s) \neq G_X(s)$ in general.

Example 1: Geometric distribution

Let $X \sim \text{Geometric}(p = 0.8)$. What is the radius of convergence of $G_X(s)$?

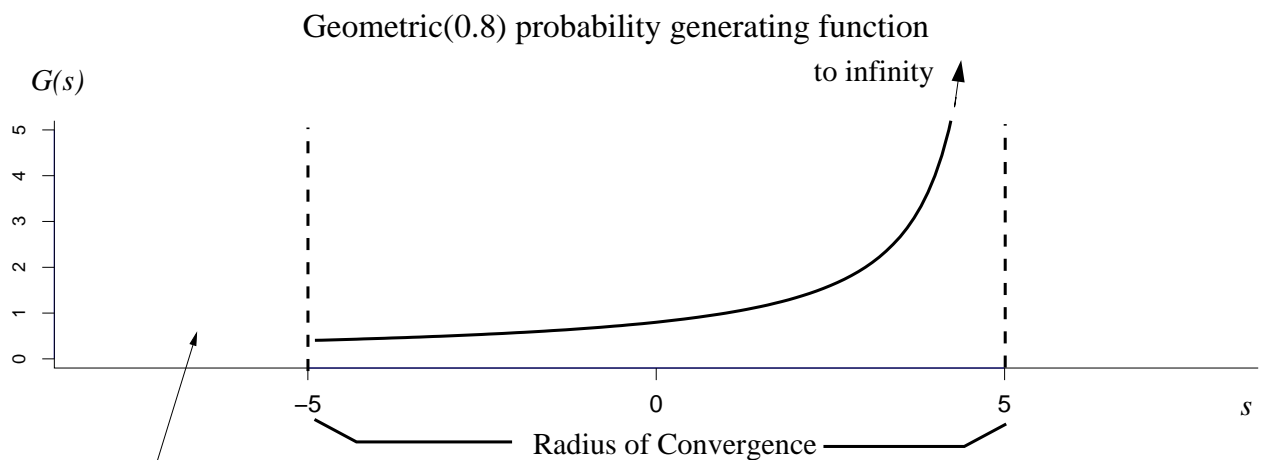
As in Section 4.2,

$$\begin{aligned} G_X(s) &= \sum_{x=0}^{\infty} s^x (0.8)(0.2)^x = 0.8 \sum_{x=0}^{\infty} (0.2s)^x \\ &= \frac{0.8}{1 - 0.2s} \quad \text{for all } s \text{ such that } |0.2s| < 1. \end{aligned}$$

This is valid for all s with $|0.2s| < 1$, so it is valid for all s with $|s| < \frac{1}{0.2} = 5$.
(i.e. $-5 < s < 5$.)

The radius of convergence is $R = 5$.

The figure shows the PGF of the Geometric($p = 0.8$) distribution, with its radius of convergence $R = 5$. Note that although the convergence set $(-5, 5)$ is symmetric about 0, the function $G_X(s) = p/(1 - qs) = 4/(5 - s)$ is not.



In this region, $p/(1 - qs)$ remains finite and well-behaved, but it is no longer equal to $E(s^X)$.

At the limits of convergence, strange things happen:

- At the positive end, as $s \uparrow 5$, both $G_X(s)$ and $p/(1 - qs)$ approach infinity. So the PGF is (left)-continuous at $+R$:

$$\lim_{s \uparrow 5} G_X(s) = G_X(5) = \infty.$$

However, the PGF does *not* converge at $s = +R$.

- At the negative end, as $s \downarrow -5$, the function $p/(1 - qs) = 4/(5 - s)$ is continuous and passes through 0.4 when $s = -5$. However, when $s \leq -5$, this function no longer represents $G_X(s) = 0.8 \sum_{x=0}^{\infty} (0.2s)^x$, because $|0.2s| \geq 1$.

Additionally, when $s = -5$, $G_X(-5) = 0.8 \sum_{x=0}^{\infty} (-1)^x$ does not exist. Unlike the positive end, this means that $G_X(s)$ is *not* (right)-continuous at $-R$:

$$\lim_{s \downarrow -5} G_X(s) = 0.4 \neq G_X(-5).$$

Like the positive end, this PGF does *not* converge at $s = -R$.

Example 2: Binomial distribution

Let $X \sim \text{Binomial}(n, p)$. What is the radius of convergence of $G_X(s)$?

As in Section 4.2,

$$\begin{aligned} G_X(s) &= \sum_{x=0}^n s^x \binom{n}{x} p^x q^{n-x} \\ &= \sum_{x=0}^n \binom{n}{x} (ps)^x q^{n-x} \\ &= (ps + q)^n \quad \text{by the Binomial Theorem: true for all } s. \end{aligned}$$

This is true for all $-\infty < s < \infty$, so the radius of convergence is $R = \infty$.

Abel's Theorem for continuity of power series at $s = 1$

Recall from above that if $X \sim \text{Geometric}(0.8)$, then $G_X(s)$ is not continuous at the negative end of its convergence ($-R$):

$$\lim_{s \downarrow -5} G_X(s) \neq G_X(-5).$$

Abel's theorem states that this sort of effect can never happen at $s = 1$ (or at $+R$). In particular, $G_X(s)$ is always left-continuous at $s = 1$:

$$\lim_{s \uparrow 1} G_X(s) = G_X(1) \quad \text{always, even if } G_X(1) = \infty.$$

Theorem 4.8: Abel's Theorem.

Let $G(s) = \sum_{i=0}^{\infty} p_i s^i$ for any p_0, p_1, p_2, \dots with $p_i \geq 0$ for all i .

Then $G(s)$ is left-continuous at $s = 1$:

$$\lim_{s \uparrow 1} G(s) = \sum_{i=0}^{\infty} p_i = G(1),$$

whether or not this sum is finite.

Note: Remember that the radius of convergence $R \geq 1$ for any PGF, so Abel's Theorem means that even in the worst-case scenario when $R = 1$, we can still trust that the PGF will be continuous at $s = 1$. (By contrast, we can not be sure that the PGF will be continuous at the lower limit $-R$).

Abel's Theorem means that for any PGF, we can write $G_X(1)$ as shorthand for $\lim_{s \uparrow 1} G_X(s)$.

It also clarifies our proof that $\mathbb{E}(X) = G'_X(1)$ from Section 4.4. If we assume that term-by-term differentiation is allowed for $G_X(s)$ (see below), then the proof on page 81 gives:

$$\begin{aligned} G_X(s) &= \sum_{x=0}^{\infty} s^x p_x, \\ \text{so } G'_X(s) &= \sum_{x=1}^{\infty} x s^{x-1} p_x \quad (\text{term-by-term differentiation: see below}). \end{aligned}$$

Abel's Theorem establishes that $\mathbb{E}(X)$ is equal to $\lim_{s \uparrow 1} G'_X(s)$:

$$\begin{aligned} \mathbb{E}(X) &= \sum_{x=1}^{\infty} x p_x \\ &= G'_X(1) \\ &= \lim_{s \uparrow 1} G'_X(s), \end{aligned}$$

because Abel's Theorem applies to $G'_X(s) = \sum_{x=1}^{\infty} x s^{x-1} p_x$, establishing that $G'_X(s)$ is left-continuous at $s = 1$. Without Abel's Theorem, we could not be sure that the limit of $G'_X(s)$ as $s \uparrow 1$ would give us the correct answer for $\mathbb{E}(X)$.

Absolute and uniform convergence for term-by-term differentiation

We have stated that the PGF converges for all $|s| < R$ for some R . In fact, the probability generating function converges *absolutely* if $|s| < R$. Absolute convergence is stronger than convergence alone: it means that the sum of absolute values, $\sum_{x=0}^{\infty} |s^x \mathbb{P}(X = x)|$, also converges. When two series both converge absolutely, the product series also converges absolutely. This guarantees that $G_X(s) \times G_Y(s)$ is absolutely convergent for any two random variables X and Y . This is useful because $G_X(s) \times G_Y(s) = G_{X+Y}(s)$ if X and Y are independent.

The PGF also converges *uniformly* on any set $\{s : |s| \leq R'\}$ where $R' < R$. Intuitively, this means that the speed of convergence does not depend upon the value of s . Thus a value n_0 can be found such that for all values of $n \geq n_0$, the *finite* sum $\sum_{x=0}^n s^x \mathbb{P}(X = x)$ is *simultaneously* close to the converged value $G_X(s)$, for all s with $|s| \leq R'$. In mathematical notation: $\forall \epsilon > 0, \exists n_0 \in \mathbb{Z}$ such that $\forall s$ with $|s| \leq R'$, and $\forall n \geq n_0$,

$$\left| \sum_{x=0}^n s^x \mathbb{P}(X = x) - G_X(s) \right| < \epsilon.$$

Uniform convergence allows us to differentiate or integrate the PGF term by term.

Fact: Let $G_X(s) = \mathbb{E}(s^X) = \sum_{x=0}^{\infty} s^x \mathbb{P}(X = x)$, and let $s < R$.

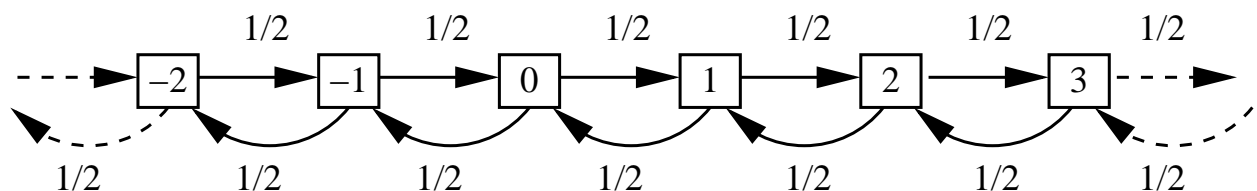
$$1. \quad G'_X(s) = \frac{d}{ds} \left(\sum_{x=0}^{\infty} s^x \mathbb{P}(X = x) \right) = \sum_{x=0}^{\infty} \frac{d}{ds} (s^x \mathbb{P}(X = x)) = \sum_{x=0}^{\infty} x s^{x-1} \mathbb{P}(X = x). \\ \text{(term by term differentiation).}$$

$$2. \quad \int_a^b G_X(s) ds = \int_a^b \left(\sum_{x=0}^{\infty} s^x \mathbb{P}(X = x) \right) ds = \sum_{x=0}^{\infty} \left(\int_a^b s^x \mathbb{P}(X = x) ds \right) \\ = \sum_{x=0}^{\infty} \frac{s^{x+1}}{x+1} \mathbb{P}(X = x) \quad \text{for } -R < a < b < R. \\ \text{(term by term integration).}$$

4.9 Special Process: the Random Walk

We briefly saw the Drunkard's Walk in Chapter 1: a drunk person staggers to left and right as he walks. This process is called the **Random Walk** in stochastic processes. Probability generating functions are particularly useful for processes such as the random walk, because the process is defined as the sum of a single repeating step. The repeating step is a move of one unit, left or right at random. The sum of the first t steps gives the position at time t .

The transition diagram below shows the *symmetric random walk* (all transitions have probability $p = 1/2$.)



Question:

What is the key difference between the random walk and the gambler's ruin?

The random walk has an INFINITE state space: it never stops. The gambler's ruin stops at both ends.

This fact has two important consequences:

- The random walk is hard to tackle using first-step analysis, because we would have to solve an *infinite* number of simultaneous equations. In this respect it might seem to be more difficult than the gambler's ruin.
- Because the random walk never stops, *all states are equal*.

In the gambler's ruin, states are not equal: the states closest to 0 are more likely to end in ruin than the states closest to winning. By contrast, the random walk has no end-points, so (for example) the distribution of the time to reach state 5 starting from state 0 is exactly the same as the distribution of the time to reach state 1005 starting from state 1000. We can exploit this fact to solve some problems for the random walk that would be much more difficult to solve for the gambler's ruin.

PGFs for finding the distribution of reaching times

For random walks, we are particularly interested in *reaching times*:

- How long will it take us to reach state j , starting from state i ?
- Is there a chance that we will **never** reach state j , starting from state i ?

In Chapter 3 we saw how to find *expected reaching times*: the expected number of steps taken to reach a particular state. We used *the law of total expectation and first-step analysis* (Section 3.5).

However, the *expected* or *average* reaching time doesn't tell the whole story. Think back to the model for gene spread in Section 3.7. If there is just one animal out of 100 with the harmful allele, the expected number of generations to fixation is quite large at 10.5: even though the allele will usually die out after one or two generations. The high average is caused by a small chance that the allele will take hold and grow, requiring a very large number of generations before it either dies out or saturates the population. In most stochastic processes, the average is of limited use by itself, without having some idea about the *variance and skew of the distribution*.

With our tool of PGFs, we can characterise the *whole distribution* of the time T taken to reach a particular state, by finding its PGF. This will give us the mean, variance, and skew by differentiation. In principle the PGF could even give us the full set of probabilities, $\mathbb{P}(T = t)$ for all possible $t = 0, 1, 2, \dots$, though in practice it may be computationally infeasible to find more than the first few probabilities by repeated differentiation.

However, there is a new and very useful piece of information that the PGF can tell us quickly and easily:

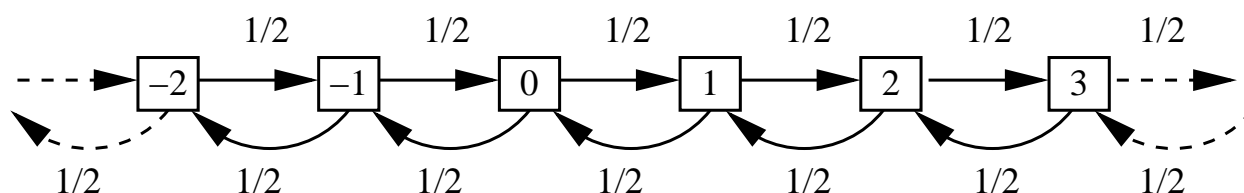
what is the probability that we NEVER reach state j , starting from state i ?

For example, imagine that the random walk represents the share value for an investment. The current share price is i dollars, and we might decide to sell when it reaches j dollars. Knowing how long this might take, and whether there is a chance we will never succeed, is fundamental to managing our investment.

To tackle this problem, we define the random variable T to be the time taken (number of steps) to reach state j , starting from state i . We find the PGF of T , and then use the PGF to discover $\mathbb{P}(T = \infty)$. If $\mathbb{P}(T = \infty) > 0$, there is a positive chance that we will NEVER reach state j , starting from state i .

We will see how to determine the probability of never reaching our goal in Section 4.11. First we will see how to calculate the PGF of a reaching time T in the random walk.

Finding the PGF of a reaching time in the random walk

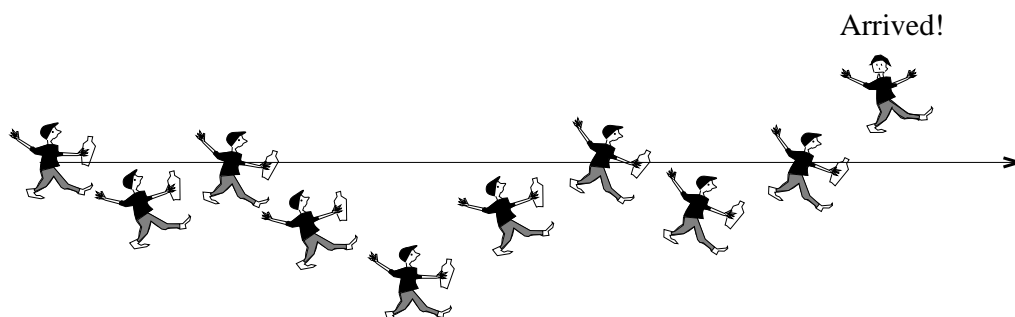


Define T_{ij} to be the *number of steps taken to reach state j , starting at state i* .

T_{ij} is called the *first reaching time from state i to state j* .

We will focus on T_{01} = *number of steps to get from state 0 to state 1*.

Problem: Let $H(s) = \mathbb{E}(s^{T_{01}})$ be the PGF of T_{01} . Find $H(s)$.



Solution:

Let Y_n be the step taken at time n : up or down. For the symmetric random walk,

$$Y_n = \begin{cases} 1 & \text{with probability } 0.5, \\ -1 & \text{with probability } 0.5, \end{cases}$$

and Y_1, Y_2, \dots are independent.

Recall T_{ij} = number of steps to get from state i to state j for any i, j ,

and $H(s) = \mathbb{E}(s^{T_{01}})$ is the PGF required.

Use first-step analysis, partitioning over the first step Y_1 :

$$\begin{aligned} H(s) &= \mathbb{E}(s^{T_{01}}) \\ &= \mathbb{E}(s^{T_{01}} | Y_1 = 1) \mathbb{P}(Y_1 = 1) + \mathbb{E}(s^{T_{01}} | Y_1 = -1) \mathbb{P}(Y_1 = -1) \\ &= \frac{1}{2} \left\{ \mathbb{E}(s^{T_{01}} | Y_1 = 1) + \mathbb{E}(s^{T_{01}} | Y_1 = -1) \right\}. \quad \spadesuit \end{aligned}$$

Now if $Y_1 = 1$, then $T_{01} = 1$ definitely, so $\mathbb{E}(s^{T_{01}} | Y_1 = 1) = s^1 = s$.

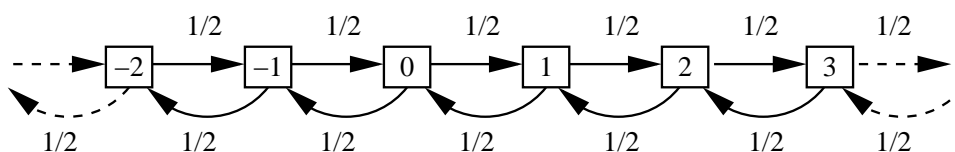
If $Y_1 = -1$, then $T_{01} = 1 + T_{-1,1}$:

→ one step from state 0 to state -1 ,

→ then $T_{-1,1}$ steps from state -1 to state 1.

But $T_{-1,1} = T_{-1,0} + T_{01}$, because the process must pass through 0 to get from -1 to 1.

Now $T_{-1,0}$ and T_{01} are independent (Markov property). Also, they have the same distribution because the process is translation invariant (i.e. all states are the same):



Thus

$$\begin{aligned}
 \mathbb{E}(s^{T_{01}} | Y_1 = -1) &= \mathbb{E}(s^{1+T_{-1,1}}) \\
 &= \mathbb{E}(s^{1+T_{-1,0}+T_{0,1}}) \\
 &= s\mathbb{E}(s^{T_{-1,0}}) \mathbb{E}(s^{T_{01}}) \quad \text{by independence} \\
 &= s(H(s))^2 \quad \text{because identically distributed.}
 \end{aligned}$$

Thus

$$H(s) = \frac{1}{2} \{s + s(H(s))^2\} \quad \text{by } \spadesuit.$$

This is a quadratic in $H(s)$:

$$\begin{aligned}
 \frac{1}{2}s(H(s))^2 - H(s) + \frac{1}{2}s &= 0 \\
 \Rightarrow H(s) &= \frac{1 \pm \sqrt{1 - 4\frac{1}{2}s\frac{1}{2}s}}{s} = \frac{1 \pm \sqrt{1 - s^2}}{s}.
 \end{aligned}$$

Which root? We know that $\mathbb{P}(T_{01} = 0) = 0$, because it must take at least one step to go from 0 to 1. With the positive root, $\lim_{s \rightarrow 0} H(s) = \lim_{s \rightarrow 0} \left(\frac{2}{s}\right) = \infty$; so we take the negative root instead.

$$\text{Thus } H(s) = \frac{1 - \sqrt{1 - s^2}}{s}.$$

Check this has $\lim_{s \rightarrow 0} H(s) = 0$ by L'Hospital's Rule:

$$\begin{aligned}
 \lim_{s \rightarrow 0} \left(\frac{f(s)}{g(s)} \right) &= \lim_{s \rightarrow 0} \left(\frac{f'(s)}{g'(s)} \right) \\
 &= \lim_{s \rightarrow 0} \left\{ \frac{\frac{1}{2}(1 - s^2)^{-1/2} \times 2s}{1} \right\} \\
 &= 0.
 \end{aligned}$$

Notation for quick solutions of first-step analysis for finding PGFs

As with first-step analysis for finding hitting probabilities and expected reaching times, setting up a good notation is extremely important. Here is a good notation for finding $H(s) = \mathbb{E}(s^{T_{01}})$.

Let $T = T_{01}$. Seek $H(s) = \mathbb{E}(s^T)$.

Now

$$T = \begin{cases} 1 & \text{with probability } 1/2, \\ 1 + T' + T'' & \text{with probability } 1/2, \end{cases}$$

where $T' \sim T'' \sim T$ and T', T'' are independent.

Taking expectations:

$$H(s) = \mathbb{E}(s^T) = \begin{cases} \mathbb{E}(s^1) & \text{w. p. } 1/2 \\ \mathbb{E}(s^{1+T'+T''}) & \text{w. p. } 1/2 \end{cases}$$

$$\Rightarrow H(s) = \begin{cases} s & \text{w. p. } 1/2 \\ s\mathbb{E}(s^{T'}) \mathbb{E}(s^{T''}) & \text{w. p. } 1/2 \end{cases} \quad (\text{by independence of } T' \text{ and } T'')$$

$$\Rightarrow H(s) = \begin{cases} s & \text{w. p. } 1/2 \\ sH(s)H(s) & \text{w. p. } 1/2 \end{cases} \quad (\text{because } T' \sim T'' \sim T)$$

$$\Rightarrow H(s) = \frac{1}{2}s + \frac{1}{2}sH(s)^2.$$

Thus:

$$sH(s)^2 - 2H(s) + s = 0.$$

Solve the quadratic and select the correct root as before, to get

$$H(s) = \frac{1 - \sqrt{1 - s^2}}{s} \quad \text{for } |s| < 1.$$

4.10 Defective random variables

A random variable is said to be *defective* if it can take the value ∞ .

In stochastic processes, a reaching time T_{ij} is defective if there is a chance that we *NEVER* reach state j , starting from state i .

The probability that we never reach state j , starting from state i , is the same as the probability that the time taken is infinite: $T_{ij} = \infty$:

$$\mathbb{P}(T_{ij} = \infty) = \mathbb{P}(\text{we NEVER reach state } j, \text{ starting from state } i).$$

In other cases, we will always reach state j eventually, starting from state i .

In that case, T_{ij} can not take the value ∞ :

$$\mathbb{P}(T_{ij} = \infty) = 0 \quad \text{if we are CERTAIN to reach state } j, \text{ starting from state } i.$$

Definition: A random variable T is defective, or improper, if it can take the value ∞ . That is,

$$T \text{ is defective if } \mathbb{P}(T = \infty) > 0.$$

Thinking of $\sum_{t=0}^{\infty} \mathbb{P}(T = t)$ as $1 - \mathbb{P}(T = \infty)$

Although it seems strange, when we write $\sum_{t=0}^{\infty} \mathbb{P}(T = t)$, *we are not including the value $t = \infty$.*

The sum $\sum_{t=0}^{\infty}$ continues without ever stopping: at no point can we say we have ‘finished’ all the finite values of t so we will now add on $t = \infty$. We simply *never get to $t = \infty$ when we take $\sum_{t=0}^{\infty}$.*

For a defective random variable T , this means that

$$\sum_{t=0}^{\infty} \mathbb{P}(T = t) < 1,$$

because we are missing the positive value of $\mathbb{P}(T = \infty)$.

All probabilities of T must still sum to 1, so we have

$$1 = \sum_{t=0}^{\infty} \mathbb{P}(T = t) + \mathbb{P}(T = \infty),$$

in other words

$$\sum_{t=0}^{\infty} \mathbb{P}(T = t) = 1 - \mathbb{P}(T = \infty).$$

PGFs for defective random variables

When T is defective, the PGF of T is *defined as* the power series

$$H(s) = \sum_{t=0}^{\infty} \mathbb{P}(T = t)s^t \quad \text{for } |s| < 1.$$

The term for $\mathbb{P}(T = \infty)s^{\infty}$ is missed out. The PGF is defined as the generating function of the probabilities for finite values only.

Because $H(s)$ is a power series satisfying the conditions of Abel's Theorem, we know that:

- $H(s)$ is left-continuous at $s = 1$, i.e. $\lim_{s \uparrow 1} H(s) = H(1)$.

This is different from the behaviour of $\mathbb{E}(s^T)$, if T is defective:

- $\mathbb{E}(s^T) = H(s)$ for $|s| < 1$ because the missing term is zero: i.e. because $s^\infty = 0$ when $|s| < 1$.
- $\mathbb{E}(s^T)$ is NOT left-continuous at $s = 1$. There is a sudden leap (discontinuity) at $s = 1$ because $s^\infty = 0$ as $s \uparrow 1$, but $s^\infty = 1$ when $s = 1$.

Thus $H(s)$ does NOT represent $\mathbb{E}(s^T)$ at $s = 1$. It is as if $H(s)$ is a 'train' that $\mathbb{E}(s^T)$ rides on between $-1 < s < 1$. At $s = 1$, the train keeps going (i.e. $H(s)$ is continuous) but $\mathbb{E}(s^T)$ jumps off the train.

We test whether T is defective by testing whether or not $\mathbb{E}(s^T)$ 'jumps off the train' — that is, we test whether or not $H(s)$ is equal to $\mathbb{E}(s^T)$ when $s = 1$.

We **know** what $\mathbb{E}(s^T)$ is when $s = 1$:

- $\mathbb{E}(s^T)$ is always 1 when $s = 1$, whether T is defective or not:

$$\mathbb{E}(1^T) = 1 \quad \text{for ANY random variable } T.$$

But the function $H(s) = \sum_{t=0}^{\infty} s^t \mathbb{P}(T = t)$ may or may not be 1 when $s = 1$:

- If T is defective, $H(s)$ is missing a term and $H(1) < 1$.
- If T is not defective, $H(s)$ is not missing anything so $H(1) = 1$.

Test for defectiveness:

Let $H(s) = \sum_{t=0}^{\infty} s^t \mathbb{P}(T = t)$ be the power series representing the PGF of T for $|s| < 1$. Then T is defective if and only if $H(1) < 1$.

Using defectiveness to find the probability we never get there

The simple test for defectiveness tells us whether there is a positive probability that we NEVER reach our goal. Here are the steps.

1. We want to know the probability that we will NEVER reach state j , starting from state i .
2. Define T to be the random variable giving the *number of steps taken* to get from state i to state j .
3. The event that we never reach state j , starting from state i , is the same as the event that $T = \infty$. (If we wait an infinite length of time, we never get there.) So

$$\mathbb{P}(\text{never reach state } j \mid \text{start at state } i) = \mathbb{P}(T = \infty).$$

4. Find $H(s) = \sum_{t=0}^{\infty} s^t \mathbb{P}(T = t)$, using a calculation like the one we did in Section 4.9. $H(s)$ is the PGF of T for $|s| < 1$. We only need to find it for $|s| < 1$. The calculation in Section 4.9 only works for $|s| \leq 1$ because the expectations are infinite or undefined when $|s| > 1$.
5. The random variable T is defective if and only if $H(1) < 1$.
6. If $H(1) < 1$, then the probability that T takes the value ∞ is the missing piece: $\mathbb{P}(T = \infty) = 1 - H(1)$.

Overall:

$$\mathbb{P}(\text{never reach state } j \mid \text{start at state } i) = \mathbb{P}(T = \infty) = 1 - H(1).$$

Expectation and variance of a defective random variable

If T is defective, there is a positive chance that $T = \infty$. This means that $\mathbb{E}(T) = \infty$, $\text{Var}(T) = \infty$, and $\mathbb{E}(T^a) = \infty$ for any power a .

$\mathbb{E}(T)$ and $\text{Var}(T)$ can not be found using the PGF when T is defective: you will get the wrong answer.

When you are asked to find $\mathbb{E}(T)$ in a context where T might be defective:

- First check whether T is defective: *is $H(1) < 1$ or $= 1$?*
- If T is defective, then $\mathbb{E}(T) = \infty$.
- If T is not defective ($H(1) = 1$), then $\mathbb{E}(T) = H'(1)$ *as usual*.

4.11 Random Walk: the probability we never reach our goal

In the random walk in Section 4.9, we defined the first reaching time T_{01} as the number of steps taken to get from state 0 to state 1.

In Section 4.9 we found the PGF of T_{01} to be:

$$\text{PGF of } T_{01} = H(s) = \frac{1 - \sqrt{1 - s^2}}{s} \text{ for } |s| < 1.$$

Questions:

- What is the probability that we *never* reach state 1, starting from state 0?
- What is expected number of steps to reach state 1, starting from state 0?

Solutions:

a) We need to know whether T_{01} is defective.

T_{01} is defective if and only if $H(1) < 1$.

Now $H(1) = \frac{1 - \sqrt{1 - 1^2}}{1} = 1$. So T_{01} is not defective.

Thus

$$\mathbb{P}(\text{never reach state 1} \mid \text{start from state 0}) = 0.$$

We will **DEFINITELY** reach state 1 eventually, even if it takes a very long time.

b) Because T_{01} is not defective, we can find $\mathbb{E}(T_{01})$ by differentiating the PGF: $\mathbb{E}(T_{01}) = H'(1)$.

$$H(s) = \frac{1 - \sqrt{1 - s^2}}{s} = s^{-1} - (s^{-2} - 1)^{1/2}$$

$$\text{So } H'(s) = -s^{-2} - \frac{1}{2}(s^{-2} - 1)^{-1/2}(-2s^{-3})$$

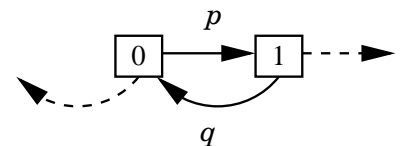
Thus

$$\mathbb{E}(T_{01}) = \lim_{s \uparrow 1} H'(s) = \lim_{s \uparrow 1} \left(-\frac{1}{s^2} + \frac{1}{s^3 \sqrt{\frac{1}{s^2} - 1}} \right) = \infty.$$

So the expected number of steps to reach state 1 starting from state 0 is infinite: $\mathbb{E}(T_{01}) = \infty$.

This result is striking. Even though we will **definitely** reach state 1, the expected time to do so is infinite! In general, we can prove the following results for random walks, starting from state 0:

Property	Reach state 1?	$\mathbb{P}(T_{01} = \infty)$	$\mathbb{E}(T_{01})$
$p > q$	Guaranteed	0	finite
$p = q = \frac{1}{2}$	Guaranteed	0	∞
$p < q$	Not guaranteed	> 0	∞



Note: (Non-examinable) If T is defective in the random walk, $\mathbb{E}(s^T)$ is not continuous at $s = 1$. In Section 4.9 we had to solve a quadratic equation to find $H(s) = \mathbb{E}(s^T)$. The negative root solution for $H(s)$ generally represents $\mathbb{E}(s^T)$ for $s < 1$. At $s = 1$, the solution for $\mathbb{E}(s^T)$ suddenly flips from the $-$ root to the $+$ root of the quadratic. This explains how $\mathbb{E}(s^T)$ can be discontinuous as $s \uparrow 1$, even though the negative root for $H(s)$ is continuous as $s \uparrow 1$ and all the working of Section 4.9 still applies for $s = 1$. The reason is that we suddenly switch from the $-$ root to the $+$ root at $s = 1$.

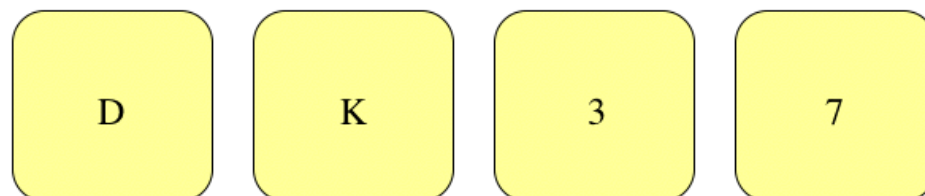
When $|s| > 1$, the conditional expectations are not finite so the working of Section 4.9 no longer applies.

The problem with being abstract . . .

- Each card has a letter on one side and a number on the other.
- We wish to test the following rule:

**If the card has a D on one side,
then it has a 3 on the other side.**

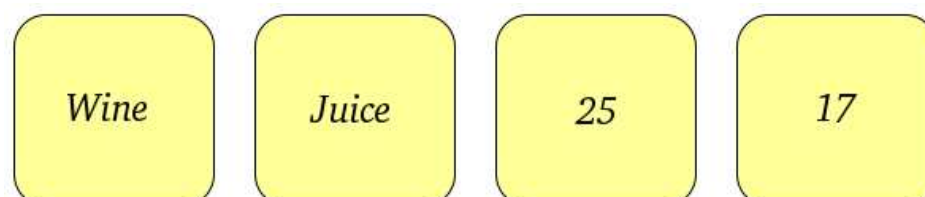
- Which card or cards should you turn over, and **ONLY** these cards, in order to test the rule?



At a party . . .

**If you are drinking alcohol,
you must be 18 or over.**

- Each card has the person's age on one side, and their drink on the other side.
- Which card or cards should you turn over, and **ONLY** these cards, in order to test the rule?



Chapter 5: Mathematical Induction

So far in this course, we have seen some techniques for dealing with stochastic processes: first-step analysis for hitting probabilities (Chapter 2), and first-step analysis for expected reaching times (Chapter 3). We now look at another tool that is often useful for exploring properties of stochastic processes: *proof by mathematical induction*.

5.1 Proving things in mathematics

There are many different ways of constructing a formal proof in mathematics. Some examples are:

- **Proof by counterexample:** a proposition is proved to be *not generally true* because a *particular example* is found for which it is not true.
- **Proof by contradiction:** this can be used either to prove a proposition is true or to prove that it is false. To prove that the proposition is *true* (say), we start by *assuming that it is false*. We then explore the consequences of this assumption until we reach a contradiction, e.g. $0 = 1$. Therefore something must have gone wrong, and the only thing we weren't sure about was our initial assumption that the proposition is false — so our initial assumption must be wrong and the proposition is proved true.

A famous proof of this sort is the proof that there are infinitely many prime numbers. We start by assuming that there are *finitely* many primes, so they can be listed as p_1, p_2, \dots, p_n , where p_n is the largest prime number. But then the number $p_1 \times p_2 \times \dots \times p_n + 1$ must also be prime, because it is not divisible by any of the smaller primes. Furthermore this number is definitely bigger than p_n . So we have contradicted the idea that there was a 'biggest' prime called p_n , and therefore there are infinitely many primes.

- **Proof by mathematical induction:** in mathematical induction, we start with a formula that we *suspect* is true. For example, I might *suspect* from

observation that $\sum_{k=1}^n k = n(n+1)/2$. I might have tested this formula for many different values of n , but of course I can never test it for *all* values of n . Therefore I need to prove that the formula is *always* true.

The idea of mathematical induction is to say: *suppose* the formula is true for all n up to the value $n = 10$ (say). Can I prove that, *if* it is true for $n = 10$, *then* it will also be true for $n = 11$? And *if* it is true for $n = 11$, then it will also be true for $n = 12$? And so on.

In practice, we usually start lower than $n = 10$. We usually take the very easiest case, $n = 1$, and prove that the formula is true for $n = 1$: $\text{LHS} = \sum_{k=1}^1 k = 1 = 1 \times 2/2 = \text{RHS}$. Then we prove that, *if* the formula is ever true for $n = x$, *then* it will always be true for $n = x + 1$. Because it is true for $n = 1$, it must be true for $n = 2$; and because it is true for $n = 2$, it must be true for $n = 3$; and so on, for all possible n . Thus the formula is proved.

Mathematical induction is therefore a bit like a *first-step analysis for proving things: prove that wherever we are now, the next step will always be OK. Then if we were OK at the very beginning, we will be OK for ever.*

The method of mathematical induction for proving results is very important in the study of Stochastic Processes. This is because a stochastic process builds up one step at a time, and mathematical induction works on the same principle.

Example: We have already seen examples of inductive-type reasoning in this course. For example, in Chapter 2 for the Gambler's Ruin problem, using the method of repeated substitution to solve for $p_x = \mathbb{P}(\text{Ruin} \mid \text{start with } \$x)$, we discovered that:

- $p_2 = 2p_1 - 1$
- $p_3 = 3p_1 - 2$
- $p_4 = 4p_1 - 3$

We deduced that $p_x = xp_1 - (x - 1)$ *in general*.

To prove this properly, we should have used the method of mathematical induction.

5.2 Mathematical Induction by example

This example explains the style and steps needed for a proof by induction.

Question: Prove by induction that $\sum_{k=1}^n k = \frac{n(n+1)}{2}$ for any integer n . (\star)

Approach: follow the steps below.

- (i) First verify that the formula is true for a *base case*: usually the smallest appropriate value of n (e.g. $n = 0$ or $n = 1$). Here, the smallest possible value of n is $n = 1$, because we can't have $\sum_{k=1}^0$.

First verify (\star) is true when $n = 1$.

$$LHS = \sum_{k=1}^1 k = 1.$$

$$RHS = \frac{1 \times 2}{2} = 1 = LHS.$$

So (\star) is proved for $n = 1$.

- (ii) Next suppose that formula (\star) is true for *all values of n up to and including some value x* . (We have already established that this is the case for $x = 1$).

Using the hypothesis that (\star) is true for all values of n up to and including x , prove that it is therefore true for the value $n = x + 1$.

Now suppose that (\star) is true for $n = 1, 2, \dots, x$ for some x .

Thus we can assume that $\sum_{k=1}^x k = \frac{x(x+1)}{2}$. (a)

((a) for 'allowed' info)

We need to show that if (\star) holds for $n = x$, then it must also hold for $n = x + 1$.

Require to prove that

$$\sum_{k=1}^{x+1} k = \frac{(x+1)(x+2)}{2} \quad (**)$$

(Obtained by putting $n = x + 1$ in $(*)$).

$$\begin{aligned} \text{LHS of } (**) &= \sum_{k=1}^{x+1} k = \sum_{k=1}^x k + (x+1) && \text{by expanding the sum} \\ &= \frac{x(x+1)}{2} + (x+1) && \text{using allowed info (a)} \\ &= (x+1) \left(\frac{x}{2} + 1 \right) && \text{rearranging} \\ &= \frac{(x+1)(x+2)}{2} \\ &= \text{RHS of } (**). \end{aligned}$$

This shows that:

$$\sum_{k=1}^n k = \frac{n(n+1)}{2} \quad \text{when } n = x + 1.$$

So, assuming $(*)$ is true for $n = x$, it is also true for $n = x + 1$.

- (iii) Refer back to the base case: if it is true for $n = 1$, then it is true for $n = 1+1 = 2$ by (ii). If it is true for $n = 2$, it is true for $n = 2 + 1 = 3$ by (ii). We could go on forever. This proves that the formula $(*)$ is true for all n .

We proved $(*)$ true for $n = 1$, thus $(*)$ is true for all integers $n \geq 1$. □

General procedure for proof by induction

The procedure above is quite standard. The inductive proof can be summarized like this:

Question: prove that $f(n) = g(n)$ for all integers $n \geq 1$. (\star)

Base case: $n = 1$. Prove that $f(1) = g(1)$ using

$$\begin{aligned} LHS &= f(1) \\ &= \vdots \\ &= g(1) = RHS. \end{aligned}$$

General case: suppose (\star) is true for $n = x$:

$$\text{so} \quad f(x) = g(x). \quad (a) \quad (\text{allowed info})$$

Prove that (\star) is therefore true for $n = x + 1$:

$$RTP \quad f(x + 1) = g(x + 1). \quad (\star\star)$$

$$\begin{aligned} LHS(\star\star) &= f(x + 1) \\ &= \left\{ \begin{array}{l} \text{some expression breaking down } f(x + 1) \\ \text{into } f(x) \text{ and an extra term in } x + 1 \end{array} \right\} \\ &= \left\{ \text{substitute } f(x) = g(x) \text{ in the line above} \right\} \quad \text{by allowed (a)} \\ &= \{ \text{do some working} \} \\ &= g(x + 1) \\ &= RHS(\star\star). \end{aligned}$$

Conclude: (\star) is proved for $n = 1$, so it is proved for $n = 2, n = 3, n = 4, \dots$

(\star) is therefore proved for all integers $n \geq 1$. □

5.3 Some harder examples of mathematical induction

Induction problems in stochastic processes are often trickier than usual. Here are some possibilities:

- Backwards induction: start with base case $n = N$ and go backwards, instead of starting at base case $n = 1$ and going forwards.
- Two-step induction, where the proof for $n = x + 1$ relies not only on the formula being true for $n = x$, but also on it being true for $n = x - 1$.

The first example below is hard probably because it is too easy. The second example is an example of a two-step induction.

Example 1: Suppose that $p_0 = 1$ and $p_x = \alpha p_{x+1}$ for all $x = 1, 2, \dots$. Prove by mathematical induction that $p_n = 1/\alpha^n$ for $n = 0, 1, 2, \dots$

Wish to prove

$$p_n = \frac{1}{\alpha^n} \quad \text{for } n = 0, 1, 2, \dots \quad (\star)$$

Information given:

$$\begin{aligned} p_{x+1} &= \frac{1}{\alpha} p_x & (G_1) \\ p_0 &= 1 & (G_2) \end{aligned}$$

Base case: $n = 0$.

$$LHS = p_0 = 1 \quad \text{by information given } (G_2).$$

$$RHS = \frac{1}{\alpha^0} = \frac{1}{1} = 1 = LHS.$$

Therefore (\star) is true for the base case $n = 0$.

General case: suppose that (\star) is true for $n = x$, so we can assume

$$p_x = \frac{1}{\alpha^x}. \quad (a)$$

Wish to prove that (\star) is also true for $n = x + 1$: i.e.

$$RTP \quad p_{x+1} = \frac{1}{\alpha^{x+1}}. \quad (\star\star)$$

$$\begin{aligned}
 \text{LHS of } (\star\star) = p_{x+1} &= \frac{1}{\alpha} \times p_x && \text{by given } (G_1) \\
 &= \frac{1}{\alpha} \times \frac{1}{\alpha^x} && \text{by allowed (a)} \\
 &= \frac{1}{\alpha^{x+1}} \\
 &= \text{RHS of } (\star\star).
 \end{aligned}$$

So if formula (\star) is true for $n = x$, it is true for $n = x + 1$. We have shown it is true for $n = 0$, so it is true for all $n = 0, 1, 2, \dots$ \square

Example 2: Gambler's Ruin. In the Gambler's Ruin problem in Section 2.7, we have the following situation:

- $p_x = \mathbb{P}(\text{Ruin} \mid \text{start with } \$x)$;
- We know from first-step analysis that $p_{x+1} = 2p_x - p_{x-1}$ (G_1)
- We know from common sense that $p_0 = 1$ (G_2)
- By direct substitution into (G_1) , we obtain:

$$\begin{aligned}
 p_2 &= 2p_1 - 1 \\
 p_3 &= 3p_1 - 2
 \end{aligned}$$

- We develop a suspicion that for all $x = 1, 2, 3, \dots$,

$$p_x = xp_1 - (x - 1) \quad (\star)$$

- We wish to prove (\star) by mathematical induction.

For this example, *our given information, in (G_1) , expresses p_{x+1} in terms of both p_x and p_{x-1} , so we need two base cases. Use $x = 1$ and $x = 2$.*

Wish to prove $p_x = xp_1 - (x - 1)$ (\star) .

Base case $x = 1$:

$$LHS = p_1.$$

$$RHS = 1 \times p_1 - 0 = p_1 = LHS.$$

\therefore formula (\star) is true for base case $x = 1$.

Base case $x = 2$:

$$LHS = p_2 = 2p_1 - 1 \quad \text{by information given } (G_1)$$

$$RHS = 2 \times p_1 - 1 = LHS.$$

\therefore formula (\star) is true for base case $x = 2$.

General case: suppose that (\star) is true for all x up to $x = k$.

So we are allowed:

$$\begin{array}{lll} (x = k) & p_k & = kp_1 - (k - 1) & (a_1) \\ (x = k - 1) & p_{k-1} & = (k - 1)p_1 - (k - 2) & (a_2) \end{array}$$

Wish to prove that (\star) is also true for $x = k + 1$, i.e.

$$RTP \quad p_{k+1} = (k + 1)p_1 - k. \quad (\star\star)$$

$$\begin{aligned} LHS \text{ of } (\star\star) &= p_{k+1} \\ &= 2p_k - p_{k-1} \quad \text{by given information } (G_1) \\ &= 2 \left\{ kp_1 - (k - 1) \right\} - \left\{ (k - 1)p_1 - (k - 2) \right\} \\ &\quad \text{by allowed } (a_1) \text{ and } (a_2) \\ &= p_1 \left\{ 2k - (k - 1) \right\} - \left\{ 2(k - 1) - (k - 2) \right\} \\ &= (k + 1)p_1 - k \\ &= RHS \quad \text{of } (\star\star) \end{aligned}$$

So if formula (\star) is true for $x = k - 1$ and $x = k$, it is true for $x = k + 1$. We have shown it is true for $x = 1$ and $x = 2$, so it is true for all $x = 1, 2, 3, \dots$ \square

Chapter 6: Branching Processes:

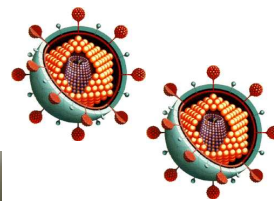
The Theory of Reproduction



Royalty



DNA



Viruses



Aphids



Although the early development of Probability Theory was motivated by problems in gambling, probabilists soon realised that, if they were to continue as a breed, they must also study *reproduction*.



Reproduction is a complicated business, but considerable insights into population growth can be gained from simplified models. The **Branching Process** is a simple but elegant model of population growth. It is also called the **Galton-Watson Process**, because some of the early theoretical results about the process derive from a correspondence between Sir Francis Galton and the Reverend Henry William Watson in 1873. Francis Galton was a cousin of Charles Darwin. In later life, he developed some less elegant ideas about reproduction — namely eugenics, or selective breeding of humans. Luckily he is better remembered for branching processes.

6.1 Branching Processes

Consider some sort of *population* consisting of reproducing individuals.

Examples: living things (animals, plants, bacteria, royal families);
diseases; computer viruses;
rumours, gossip, lies (one lie always leads to another!)

Start conditions: start at time $n = 0$, with a single individual.

Each individual: lives for 1 unit of time. At time $n = 1$, it produces a family of offspring, and immediately dies.

How many offspring? Could be 0, 1, 2, This is the family size, Y . (“Y” stands for “number of Young”).

Each offspring: lives for 1 unit of time. At time $n = 2$, it produces its own family of offspring, and immediately dies.

and so on...

Assumptions

1. All individuals reproduce independently of each other.
2. The family sizes of different individuals are independent, identically distributed random variables. Denote the family size by Y (number of Young).

Family size distribution, Y $\mathbb{P}(Y = k) = p_k$.

y	0	1	2	3	4	...
$P(Y=y)$	p_0	p_1	p_2	p_3	p_4	...

Definition: A branching process is defined as follows.

- Single individual at time $n = 0$.
- Every individual lives exactly one unit of time, then produces Y offspring, and dies.
- The number of offspring, Y , takes values $0, 1, 2, \dots$, and the probability of producing k offspring is $\mathbb{P}(Y = k) = p_k$.
- All individuals reproduce independently. Individuals $1, 2, \dots, n$ have family sizes Y_1, Y_2, \dots, Y_n , where *each Y_i has the same distribution as Y* .
- Let Z_n be the *number of individuals born at time n , for $n = 0, 1, 2, \dots$* . Interpret ' Z_n ' as the 'siZe' of generation n .
- Then the branching process is $\{Z_0, Z_1, Z_2, Z_3, \dots\} = \{Z_n : n \in \mathbb{N}\}$.

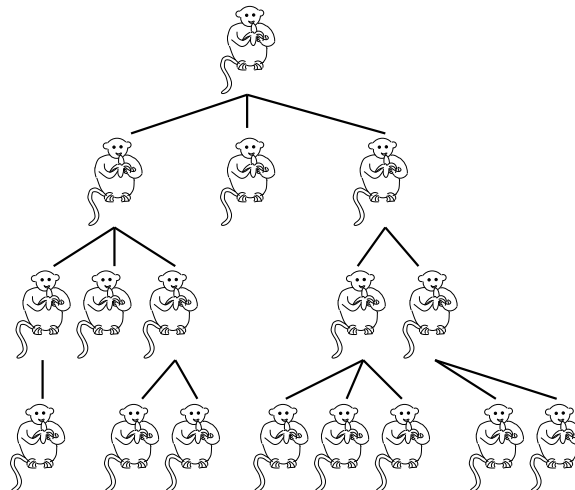
Definition: The state of the branching process at time n is z_n , where *each z_n can take values $0, 1, 2, 3, \dots$* . Note that $z_0 = 1$ always.
 z_n represents the size of the population at time n .

Note: When we want to say that two random variables X and Y have the same distribution, we write: $X \sim Y$.

For example: $Y_i \sim Y$, where Y_i is the family size of any individual i .

Note: The definition of the branching process is easily generalized to start with more than one individual at time $n = 0$.

Branching Process



6.2 Questions about the Branching Process

When we have a situation that can be modelled by a branching process, there are several questions we might want to answer.

If the branching process is just beginning, what will happen in the future?

1. What can we find out about the distribution of Z_n (the population size at generation n)?
 - can we find the mean and variance of Z_n ?
— *yes, using the probability generating function of family size, Y ;*
 - can we find the whole distribution of Z_n ?
— *for special cases of the family size distribution Y , we can find the PGF of Z_n explicitly;*
 - can we find the probability that the population has become extinct by generation n , $\mathbb{P}(Z_n = 0)$?
— *for special cases where we can find the PGF of Z_n (as above).*
2. What can we find out about eventual extinction?
 - can we find the probability of eventual extinction, $\mathbb{P}\left(\lim_{n \rightarrow \infty} Z_n = 0\right)$?
— *yes, always: using the PGF of Y .*
 - can we find general conditions for eventual extinction?
— *yes: we can find conditions that guarantee that extinction will occur with probability 1.*
 - if eventual extinction is definite, can we find the distribution of the time to extinction?
— *for special cases where we can find the PGF of Z_n (as above).*

Example: Modelling cancerous growths. Will a colony of cancerous cells become extinct before it is sufficiently large to overgrow the surrounding tissue?

If the branching process is already in progress, what happened in the past?

1. How long has the process been running?
 - *how many generations do we have to go back to get to the single common ancestor?*
2. What has been the distribution of family size over the generations?
3. What is the total number of individuals (over all generations) up to the present day?

Example: It is believed that all humans are descended from a single female ancestor, who lived in Africa. How long ago?

— *estimated at approximately 200,000 years.*

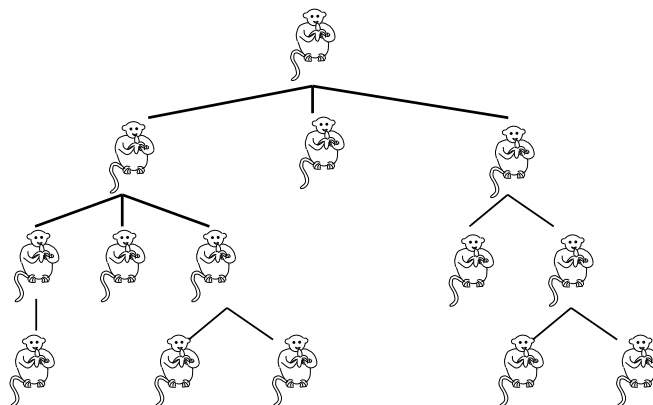
What has been the mean family size over that period?

— *probably very close to 1 female offspring per female adult: e.g. estimate = 1.002.*



6.3 Analysing the Branching Process

Key Observation: *every individual in every generation starts a new, independent branching process, as if the whole process were starting at the beginning again.*



Z_n as a randomly stopped sum

Most of the interesting properties of the branching process centre on the distribution of Z_n (the population size at time n). Using the Key Observation from overleaf, we can find an expression for the probability generating function of Z_n .

Consider the following.

- *The population size at time $n - 1$ is given by Z_{n-1} .*
- *Label the individuals at time $n - 1$ as $1, 2, 3, \dots, Z_{n-1}$.*
- *Each individual $1, 2, \dots, Z_{n-1}$ starts a new branching process. Let $Y_1, Y_2, \dots, Y_{Z_{n-1}}$ be the random family sizes of the individuals $1, 2, \dots, Z_{n-1}$.*
- *The number of individuals at time n , Z_n , is equal to the total number of offspring of the individuals $1, 2, \dots, Z_{n-1}$. That is,*

$$Z_n = \sum_{i=1}^{Z_{n-1}} Y_i.$$

Thus Z_n is a randomly stopped sum: a sum of Y_1, Y_2, \dots , randomly stopped by the random variable Z_{n-1} .

Note: 1. *Each $Y_i \sim Y$: that is, each individual $i = 1, \dots, Z_{n-1}$ has the same family size distribution.*

2. *$Y_1, Y_2, \dots, Y_{Z_{n-1}}$ are independent.*

Probability Generating Function of Z_n

Let $G_Y(s) = \mathbb{E}(s^Y)$ be the probability generating function of Y .

(Recall that Y is the *number of Young of an individual: the family size.*)

Now Z_n is a randomly stopped sum: it is the sum of Y_1, Y_2, \dots , stopped by the random variable Z_{n-1} . So we can use Theorem 4.6 (Chapter 4) to express the PGF of Z_n directly in terms of the PGFs of Y and Z_{n-1} .

By Theorem 4.6, if $Z_n = Y_1 + Y_2 + \dots + Y_{Z_{n-1}}$, and Z_{n-1} is itself random, then the PGF of Z_n is given by:

$$G_{Z_n}(s) = G_{Z_{n-1}}(G_Y(s)), \quad (\clubsuit)$$

where $G_{Z_{n-1}}$ is the PGF of the random variable Z_{n-1} .

For ease of notation, we can write:

$$G_{Z_n}(s) = G_n(s), \quad G_{Z_{n-1}}(s) = G_{n-1}(s), \quad \text{and so on.}$$

Note that $Z_1 = Y$ (the *number of individuals born at time $n = 1$*), so we can also write:

$$G_Y(s) = G_1(s) = G(s) \quad (\text{for simplicity}).$$

Thus, *from* (\clubsuit) ,

$$G_n(s) = G_{n-1}(G(s)) \quad (\text{Branching Process Recursion Formula.})$$

Note:

1. $G_n(s) = \mathbb{E}(s^{Z_n})$, the PGF of the population size at time n , Z_n .
2. $G_{n-1}(s) = \mathbb{E}(s^{Z_{n-1}})$, the PGF of the population size at time $n - 1$, Z_{n-1} .
3. $G(s) = \mathbb{E}(s^Y) = \mathbb{E}(s^{Z_1})$, the PGF of the family size, Y .

We are trying to find the PGF of Z_n , the population size at time n .

So far, we have: $G_n(s) = G_{n-1}(G(s)). \quad (\star)$

But by the same argument,

$$G_{n-1}(r) = G_{n-2}(G(r)).$$

(use r instead of s to avoid confusion in the next line.)

Substituting in (\star) ,

$$\begin{aligned} G_n(s) &= G_{n-1}(G(s)) \\ &= G_{n-1}(r) \quad \text{where } r = G(s) \\ &= G_{n-2}(G(r)) \\ &= G_{n-2}(G(G(s))) \quad \text{replacing } r = G(s). \end{aligned}$$

By the same reasoning, we will obtain:

$$G_n(s) = G_{\underbrace{n-3}_{n-3}} \left(\underbrace{G(G(G(s)))}_{3 \text{ times}} \right),$$

and so on, until we finally get:

$$\begin{aligned} G_n(s) &= G_{n-(n-1)} \left(\underbrace{G(G(G(\dots G(s) \dots)))}_{n-1 \text{ times}} \right) \\ &= \underbrace{G_1}_{=G} \left(\underbrace{G(G(G(\dots G(s) \dots)))}_{n-1 \text{ times}} \right) \\ &= \underbrace{G(G(G(\dots G(s) \dots)))}_{n \text{ times}}. \end{aligned}$$

We have therefore proved the following Theorem.

Theorem 6.3: Let $G(s) = \mathbb{E}(s^Y) = \sum_{y=0}^{\infty} p_y s^y$ be the PGF of the family size distribution, Y . Let $Z_0 = 1$ (start from a single individual at time 0), and let Z_n be the population size at time n ($n = 0, 1, 2, \dots$). Let $G_n(s)$ be the PGF of the random variable Z_n . Then

$$G_n(s) = \underbrace{G\left(G\left(G\left(\dots G(s)\dots\right)\right)\right)}_{n \text{ times}}. \quad \square$$

Note: $G_n(s) = \underbrace{G\left(G\left(G\left(\dots G(s)\dots\right)\right)\right)}_{n \text{ times}}$ is called the *n-fold iterate* of G .

We have therefore found an expression for the PGF of the population size at generation n , although there is no guarantee that it is possible to write it down or manipulate it very easily for large n . For example, if Y has a $\text{Poisson}(\lambda)$ distribution, then $G(s) = e^{\lambda(s-1)}$, and already by generation $n = 3$ we have the following fearsome expression for $G_3(s)$:

$$G_3(s) = e^{\lambda\left(e^{\lambda\left(e^{\lambda(s-1)}-1\right)}-1\right)}. \quad (\text{Or something like that!})$$

However, in some circumstances we can find quite reasonable closed-form expressions for $G_n(s)$, notably when Y has a Geometric distribution. In addition, for any distribution of Y we can use the expression $G_n(s) = G_{n-1}(G(s))$ to derive properties such as the mean and variance of Z_n , and the probability of eventual extinction ($\mathbb{P}(Z_n = 0)$ for some n).

6.4 What does the distribution of Z_n look like?

Before deriving the mean and the variance of Z_n , it is helpful to get some intuitive idea of how the branching process behaves. For example, it seems reasonable to calculate the mean, $\mathbb{E}(Z_n)$, to find out what we expect the population size to be in n generations time, but why are we interested in $\text{Var}(Z_n)$?

The answer is that Z_n usually has a “boom-or-bust” distribution: either the population will take off (boom), and the population size grows quickly, or the population will fail altogether (bust). In fact, if the population fails, it is likely to do so very quickly, within the first few generations. This explains why we are

interested in $\text{Var}(Z_n)$. A huge variance will alert us to the fact that the process does not cluster closely around its mean values. In fact, the mean might be almost useless as a measure of what to expect from the process.

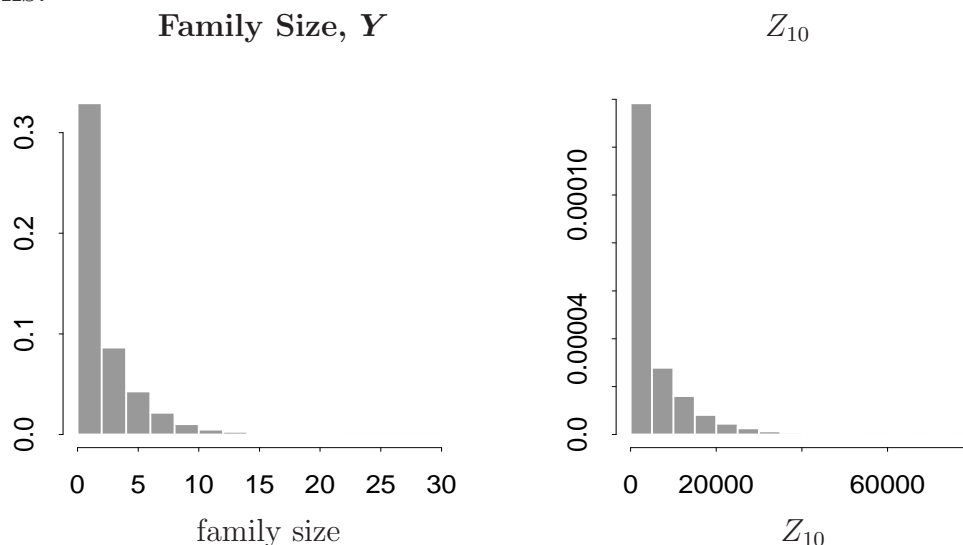
Simulation 1: $Y \sim \text{Geometric}(p = 0.3)$

The following table shows the results from 10 simulations of a branching process, where the family size distribution is $Y \sim \text{Geometric}(p = 0.3)$.

Simulation	Z_0	Z_1	Z_2	Z_3	Z_4	Z_5	Z_6	Z_7	Z_8	Z_9	Z_{10}
1	1	0	0	0	0	0	0	0	0	0	0
2	1	1	0	0	0	0	0	0	0	0	0
3	1	4	19	42	81	181	433	964	2276	5383	12428
4	1	3	3	5	3	15	29	86	207	435	952
5	1	0	0	0	0	0	0	0	0	0	0
6	1	1	0	0	0	0	0	0	0	0	0
7	1	2	8	26	68	162	360	845	2039	4746	10941
8	1	1	0	0	0	0	0	0	0	0	0
9	1	1	0	0	0	0	0	0	0	0	0
10	1	1	4	13	18	39	104	294	690	1566	3534

Often, the population is extinct by generation 10. However, when it is not extinct, it can take enormous values (12428, 10941, ...).

The same simulation was repeated 5000 times to find the empirical distribution of the population size at generation 10 (Z_{10}). The figures below show the distribution of family size, Y , and the distribution of Z_{10} from the 5000 simulations.



Proportion of samples extinct by generation 10: 0.436

Min	1st Qu	Median	Mean	3rd Qu	Max
0	0	1003	4617	6656	82486

Variance of Zn: 53937785.7

For interest, out of the 5000 simulations, there were only 35 (0.7%) that had a value for Z_{10} greater than 0 but less than 100. This emphasizes the “boom-or-bust” nature of the distribution of Z_n .

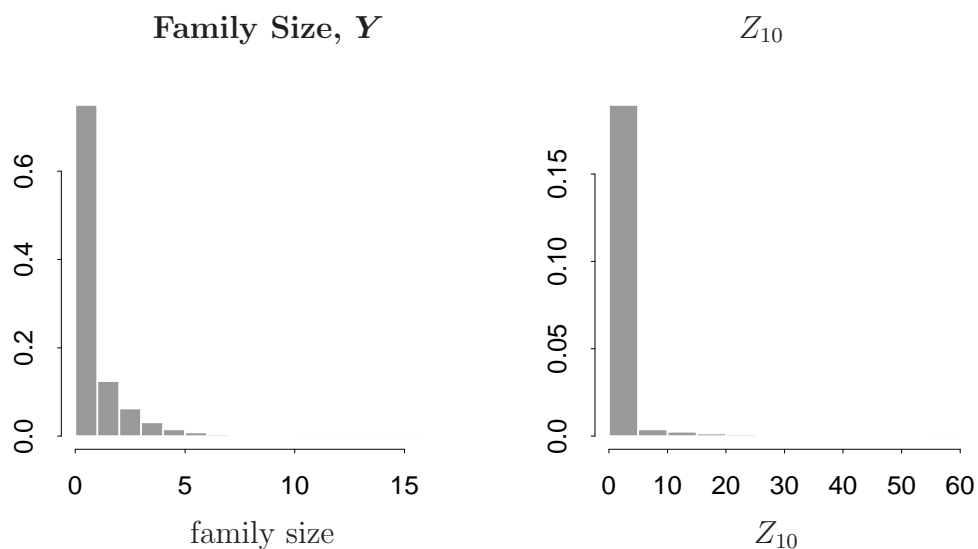
Simulation 2: $Y \sim \text{Geometric}(p = 0.5)$

We repeat the simulation above with a different value for p in the Geometric family size distribution: this time, $p = 0.5$. The family size distribution is therefore $Y \sim \textit{\textbf{Geometric}}(p = 0.5)$.

[illegible]

This time, almost all the populations become extinct. We will see later that this value of p (just) guarantees eventual extinction with probability 1.

The family size distribution, $Y \sim \text{Geometric}(p = 0.5)$, and the results for Z_{10} from 5000 simulations, are shown below. Family sizes are often zero, but families of size 2 and 3 are not uncommon. It seems that this is not enough to save the process from extinction. This time, the maximum population size observed for Z_{10} from 5000 simulations was only 56, and the mean and variance of Z_{10} are much smaller than before.



Proportion of samples extinct by generation 10: 0.9108

Summary of Z_n :

Min	1st Qu	Median	Mean	3rd Qu	Max
0	0	0	0.965	0	56

Mean of Z_n : 0.965

Variance of Z_n : 19.497

What happens for larger values of p ?

It was mentioned above that $Y \sim \text{Geometric}(p = 0.5)$ just guarantees eventual extinction with probability 1. For $p > 0.5$, extinction is also guaranteed, and tends to happen quickly. For example, when $p = 0.55$, over 97% of simulated populations are already extinct by generation 10.

6.5 Mean and variance of Z_n

The previous section has given us a good idea of the significance and interpretation of $\mathbb{E}(Z_n)$ and $\text{Var}(Z_n)$. We now proceed to calculate them. Both $\mathbb{E}(Z_n)$ and $\text{Var}(Z_n)$ can be expressed in terms of *the mean and variance of the family size distribution, Y* .

Thus, *let $\mathbb{E}(Y) = \mu$ and let $\text{Var}(Y) = \sigma^2$. These are the mean and variance of the number of offspring of a single individual.*

Theorem 6.5: Let $\{Z_0, Z_1, Z_2, \dots\}$ be a branching process with $Z_0 = 1$ (start with a single individual). Let Y denote the family size distribution, and suppose that $\mathbb{E}(Y) = \mu$. Then

$$\mathbb{E}(Z_n) = \mu^n.$$

Proof:

By page 121, $Z_n = Y_1 + Y_2 + \dots + Y_{Z_{n-1}}$ is a randomly stopped sum:

$$Z_n = \sum_{i=1}^{Z_{n-1}} Y_i$$

Thus, from Section 3.4 (page 62),

$$\begin{aligned} \mathbb{E}(Z_n) &= \mathbb{E}(Y_i) \times \mathbb{E}(Z_{n-1}) \\ &= \mu \times \mathbb{E}(Z_{n-1}) \\ &= \mu \{ \mu \mathbb{E}(Z_{n-2}) \} \\ &= \mu^2 \mathbb{E}(Z_{n-2}) \\ &= \vdots \\ &= \mu^{n-1} \mathbb{E}(Z_1) \\ &= \mu^{n-1} \times \mu \\ &= \mu^n. \quad \square \end{aligned}$$

Examples: Consider the simulations of Section 6.4.

1. Family size $Y \sim \text{Geometric}(p = 0.3)$. So $\mu = \mathbb{E}(Y) = \frac{q}{p} = \frac{0.7}{0.3} = 2.33$.

Expected population size by generation $n = 10$ is:

$$\mathbb{E}(Z_{10}) = \mu^{10} = (2.33)^{10} = 4784.$$

The theoretical value, 4784, compares well with the sample mean from 5000 simulations, 4617 (page 126).

2. Family size $Y \sim \text{Geometric}(p = 0.5)$. So $\mu = \mathbb{E}(Y) = \frac{q}{p} = \frac{0.5}{0.5} = 1$, and

$$\mathbb{E}(Z_{10}) = \mu^{10} = (1)^{10} = 1.$$

Compares well with the sample mean of 0.965 (page 127).

Variance of Z_n

Theorem 6.5: Let $\{Z_0, Z_1, Z_2, \dots\}$ be a branching process with $Z_0 = 1$ (start with a single individual). Let Y denote the family size distribution, and suppose that $\mathbb{E}(Y) = \mu$ and $\text{Var}(Y) = \sigma^2$. Then

$$\text{Var}(Z_n) = \begin{cases} \sigma^2 n & \text{if } \mu = 1, \\ \sigma^2 \mu^{n-1} \left(\frac{1 - \mu^n}{1 - \mu} \right) & \text{if } \mu \neq 1 \quad (> 1 \text{ or } < 1). \end{cases}$$

Proof:

Write $V_n = \text{Var}(Z_n)$. The proof works by finding a recursive formula for V_n .

Using the Law of Total Variance for randomly stopped sums from Section 3.4 (page 62),

$$\begin{aligned}
 Z_n &= \sum_{i=1}^{Z_{n-1}} Y_i \\
 \Rightarrow \text{Var}(Z_n) &= \{\mathbb{E}(Y_i)\}^2 \times \text{Var}(Z_{n-1}) + \text{Var}(Y_i) \times \mathbb{E}(Z_{n-1}) \\
 \Rightarrow V_n &= \mu^2 V_{n-1} + \sigma^2 \mathbb{E}(Z_{n-1}) \\
 \Rightarrow V_n &= \mu^2 V_{n-1} + \sigma^2 \mu^{n-1},
 \end{aligned}$$

using $\mathbb{E}(Z_{n-1}) = \mu^{n-1}$ as above.

Also,

$$V_1 = \text{Var}(Z_1) = \text{Var}(Y) = \sigma^2.$$

Find V_n by repeated substitution:

$$V_1 = \sigma^2$$

$$V_2 = \mu^2 V_1 + \sigma^2 \mu = \mu^2 \sigma^2 + \mu \sigma^2 = \mu \sigma^2 (1 + \mu)$$

$$V_3 = \mu^2 V_2 + \sigma^2 \mu^2 = \mu^2 \sigma^2 (1 + \mu + \mu^2)$$

$$V_4 = \mu^2 V_3 + \sigma^2 \mu^3 = \mu^3 \sigma^2 (1 + \mu + \mu^2 + \mu^3)$$

\vdots etc.

Completing the pattern,

$$\begin{aligned}
 V_n &= \mu^{n-1} \sigma^2 (1 + \mu + \mu^2 + \dots + \mu^{n-1}) \\
 &= \mu^{n-1} \sigma^2 \sum_{r=0}^{n-1} \mu^r \\
 &= \mu^{n-1} \sigma^2 \left(\frac{1 - \mu^n}{1 - \mu} \right). \quad \text{Valid for } \mu \neq 1. \\
 &\quad \text{(sum of first } n \text{ terms of Geometric series)}
 \end{aligned}$$

When $\mu = 1$:

$$V_n = 1^{n-1} \sigma^2 \underbrace{(1^0 + 1^1 + \dots + 1^{n-1})}_{n \text{ times}} = \sigma^2 n.$$

Hence the result:

$$\text{Var}(Z_n) = \begin{cases} \sigma^2 n & \text{if } \mu = 1, \\ \sigma^2 \mu^{n-1} \left(\frac{1 - \mu^n}{1 - \mu} \right) & \text{if } \mu \neq 1. \end{cases} \quad \square$$

Examples: Again consider the simulations of Section 6.4.

1. Family size $Y \sim \text{Geometric}(p = 0.3)$. So $\mu = \mathbb{E}(Y) = \frac{q}{p} = \frac{0.7}{0.3} = 2.33$.

$$\sigma^2 = \text{Var}(Y) = \frac{q}{p^2} = \frac{0.7}{(0.3)^2} = 7.78.$$

$$\text{Var}(Z_{10}) = \sigma^2 \mu^9 \left(\frac{1 - \mu^{10}}{1 - \mu} \right) = 5.72 \times 10^7.$$

Compares well with the sample variance from 5000 simulations, 5.39×10^7 (page 126).

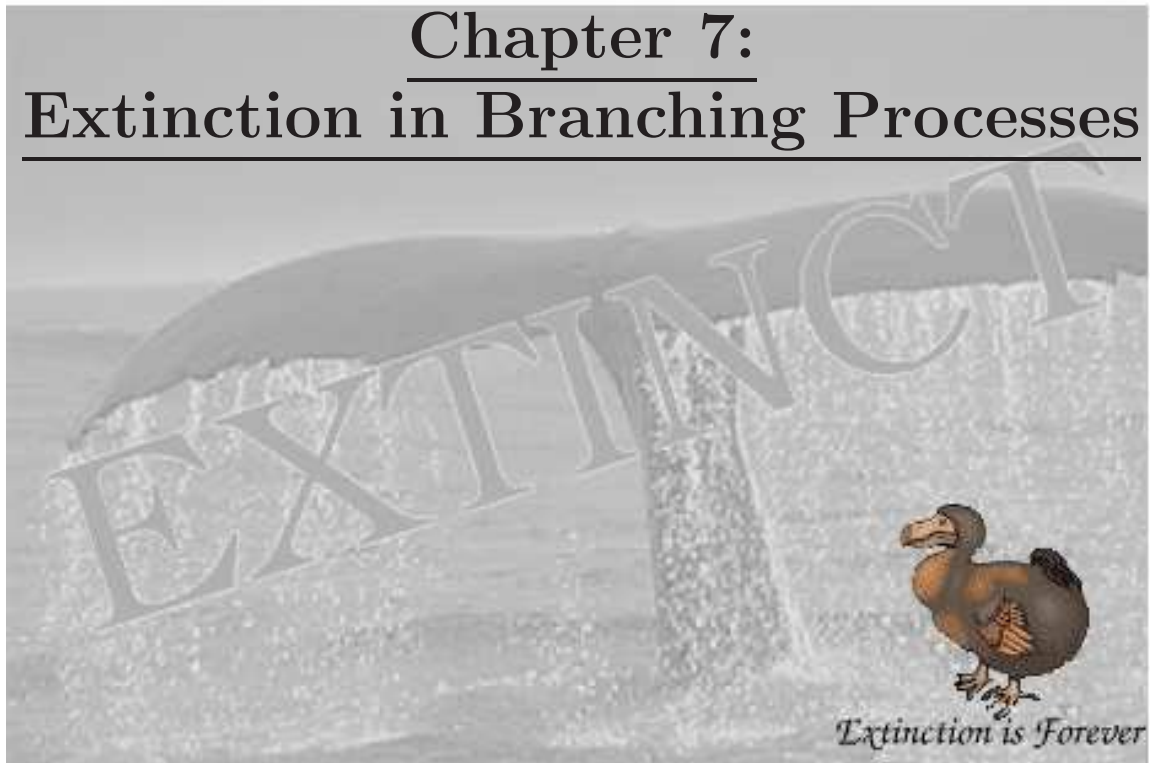
2. Family size $Y \sim \text{Geometric}(p = 0.5)$. So $\mu = \mathbb{E}(Y) = \frac{q}{p} = \frac{0.5}{0.5} = 1$.

$\sigma^2 = \text{Var}(Y) = \frac{q}{p^2} = \frac{0.5}{(0.5)^2} = 2$. Using the formula for $\text{Var}(Z_n)$ when $\mu = 1$, we have:

$$\text{Var}(Z_{10}) = \sigma^2 n = 2 \times 10 = 20.$$

Compares well with the sample variance of 19.5 (page 127).

Chapter 7: Extinction in Branching Processes



Revision: a branching process consists of reproducing individuals.

- All individuals are independent.
- Start with a single individual at time 0: $Z_0 = 1$.
- Each individual lives a single unit of time, then has Y offspring and dies.
- Let Z_n be the size of generation n : the number of individuals born at time n .
- The branching process is $\{Z_0 = 1, Z_1, Z_2, \dots\}$.

Branching Process Recursion Formula

This is the fundamental formula for branching processes. Let $G_n(s) = \mathbb{E}(s^{Z_n})$ be the PGF of Z_n , the population size at time n . Let $G(s) = G_1(s)$, the PGF of the family size distribution Y , or equivalently, of Z_1 . Then:

$$G_n(s) = G_{n-1}(G(s)) = \underbrace{G(G(\dots G(s)\dots))}_{n \text{ times}} = G(G_{n-1}(s)).$$

7.1 Extinction Probability

One of the most interesting applications of branching processes is calculating the probability of eventual extinction. For example, what is the probability that a colony of cancerous cells becomes extinct before it overgrows the surrounding tissue? What is the probability that an infectious disease dies out before reaching an epidemic? What is the probability that a family line (e.g. for royal families) becomes extinct?

It is possible to find several results about the probability of eventual extinction.

Extinction by generation n

The population is extinct by generation n if $Z_n = 0$
(no individuals at time n).

If $Z_n = 0$, then *the population is extinct for ever*: $Z_t = 0$ for all $t \geq n$.



Definition: Define event E_n to be the event
 $E_n = \{Z_n = 0\}$ (event that the population is extinct by generation n).

Note: $E_0 \subseteq E_1 \subseteq E_2 \subseteq E_3 \subseteq E_4 \subseteq \dots$

This is because event E_i forces E_j to be true for all $j \geq i$, so E_i is a ‘part’ or subset of E_j for $j \geq i$.

Ultimate extinction

At the start of the branching process, we are interested in the probability of *ultimate extinction*: *the probability that the population will be extinct by generation n , for any value of n .*

We can express this probability in different ways:

$$\mathbb{P}(\text{ultimate extinction}) = \mathbb{P}\left(\bigcup_{n=0}^{\infty} E_n\right) \left(\begin{array}{l} \text{i.e. extinct by generation 0 \underline{or}} \\ \text{extinct by generation 1 \underline{or}} \\ \text{extinct by generation 2 \underline{or} \dots} \end{array} \right)$$

Or: $\mathbb{P}(\text{ultimate extinction}) = \mathbb{P}\left(\lim_{n \rightarrow \infty} E_n\right)$. (i.e. $\mathbb{P}(\text{extinct by generation } \infty)$).

Note: By the Continuity Theorem (Chapter 2), and because $E_0 \subseteq E_1 \subseteq E_2 \subseteq \dots$, we have:

$$\mathbb{P}(\text{ultimate extinction}) = \mathbb{P}\left(\lim_{n \rightarrow \infty} E_n\right) = \lim_{n \rightarrow \infty} \mathbb{P}(E_n).$$

Thus the probability of eventual extinction is the limit as $n \rightarrow \infty$ of the probability of extinction by generation n .

We will use the Greek letter Gamma (γ) for the probability of extinction: think of Gamma for ‘all Gone’!

$$\gamma_n = \mathbb{P}(E_n) = \mathbb{P}(\text{extinct by generation } n).$$

$$\gamma = \mathbb{P}(\text{ultimate extinction}).$$

By the Note above, we have established that we are looking for:



$$\mathbb{P}(\text{ultimate extinction}) = \gamma = \lim_{n \rightarrow \infty} \gamma_n.$$



Extinction is Forever

Theorem 7.1: Let γ be the probability of ultimate extinction. Then

γ is the smallest non-negative solution of the equation

$G(s) = s$, where $G(s)$ is the PGF of the family size distribution, Y .

To find the probability of ultimate extinction, we therefore:

- find the PGF of family size, Y : $G(s) = \mathbb{E}(s^Y)$;
- find values of s that satisfy $G(s) = s$;
- find the smallest of these values that is ≥ 0 . This is the required value γ .

$$G(\gamma) = \gamma, \text{ and } \gamma \text{ is the smallest value } \geq 0 \text{ for which this holds.}$$

Note: Recall that, for any (non-defective) random variable Y with PGF $G(s)$,

$$G(1) = \mathbb{E}(1^Y) = \sum_y 1^y \mathbb{P}(Y = y) = \sum_y \mathbb{P}(Y = y) = 1.$$

So $G(1) = 1$ always, and therefore *there always exists a solution for $G(s) = s$ in $[0, 1]$.*

The required value γ is the smallest such solution ≥ 0 .

Before proving Theorem 7.1 we prove the following Lemma.

Lemma: Let $\gamma_n = \mathbb{P}(Z_n = 0)$. Then $\gamma_n = G(\gamma_{n-1})$.

Proof: If $G_n(s)$ is the PGF of Z_n , then $\mathbb{P}(Z_n = 0) = G_n(0)$. (Chapter 4.)

So $\gamma_n = G_n(0)$. Similarly, $\gamma_{n-1} = G_{n-1}(0)$.

$$\text{Now } G_n(0) = \underbrace{G\left(G\left(G\left(\dots G(0)\dots\right)\right)\right)}_{n \text{ times}} = G\left(G_{n-1}(0)\right).$$

$$\text{So } \gamma_n = G\left(G_{n-1}(0)\right) = G\left(\gamma_{n-1}\right). \quad \square$$

Proof of Theorem 7.1: We need to prove:

(i) $G(\gamma) = \gamma$;

(ii) γ is the smallest non-negative value for which $G(\gamma) = \gamma$.

That is, if $s \geq 0$ and $G(s) = s$, then $\gamma \leq s$.

Proof of (i):

$$\begin{aligned} \text{From } \text{overleaf}, \quad \gamma &= \lim_{n \rightarrow \infty} \gamma_n = \lim_{n \rightarrow \infty} G\left(\gamma_{n-1}\right) && \text{(by Lemma)} \\ &= G\left(\lim_{n \rightarrow \infty} \gamma_{n-1}\right) && (G \text{ is continuous}) \\ &= G(\gamma). \end{aligned}$$

So $G(\gamma) = \gamma$, as required.

Proof of (ii):

First note that $G(s)$ is an increasing function on $[0, 1]$:

$$\begin{aligned} G(s) = \mathbb{E}(s^Y) &= \sum_{y=0}^{\infty} s^y \mathbb{P}(Y = y) \\ \Rightarrow G'(s) &= \sum_{y=0}^{\infty} y s^{y-1} \mathbb{P}(Y = y) \\ \Rightarrow G'(s) &\geq 0 \quad \text{for } 0 \leq s \leq 1, \quad \text{so } G \text{ is increasing on } [0, 1]. \end{aligned}$$

$G(s)$ is increasing on $[0, 1]$ means that:

$$s_1 \leq s_2 \quad \Rightarrow \quad G(s_1) \leq G(s_2) \quad \text{for any } s_1, s_2 \in [0, 1]. \quad \clubsuit$$

The branching process begins with $Z_0 = 1$, so

$$\mathbb{P}(\text{extinct by generation } 0) = \gamma_0 = 0.$$

At any later generation, $\gamma_n = G(\gamma_{n-1})$ by Lemma.

Now suppose that $s \geq 0$ and $G(s) = s$. Then we have:

$$\begin{aligned} 0 \leq s &\Rightarrow \gamma_0 \leq s && \text{(because } \gamma_0 = 0) \\ &\Rightarrow G(\gamma_0) \leq G(s) && \text{(by } \clubsuit) \\ \text{i.e.} &\quad \gamma_1 \leq s \\ &\Rightarrow G(\gamma_1) \leq G(s) && \text{(by } \clubsuit) \\ \text{i.e.} &\quad \gamma_2 \leq s \\ &\quad \vdots \end{aligned}$$

Thus $\gamma_n \leq s$ for all n .

So if $s \geq 0$ and $G(s) = s$, then $\gamma = \lim_{n \rightarrow \infty} \gamma_n \leq s$. □

Example 1: Let $\{Z_0 = 1, Z_1, Z_2, \dots\}$ be a branching process with family size distribution $Y \sim \text{Binomial}(2, \frac{1}{4})$. Find the probability that the process will eventually die out.

Solution:

Let $G(s) = \mathbb{E}(s^Y)$. The probability of ultimate extinction is γ , where γ is the smallest solution ≥ 0 to the equation $G(s) = s$.

For $Y \sim \text{Binomial}(n, p)$, the PGF is $G(s) = (ps + q)^n$ (Chapter 4).

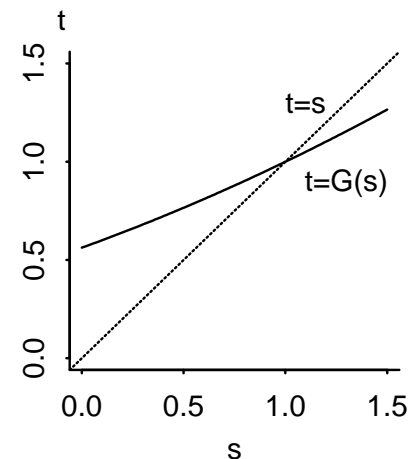
So if $Y \sim \text{Binomial}(2, \frac{1}{4})$ then $G(s) = (\frac{1}{4}s + \frac{3}{4})^2$.

We need to solve $G(s) = s$:

$$G(s) = (\frac{1}{4}s + \frac{3}{4})^2 = s$$

$$\frac{1}{16}s^2 + \frac{6}{16}s + \frac{9}{16} = s$$

$$\frac{1}{16}s^2 - \frac{10}{16}s + \frac{9}{16} = 0$$



Trick: we know that $G(1) = 1$, so $s = 1$ has got to be a solution. Use this for a quick factorization.

$$(s - 1) \left(\frac{1}{16}s - \frac{9}{16} \right) = 0.$$

Thus

$$s = 1$$

or

$$\frac{1}{16}s = \frac{9}{16} \Rightarrow s = 9.$$

The smallest solution ≥ 0 is $s = 1$.

Thus the probability of ultimate extinction is $\gamma = 1$.



Extinction is definite when the family size distribution is $Y \sim \text{Binomial}(2, \frac{1}{4})$.

Example 2: Let $\{Z_0 = 1, Z_1, Z_2, \dots\}$ be a branching process with family size distribution $Y \sim \text{Geometric}(\frac{1}{4})$. Find the probability that the process will eventually die out.

Solution:

Let $G(s) = \mathbb{E}(s^Y)$. Then $\mathbb{P}(\text{ultimate extinction}) = \gamma$, where γ is the smallest solution ≥ 0 to the equation $G(s) = s$.

For $Y \sim \text{Geometric}(p)$, the PGF is $G(s) = \frac{p}{1-qs}$ (Chapter 4).

So if $Y \sim \text{Geometric}(\frac{1}{4})$ then $G(s) = \frac{1/4}{1 - (3/4)s} = \frac{1}{4 - 3s}$.

We need to solve $G(s) = s$:

$$G(s) = \frac{1}{4-3s} = s$$

$$4s - 3s^2 = 1$$

$$3s^2 - 4s + 1 = 0$$

Trick: know that $s = 1$ is a solution.

$$(s - 1)(3s - 1) = 0.$$

Thus

$$s = 1$$

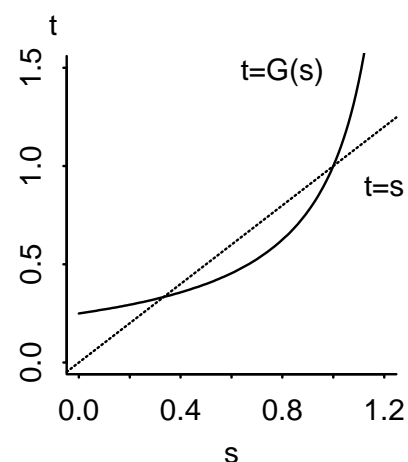
or

$$3s = 1 \Rightarrow s = \frac{1}{3}.$$

The smallest solution ≥ 0 is $s = \frac{1}{3}$.

Thus the probability of ultimate extinction is $\gamma = \frac{1}{3}$.

Extinction is possible but not definite when the family size distribution is $Y \sim \text{Geometric}(\frac{1}{4})$.



7.2 Conditions for ultimate extinction

It turns out that the probability of extinction depends crucially on the value of μ , *the mean of the family size distribution* Y .

Some values of μ *guarantee* that the branching process will die out with probability 1. Other values guarantee that the probability of extinction will be strictly less than 1. We will see below that the threshold value is $\mu = 1$.

If the mean number of offspring per individual μ is more than 1 (so on average, individuals replace themselves plus a bit extra), then the branching process is not *guaranteed* to die out — although it might do. However, if the mean number of offspring per individual μ is 1 or less, the process is *guaranteed* to become extinct (unless $Y = 1$ with probability 1). The result is not too surprising for $\mu > 1$ or $\mu < 1$, but it is a little surprising that extinction is generally guaranteed if $\mu = 1$.

Theorem 7.2: Let $\{Z_0 = 1, Z_1, Z_2, \dots\}$ be a branching process with family size distribution Y . Let $\mu = \mathbb{E}(Y)$ be the mean family size distribution, and let γ be the probability of ultimate extinction. Then

- (i) If $\mu > 1$, then $\gamma < 1$: extinction is not guaranteed if $\mu > 1$.
- (ii) If $\mu < 1$, then $\gamma = 1$: extinction is guaranteed if $\mu < 1$.
- (iii) If $\mu = 1$, then $\gamma = 1$ unless the family size is always constant at $Y = 1$.

Lemma: Let $G(s)$ be the PGF of family size Y . Then $G(s)$ and $G'(s)$ are strictly increasing for $0 < s < 1$, as long as Y can take values ≥ 2 .

Proof: $G(s) = \mathbb{E}(s^Y) = \sum_{y=0}^{\infty} s^y \mathbb{P}(Y = y)$.

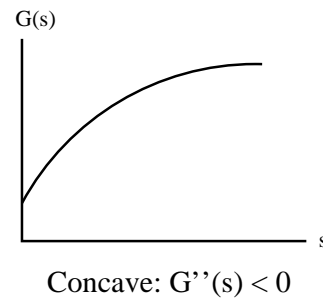
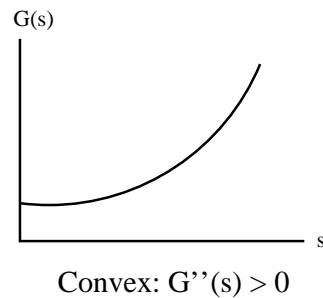
So $G'(s) = \sum_{y=1}^{\infty} y s^{y-1} \mathbb{P}(Y = y) > 0$ for $0 < s < 1$,

because all terms are ≥ 0 and at least 1 term is > 0 (if $\mathbb{P}(Y \geq 2) > 0$).

Similarly, $G''(s) = \sum_{y=2}^{\infty} y(y-1) s^{y-2} \mathbb{P}(Y = y) > 0$ for $0 < s < 1$.

So $G(s)$ and $G'(s)$ are strictly increasing for $0 < s < 1$. \square

Note: When $G''(s) > 0$ for $0 < s < 1$, the function G is said to be convex on that interval.

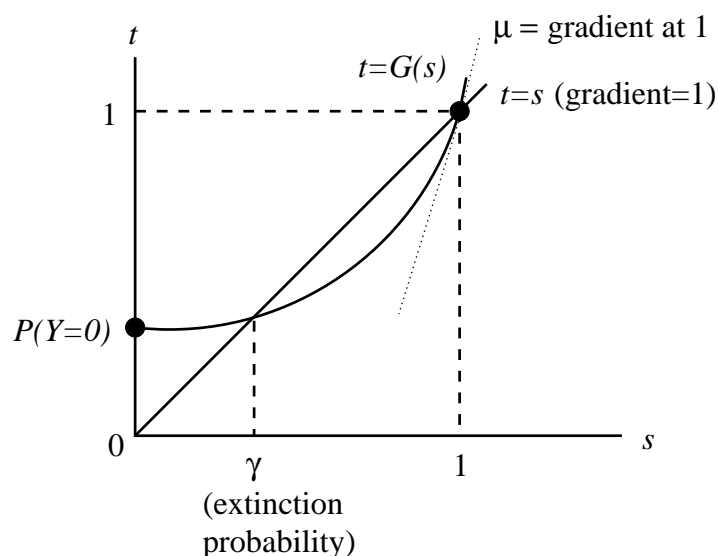


$G''(s) > 0$ means that the gradient of G is constantly increasing for $0 < s < 1$.

Proof of Theorem 7.2: This is usually done graphically.

The graph of $G(s)$ satisfies the following conditions:

1. $G(s)$ is increasing and strictly convex (as long as Y can be ≥ 2).
2. $G(0) = \mathbb{P}(Y = 0) \geq 0$.
3. $G(1) = 1$.
4. $G'(1) = \mu$, so the slope of $G(s)$ at $s = 1$ gives the value μ .
5. The extinction probability γ is the smallest value ≥ 0 for which $G(s) = s$.

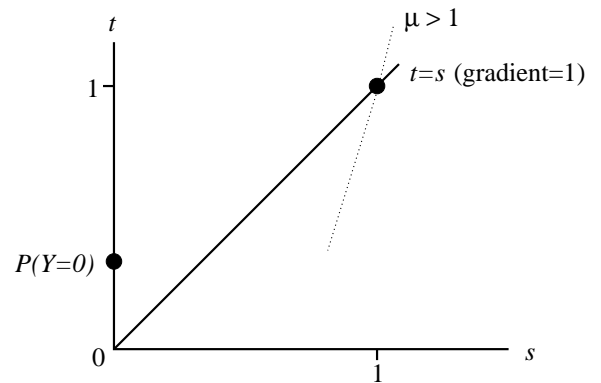


Case (i): $\mu > 1$

When $\mu > 1$, the curve $G(s)$ is forced beneath the line $t = s$ at $s = 1$.

The curve $G(s)$ has to cross the line $t = s$ again to meet the t -axis at $\mathbb{P}(Y = 0)$.

Thus there must be a solution $\gamma < 1$ to the equation $G(s) = s$.



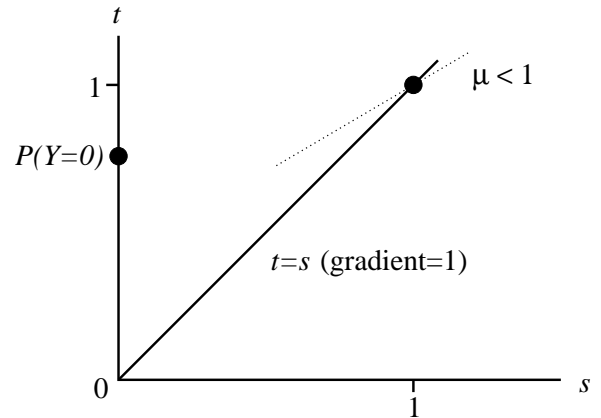
Case (ii): $\mu < 1$

When $\mu < 1$, the curve $G(s)$ is forced above the line $t = s$ for $s < 1$.

There is no possibility for the curve $G(s)$ to cross the line $t = s$ again before meeting the t -axis.

Thus there can be no solution < 1 to the equation $G(s) = s$, so $\gamma = 1$.

The exception is where Y can take only values 0 and 1, so $G(s)$ is not strictly convex (see Lemma). However, in that case $G(s) = p_0 + p_1 s$ is a straight line, giving the same result $\gamma = 1$.

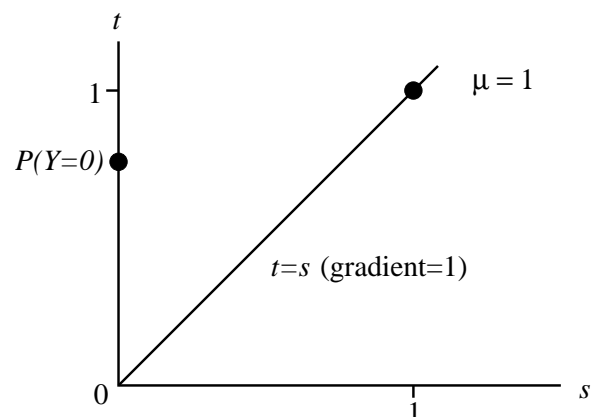


Case (iii): $\mu = 1$

When $\mu = 1$, the situation is the same as for $\mu < 1$.

The exception is where Y takes only the value 1. Then $G(s) = s$ for all $0 \leq s \leq 1$, so the smallest solution ≥ 0 is $\gamma = 0$.

Thus extinction is guaranteed for $\mu = 1$, **unless** $Y = 1$ with probability 1.



Example 1: Let $\{Z_0 = 1, Z_1, Z_2, \dots\}$ be a branching process with family size distribution $Y \sim \text{Binomial}(2, \frac{1}{4})$, as in Section 7.1. Find the probability of eventual extinction.

Solution:

Consider $Y \sim \text{Binomial}(2, \frac{1}{4})$. The mean of Y is $\mu = 2 \times \frac{1}{4} = \frac{1}{2} < 1$. Thus, by Theorem 7.2,

$$\gamma = \mathbb{P}(\text{ultimate extinction}) = 1.$$

(The longer calculation in Section 7.1 was not necessary.)

Example 2: Let $\{Z_0 = 1, Z_1, Z_2, \dots\}$ be a branching process with family size distribution $Y \sim \text{Geometric}(\frac{1}{4})$, as in Section 7.1. Find the probability of eventual extinction.

Solution:

Consider $Y \sim \text{Geometric}(\frac{1}{4})$. The mean of Y is $\mu = \frac{1-1/4}{1/4} = 3 > 1$. Thus, by Theorem 7.2,

$$\gamma = \mathbb{P}(\text{ultimate extinction}) < 1.$$

To find the value of γ , we still need to go through the calculation presented in Section 7.1. (Answer: $\gamma = \frac{1}{3}$.)

Note: The mean μ of the offspring distribution Y is known as the *criticality parameter*.

- If $\mu < 1$, extinction is definite ($\gamma = 1$). The process is called subcritical. Note that $\mathbb{E}(Z_n) = \mu^n \rightarrow 0$ as $n \rightarrow \infty$.
- If $\mu = 1$, extinction is definite unless $Y \equiv 1$. The process is called critical. Note that $\mathbb{E}(Z_n) = \mu^n = 1 \forall n$, even though extinction is definite.
- If $\mu > 1$, extinction is not definite ($\gamma < 1$). The process is called supercritical. Note that $\mathbb{E}(Z_n) = \mu^n \rightarrow \infty$ as $n \rightarrow \infty$.



But how long have you got...?

7.3 Time to Extinction

Suppose the population is doomed to extinction — or maybe it isn't. Either way, it is useful to know how long it will take for the population to become extinct. This is the distribution of T , the number of generations before extinction. For example, how long do we expect a disease epidemic like SARS to continue? How long have we got to organize ourselves to save the kakapo or the tuatara before they become extinct before our very eyes?



1. Extinction by time n

The branching process is extinct by time n if $Z_n = 0$.

Thus the probability that the process has become extinct by time n is:

$$\mathbb{P}(Z_n = 0) = G_n(0) = \gamma_n.$$

Note: Recall that $G_n(s) = \mathbb{E}(s^{Z_n}) = \underbrace{G\left(G\left(G\left(\dots G(s)\dots\right)\right)\right)}_{n \text{ times}}.$

There is no guarantee that the PGF $G_n(s)$ or the value $G_n(0)$ can be calculated easily. However, we can build up $G_n(0)$ in steps:

e.g. $G_2(0) = G(G(0));$ then $G_3(0) = G(G_2(0)),$ or even $G_4(0) = G_2(G_2(0)).$

2. Extinction at time n

Let T be the exact time of extinction. That is, $T = n$ if generation n is the first generation with no individuals:

$$T = n \iff Z_n = 0 \text{ AND } Z_{n-1} > 0.$$

Now by the Partition Rule,

$$\mathbb{P}(Z_n = 0 \cap Z_{n-1} > 0) + \mathbb{P}(Z_n = 0 \cap Z_{n-1} = 0) = \mathbb{P}(Z_n = 0). \quad (\star)$$

But the event $\{Z_n = 0 \cap Z_{n-1} = 0\}$ is the event that the process is extinct by generation $n - 1$ AND it is extinct by generation n . However, we know it will always be extinct by generation n if it is extinct by generation $n - 1$, so the $Z_n = 0$ part is redundant. So

$$\mathbb{P}(Z_n = 0 \cap Z_{n-1} = 0) = \mathbb{P}(Z_{n-1} = 0) = G_{n-1}(0).$$

Similarly,

$$\mathbb{P}(Z_n = 0) = G_n(0).$$

So (\star) gives:

$$\mathbb{P}(T = n) = \mathbb{P}(Z_n = 0 \cap Z_{n-1} > 0) = G_n(0) - G_{n-1}(0) = \gamma_n - \gamma_{n-1}.$$

This gives the distribution of T , the exact time at which extinction occurs.

Example: Binary splitting. Suppose that the family size distribution is

$$Y = \begin{cases} 0 & \text{with probability } q = 1 - p, \\ 1 & \text{with probability } p. \end{cases}$$

Find the distribution of the time to extinction.

Solution:

Consider

$$G(s) = \mathbb{E}(s^Y) = qs^0 + ps^1 = q + ps.$$

$$G_2(s) = G(G(s)) = q + p(q + ps) = q(1 + p) + p^2s.$$

$$G_3(s) = G(G_2(s)) = q + p(q + pq + p^2s) = q(1 + p + p^2) + p^3s.$$

\vdots

$$G_n(s) = q(1 + p + p^2 + \dots + p^{n-1}) + p^n s.$$

Thus time to extinction, T , satisfies

$$\begin{aligned} \mathbb{P}(T = n) &= G_n(0) - G_{n-1}(0) \\ &= q(1 + p + p^2 + \dots + p^{n-1}) - q(1 + p + p^2 + \dots + p^{n-2}) \\ &= qp^{n-1} \quad \text{for } n = 1, 2, \dots \end{aligned}$$

Thus

$$T - 1 \sim \text{Geometric}(q).$$

It follows that $\mathbb{E}(T - 1) = \frac{p}{q}$, so

$$\mathbb{E}(T) = 1 + \frac{p}{q} = \frac{1 - p + p}{q} = \frac{1}{q}.$$

Note: The expected time to extinction, $\mathbb{E}(T)$, is:

- finite if $\mu < 1$;
- infinite if $\mu = 1$ (despite extinction being definite), if σ^2 is finite;
- infinite if $\mu > 1$ (because with positive probability, extinction never happens).

(Results not proved here.)

7.4 Case Study: Geometric Branching Processes

Recall that $G_n(s) = \mathbb{E}(s^{Z_n}) = \underbrace{G\left(G\left(G\left(\dots G(s)\dots\right)\right)\right)}_{n \text{ times}}.$

In general, it is not possible to find a closed-form expression for $G_n(s)$. We achieved a closed-form $G_n(s)$ in the Binary Splitting example (page 144), but binary splitting only allows family size Y to be 0 or 1, which is a very restrictive model.

The only non-trivial family size distribution that allows us to find a closed-form expression for $G_n(s)$ is the **Geometric distribution**.

When family size $Y \sim \text{Geometric}(p)$, we can do the following:

- Derive a closed-form expression for $G_n(s)$, the PGF of Z_n .
- Find the probability distribution of the **exact time of extinction, T** : not just the probability that extinction will occur at some unspecified time (γ).
- Find the **full probability distribution of Z_n** : probabilities $\mathbb{P}(Z_n = 0)$, $\mathbb{P}(Z_n = 1)$, $\mathbb{P}(Z_n = 2)$, \dots .

With $Y \sim \text{Geometric}(p)$, we can therefore calculate just about every quantity we might be interested in for the branching process.

1. Closed form expression for $G_n(s)$

Theorem 7.4: Let $\{Z_0 = 1, Z_1, Z_2, \dots\}$ be a branching process with family size distribution $Y \sim \text{Geometric}(p)$. The PGF of Z_n is given by:

$$G_n(s) = \mathbb{E}(s^{Z_n}) = \begin{cases} \frac{n - (n-1)s}{n+1 - ns} & \text{if } p = q = 0.5, \\ \frac{(\mu^n - 1) - \mu(\mu^{n-1} - 1)s}{(\mu^{n+1} - 1) - \mu(\mu^n - 1)s} & \text{if } p \neq q, \text{ where } \mu = \frac{q}{p}. \end{cases}$$

Proof (sketch):

The proof for both $p = q$ and $p \neq q$ proceed by mathematical induction. We will give a sketch of the proof when $p = q = 0.5$. The proof for $p \neq q$ works in the same way but is trickier.

Consider $p = q = \frac{1}{2}$. Then

$$G(s) = \frac{p}{1 - qs} = \frac{\frac{1}{2}}{1 - \frac{s}{2}} = \frac{1}{2 - s}.$$

Using the Branching Process Recursion Formula (Chapter 6),

$$G_2(s) = G(G(s)) = \frac{1}{2 - G(s)} = \frac{1}{2 - \frac{1}{2-s}} = \frac{2 - s}{2(2 - s) - 1} = \frac{2 - s}{3 - 2s}.$$

The inductive hypothesis is that $G_n(s) = \frac{n - (n - 1)s}{n + 1 - ns}$, and it holds for $n = 1$ and $n = 2$. Suppose it holds for n . Then

$$\begin{aligned} G_{n+1}(s) &= G_n(G(s)) = \frac{n - (n - 1)G(s)}{n + 1 - nG(s)} = \frac{n - (n - 1)\left(\frac{1}{2-s}\right)}{n + 1 - n\left(\frac{1}{2-s}\right)} \\ &= \frac{(2 - s)n - (n - 1)}{(2 - s)(n + 1) - n} \\ &= \frac{n + 1 - ns}{n + 2 - (n + 1)s}. \end{aligned}$$

Therefore, if the hypothesis holds for n , it also holds for $n + 1$. Thus the hypothesis is proved for all n . \square

2. Exact time of extinction, T

Let $Y \sim \text{Geometric}(p)$, and let T be the exact generation of extinction.

From Section 7.3,

$$\mathbb{P}(T = n) = \mathbb{P}(Z_n = 0) - \mathbb{P}(Z_{n-1} = 0) = G_n(0) - G_{n-1}(0).$$

By using the closed-form expressions overleaf for $G_n(0)$ and $G_{n-1}(0)$, we can find $\mathbb{P}(T = n)$ for any n .

3. Whole distribution of Z_n

From Chapter 4, $\mathbb{P}(Z_n = r) = \frac{1}{r!} G_n^{(r)}(0)$.

Now our closed-form expression for $G_n(s)$ has the same format regardless of whether $\mu = 1$ ($p = 0.5$), or $\mu \neq 1$ ($p \neq 0.5$):

$$G_n(s) = \frac{A - Bs}{C - Ds}.$$

(For example, when $\mu = 1$, we have $A = D = n$, $B = n - 1$, $C = n + 1$.) Thus:

$$\mathbb{P}(Z_n = 0) = G_n(0) = \frac{A}{C}$$

$$G'_n(s) = \frac{(C - Ds)(-B) + (A - Bs)D}{(C - Ds)^2} = \frac{AD - BC}{(C - Ds)^2}$$

$$\Rightarrow \mathbb{P}(Z_n = 1) = \frac{1}{1!} G'_n(0) = \frac{AD - BC}{C^2}$$

$$G''_n(s) = \frac{(-2)(-D)(AD - BC)}{(C - Ds)^3} = \frac{2D(AD - BC)}{(C - Ds)^3}$$

$$\Rightarrow \mathbb{P}(Z_n = 2) = \frac{1}{2!} G''_n(0) = \left(\frac{AD - BC}{CD} \right) \left(\frac{D}{C} \right)^2$$

\vdots

$$\Rightarrow \mathbb{P}(Z_n = r) = \frac{1}{r!} G_n^{(r)}(0) = \left(\frac{AD - BC}{CD} \right) \left(\frac{D}{C} \right)^r \quad \text{for } r = 1, 2, \dots$$

(Exercise)

This is very simple and powerful: we can substitute the values of A, B, C , and D to find $\mathbb{P}(Z_n = r)$ or $\mathbb{P}(Z_n \leq r)$ for any r and n .

Note: A Java applet that simulates branching processes can be found at:

http://www.dartmouth.edu/~chance/teaching_aids/books_articles/probability_book/bookapplets/chapter10/Branch/Branch.html

Chapter 8: Markov Chains

it only matters where you are, not where you've been...

8.1 Introduction

So far, we have examined several stochastic processes using transition diagrams and First-Step Analysis.

The processes can be written as $\{X_0, X_1, X_2, \dots\}$, where X_t is the *state at time t* .

On the transition diagram, X_t corresponds to *which box we are in at step t* .



A.A. Markov
1856-1922

In the Gambler's Ruin (Section 2.7), X_t is the amount of money the gambler possesses after toss t . In the model for gene spread (Section 3.7), X_t is the number of animals possessing the harmful allele A in generation t .

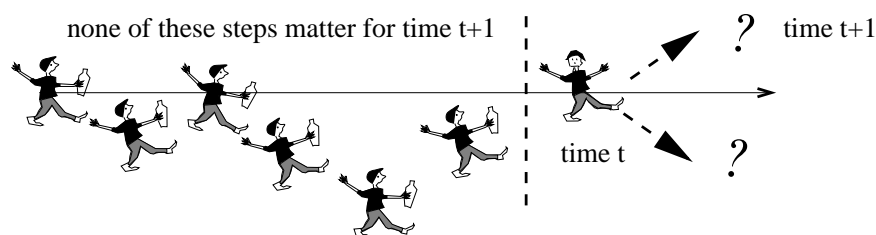
The processes that we have looked at via the transition diagram have a crucial property in common:

X_{t+1} *depends only on X_t* .

It does not depend upon X_0, X_1, \dots, X_{t-1} .

Processes like this are called *Markov Chains*.

Example: Random Walk (see Chapter 4)

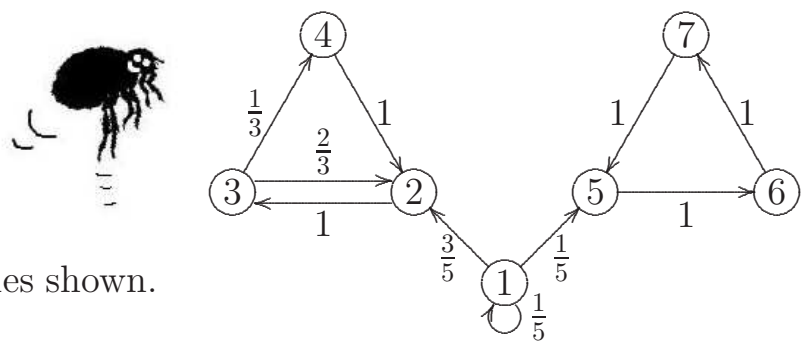


*In a Markov chain, the
future depends only
upon the present:
NOT upon the past.*

Meet... the *Markov fleas*!!



The text-book image of a Markov chain has a flea hopping about at random on the vertices of the transition diagram, according to the probabilities shown.



The transition diagram above shows a system with 7 possible states:

$$\text{state space } S = \{1, 2, 3, 4, 5, 6, 7\}.$$

Questions of interest

- Starting from state 1, what is the probability of ever reaching state 7?
- Starting from state 2, what is the expected time taken to reach state 4?
- Starting from state 2, what is the long-run proportion of time spent in state 3?
- Starting from state 1, what is the probability of being in state 2 at time t ? Does the probability converge as $t \rightarrow \infty$, and if so, to what?

We have been answering questions like the first two using first-step analysis since the start of STATS 325. In this chapter we develop a unified approach to all these questions using the matrix of transition probabilities, called the *transition matrix*.

8.2 Definitions

The Markov chain is the process X_0, X_1, X_2, \dots

Definition: The state of a Markov chain at time t is the *value of X_t* .

For example, if $X_t = 6$, we say *the process is in state 6 at time t* .

Definition: The state space of a Markov chain, S , is the set of values that each X_t can take. For example, $S = \{1, 2, 3, 4, 5, 6, 7\}$.

Let S have size N (possibly infinite).

Definition: A trajectory of a Markov chain is *a particular set of values for X_0, X_1, X_2, \dots*

For example, if $X_0 = 1$, $X_1 = 5$, and $X_2 = 6$, then the trajectory up to time $t = 2$ is 1, 5, 6.

More generally, if we refer to the trajectory $s_0, s_1, s_2, s_3, \dots$, we mean that $X_0 = s_0, X_1 = s_1, X_2 = s_2, X_3 = s_3, \dots$

‘Trajectory’ is just a word meaning ‘*path*’.

Markov Property

The basic property of a Markov chain is that *only the most recent point in the trajectory affects what happens next*.

This is called the *Markov Property*.

It means that X_{t+1} *depends upon X_t , but it does not depend upon X_{t-1}, \dots, X_1, X_0* .

We formulate the Markov Property in mathematical notation as follows:

$$\mathbb{P}(X_{t+1} = s \mid X_t = s_t, X_{t-1} = s_{t-1}, \dots, X_0 = s_0) = \mathbb{P}(X_{t+1} = s \mid X_t = s_t),$$

for all $t = 1, 2, 3, \dots$ and for all states s_0, s_1, \dots, s_t, s .

Explanation:

$$\begin{array}{ccccccc} \mathbb{P}(X_{t+1} = s & | & X_t = s_t, & \cancel{X_{t-1} = s_{t-1}, X_{t-2} = s_{t-2}, \dots, X_1 = s_1, X_0 = s_0}) \\ \uparrow & & \uparrow & \underbrace{\hspace{10em}} \\ \text{distribution} & & \text{depends} & \uparrow \\ \text{of } X_{t+1} & & \text{on } X_t & \text{but whatever happened before time } t \\ & & & \text{doesn't matter.} \end{array}$$

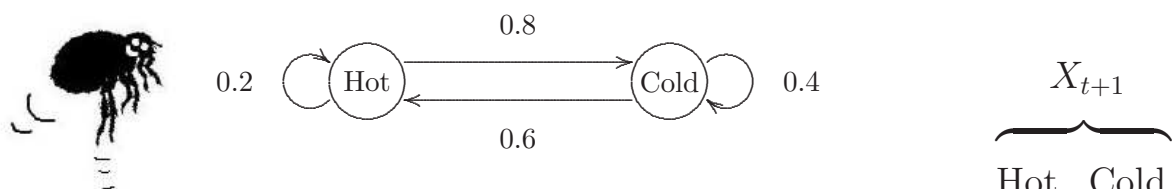
Definition: Let $\{X_0, X_1, X_2, \dots\}$ be a sequence of discrete random variables. Then $\{X_0, X_1, X_2, \dots\}$ is a Markov chain if it satisfies the Markov property:

$$\mathbb{P}(X_{t+1} = s \mid X_t = s_t, \dots, X_0 = s_0) = \mathbb{P}(X_{t+1} = s \mid X_t = s_t),$$

for all $t = 1, 2, 3, \dots$ and for all states s_0, s_1, \dots, s_t, s .

8.3 The Transition Matrix

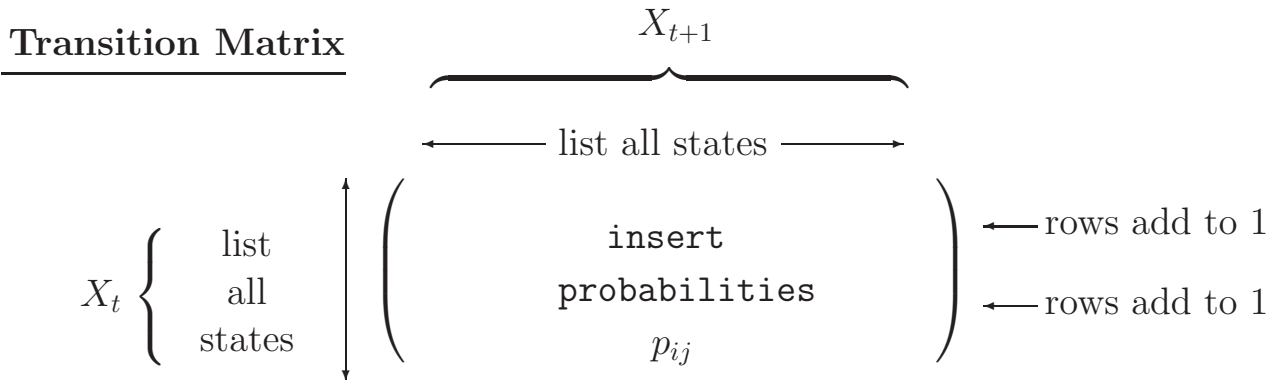
We have seen many examples of transition diagrams to describe Markov chains. The transition diagram is so-called because it shows the transitions between different states.



We can also summarize the probabilities in a matrix:

$$X_t \left\{ \begin{array}{l} \text{Hot} \\ \text{Cold} \end{array} \right. \left(\begin{array}{cc} 0.2 & 0.8 \\ 0.6 & 0.4 \end{array} \right)$$

The matrix describing the Markov chain is called the *transition matrix*. It is the most important tool for analysing Markov chains.



The transition matrix is usually given the symbol $P = (p_{ij})$.

In the transition matrix P :

- the ROWS represent NOW, or *FROM* (X_t);
- the COLUMNS represent NEXT, or *TO* (X_{t+1});
- entry (i, j) is the *CONDITIONAL* probability that *NEXT* = j , given that *NOW* = i : the probability of going *FROM* state i *TO* state j .

$$p_{ij} = \mathbb{P}(X_{t+1} = j \mid X_t = i).$$

- Notes:**
1. The transition matrix P must list *all* possible states in the state space S .
 2. P is a *square matrix* ($N \times N$), because X_{t+1} and X_t both take values in the same state space S (of size N).
 3. The rows of P should each *sum to 1*:

$$\sum_{j=1}^N p_{ij} = \sum_{j=1}^N \mathbb{P}(X_{t+1} = j \mid X_t = i) = \sum_{j=1}^N \mathbb{P}_{\{X_t = i\}}(X_{t+1} = j) = 1.$$

This simply states that X_{t+1} *must* take one of the listed values.

4. The columns of P do not in general sum to 1.

Definition: Let $\{X_0, X_1, X_2, \dots\}$ be a Markov chain with state space S , where S has size N (possibly infinite). The transition probabilities of the Markov chain are

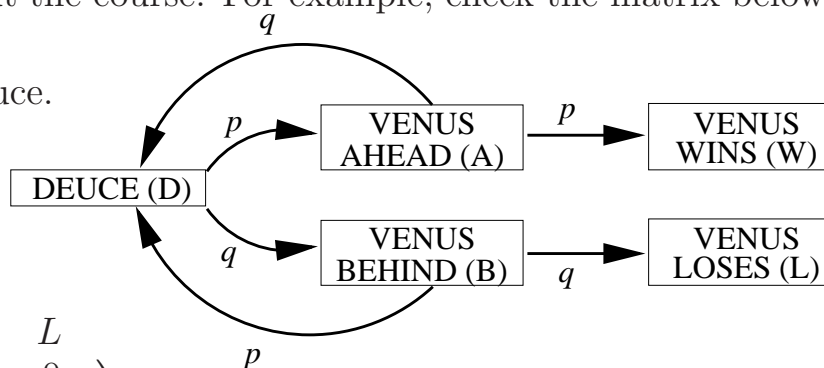
$$p_{ij} = \mathbb{P}(X_{t+1} = j \mid X_t = i) \quad \text{for } i, j \in S, \quad t = 0, 1, 2, \dots$$

Definition: The transition matrix of the Markov chain is $P = (p_{ij})$.

8.4 Example: setting up the transition matrix

We can create a transition matrix for any of the transition diagrams we have seen in problems throughout the course. For example, check the matrix below.

Example: Tennis game at Deuce.



$$\begin{array}{c} D \\ A \\ B \\ W \\ L \end{array}
 \begin{pmatrix}
 D & A & B & W & L \\
 \begin{pmatrix} 0 & p & q & 0 & 0 \\ q & 0 & 0 & p & 0 \\ p & 0 & 0 & 0 & q \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}
 \end{pmatrix}$$

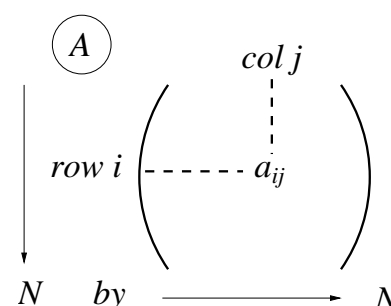
8.5 Matrix Revision

Notation

Let A be an $N \times N$ matrix.

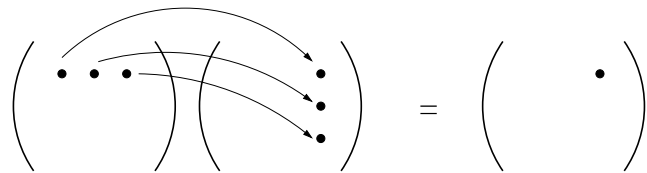
We write $A = (a_{ij})$,
i.e. A comprises elements a_{ij} .

The (i, j) element of A is written both as a_{ij} and $(A)_{ij}$:
e.g. for matrix A^2 we might write $(A^2)_{ij}$.



Matrix multiplication

Let $A = (a_{ij})$ and $B = (b_{ij})$
be $N \times N$ matrices.



$$\left(\begin{array}{ccc} \bullet & \bullet & \bullet \end{array} \right) \left(\begin{array}{ccc} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array} \right) = \left(\begin{array}{ccc} \bullet & \bullet & \bullet \end{array} \right)$$

The product matrix is $A \times B = AB$, with elements $(AB)_{ij} = \sum_{k=1}^N a_{ik}b_{kj}$.

Summation notation for a matrix squared

Let A be an $N \times N$ matrix. Then

$$(A^2)_{ij} = \sum_{k=1}^N (A)_{ik}(A)_{kj} = \sum_{k=1}^N a_{ik}a_{kj}.$$

Pre-multiplication of a matrix by a vector

Let A be an $N \times N$ matrix, and let $\boldsymbol{\pi}$ be an $N \times 1$ column vector: $\boldsymbol{\pi} = \begin{pmatrix} \pi_1 \\ \vdots \\ \pi_N \end{pmatrix}$.

We can pre-multiply A by $\boldsymbol{\pi}^T$ to get a $1 \times N$ row vector,
 $\boldsymbol{\pi}^T A = ((\boldsymbol{\pi}^T A)_1, \dots, (\boldsymbol{\pi}^T A)_N)$, with elements

$$(\boldsymbol{\pi}^T A)_j = \sum_{i=1}^N \pi_i a_{ij}.$$

8.6 The t -step transition probabilities

Let $\{X_0, X_1, X_2, \dots\}$ be a Markov chain with state space $S = \{1, 2, \dots, N\}$.

Recall that the elements of the transition matrix P are defined as:

$$(P)_{ij} = p_{ij} = \mathbb{P}(X_1 = j \mid X_0 = i) = \mathbb{P}(X_{n+1} = j \mid X_n = i) \quad \text{for any } n.$$

p_{ij} is the probability of making a transition FROM state i TO state j in a SINGLE step.

Question: what is the probability of making a transition from state i to state j over two steps? *I.e. what is $\mathbb{P}(X_2 = j \mid X_0 = i)$?*

We are seeking $\mathbb{P}(X_2 = j \mid X_0 = i)$. Use the *Partition Theorem*:

$$\begin{aligned}
 \mathbb{P}(X_2 = j \mid X_0 = i) &= \mathbb{P}_i(X_2 = j) \quad (\text{notation of Ch 2}) \\
 &= \sum_{k=1}^N \mathbb{P}_i(X_2 = j \mid X_1 = k) \mathbb{P}_i(X_1 = k) \quad (\text{Partition Thm}) \\
 &= \sum_{k=1}^N \mathbb{P}(X_2 = j \mid X_1 = k, X_0 = i) \mathbb{P}(X_1 = k \mid X_0 = i) \\
 &= \sum_{k=1}^N \mathbb{P}(X_2 = j \mid X_1 = k) \mathbb{P}(X_1 = k \mid X_0 = i) \\
 &\hspace{15em} (\text{Markov Property}) \\
 &= \sum_{k=1}^N p_{kj} p_{ik} \quad (\text{by definitions}) \\
 &= \sum_{k=1}^N p_{ik} p_{kj} \quad (\text{rearranging}) \\
 &= (P^2)_{ij}. \quad (\text{see Matrix Revision})
 \end{aligned}$$

The two-step transition probabilities are therefore given by *the matrix* P^2 :

$$\mathbb{P}(X_2 = j \mid X_0 = i) = \mathbb{P}(X_{n+2} = j \mid X_n = i) = (P^2)_{ij} \quad \text{for any } n.$$

3-step transitions: We can find $\mathbb{P}(X_3 = j \mid X_0 = i)$ similarly, but conditioning on the state at time 2:

$$\begin{aligned}
 \mathbb{P}(X_3 = j \mid X_0 = i) &= \sum_{k=1}^N \mathbb{P}(X_3 = j \mid X_2 = k) \mathbb{P}(X_2 = k \mid X_0 = i) \\
 &= \sum_{k=1}^N p_{kj} (P^2)_{ik} \\
 &= (P^3)_{ij}.
 \end{aligned}$$

The three-step transition probabilities are therefore given by the matrix P^3 :

$$\mathbb{P}(X_3 = j \mid X_0 = i) = \mathbb{P}(X_{n+3} = j \mid X_n = i) = (P^3)_{ij} \quad \text{for any } n.$$

General case: t -step transitions

The above working extends to show that the t -step transition probabilities are given by the matrix P^t for any t :

$$\mathbb{P}(X_t = j \mid X_0 = i) = \mathbb{P}(X_{n+t} = j \mid X_n = i) = (P^t)_{ij} \quad \text{for any } n.$$

We have proved the following Theorem.

Theorem 8.6: Let $\{X_0, X_1, X_2, \dots\}$ be a Markov chain with $N \times N$ transition matrix P . Then the t -step transition probabilities are given by the matrix P^t . That is,

$$\mathbb{P}(X_t = j \mid X_0 = i) = (P^t)_{ij}.$$

It also follows that

$$\mathbb{P}(X_{n+t} = j \mid X_n = i) = (P^t)_{ij} \quad \text{for any } n. \quad \square$$

8.7 Distribution of X_t

Let $\{X_0, X_1, X_2, \dots\}$ be a Markov chain with state space $S = \{1, 2, \dots, N\}$.

Now each X_t is a random variable, so it has a *probability distribution*.

We can write the probability distribution of X_t as an $N \times 1$ *vector*.

For example, consider X_0 . Let $\boldsymbol{\pi}$ be an $N \times 1$ vector denoting the probability distribution of X_0 :

$$\boldsymbol{\pi} = \begin{pmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_N \end{pmatrix} = \begin{pmatrix} \mathbb{P}(X_0 = 1) \\ \mathbb{P}(X_0 = 2) \\ \vdots \\ \mathbb{P}(X_0 = N) \end{pmatrix}$$

In the flea model, this corresponds to *the flea choosing at random which vertex it starts off from at time 0, such that*

$$\mathbb{P}(\text{flea chooses vertex } i \text{ to start}) = \pi_i.$$

Notation: we will write $X_0 \sim \boldsymbol{\pi}^T$ to denote that the row vector of probabilities is given by the row vector $\boldsymbol{\pi}^T$.

Probability distribution of X_1

Use the Partition Rule, conditioning on X_0 :

$$\begin{aligned} \mathbb{P}(X_1 = j) &= \sum_{i=1}^N \mathbb{P}(X_1 = j \mid X_0 = i) \mathbb{P}(X_0 = i) \\ &= \sum_{i=1}^N p_{ij} \pi_i \quad \text{by definitions} \\ &= \sum_{i=1}^N \pi_i p_{ij} \\ &= (\boldsymbol{\pi}^T P)_j. \end{aligned}$$

(pre-multiplication by a vector from Section 8.5).

This shows that $\mathbb{P}(X_1 = j) = (\boldsymbol{\pi}^T P)_j$ for all j .

The row vector $\boldsymbol{\pi}^T P$ is therefore *the probability distribution of X_1* :

$\begin{aligned} X_0 &\sim \boldsymbol{\pi}^T \\ X_1 &\sim \boldsymbol{\pi}^T P. \end{aligned}$

Probability distribution of X_2

Using the Partition Rule as before, conditioning again on X_0 :

$$\mathbb{P}(X_2 = j) = \sum_{i=1}^N \mathbb{P}(X_2 = j \mid X_0 = i) \mathbb{P}(X_0 = i) = \sum_{i=1}^N (P^2)_{ij} \pi_i = (\boldsymbol{\pi}^T P^2)_j.$$

The row vector $\pi^T P^2$ is therefore the probability distribution of X_2 :

$$\begin{array}{l} X_0 \sim \pi^T \\ X_1 \sim \pi^T P \\ X_2 \sim \pi^T P^2 \\ \vdots \\ X_t \sim \pi^T P^t. \end{array}$$

These results are summarized in the following Theorem.

Theorem 8.7: Let $\{X_0, X_1, X_2, \dots\}$ be a Markov chain with $N \times N$ transition matrix P . If the probability distribution of X_0 is given by the $1 \times N$ row vector π^T , then the probability distribution of X_t is given by the $1 \times N$ row vector $\pi^T P^t$. That is,

$$X_0 \sim \pi^T \Rightarrow X_t \sim \pi^T P^t.$$

Note: The distribution of X_t is $X_t \sim \pi^T P^t$.

The distribution of X_{t+1} is $X_{t+1} \sim \pi^T P^{t+1}$.

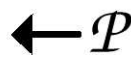
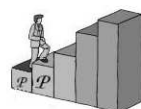
Taking one step in the Markov chain corresponds to *multiplying by P on the right*.

Note: The t -step transition matrix is P^t (Theorem 8.6).

The $(t+1)$ -step transition matrix is P^{t+1} .

Again, taking one step in the Markov chain corresponds to *multiplying by P on the right*.

take 1 step...

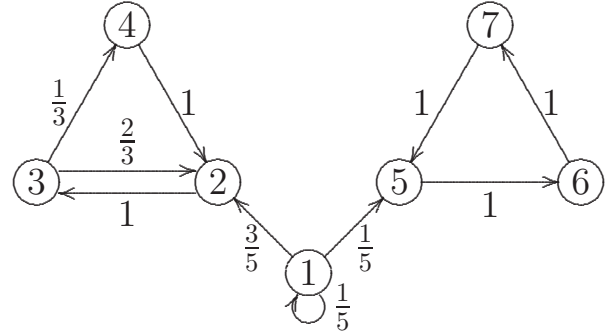


\equiv ...multiply by P
on the right

8.8 Trajectory Probability

Recall that a trajectory is a sequence of values for X_0, X_1, \dots, X_t .

Because of the Markov Property, we can find the probability of any trajectory by multiplying together the starting probability and all subsequent single-step probabilities.



Example: Let $X_0 \sim (\frac{3}{4}, 0, \frac{1}{4}, 0, 0, 0, 0)$. What is the probability of the trajectory 1, 2, 3, 2, 3, 4?

$$\begin{aligned} \mathbb{P}(1, 2, 3, 2, 3, 4) &= \mathbb{P}(X_0 = 1) \times p_{12} \times p_{23} \times p_{32} \times p_{23} \times p_{34} \\ &= \frac{3}{4} \times \frac{3}{5} \times 1 \times \frac{2}{3} \times 1 \times \frac{1}{3} \\ &= \frac{1}{10}. \end{aligned}$$

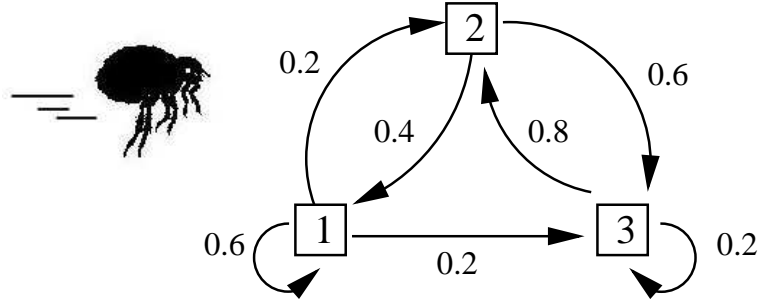
Proof in formal notation using the Markov Property:

Let $X_0 \sim \boldsymbol{\pi}^T$. We wish to find the probability of the trajectory $s_0, s_1, s_2, \dots, s_t$.

$$\begin{aligned} &\mathbb{P}(X_0 = s_0, X_1 = s_1, \dots, X_t = s_t) \\ &= \mathbb{P}(X_t = s_t \mid X_{t-1} = s_{t-1}, \dots, X_0 = s_0) \times \mathbb{P}(X_{t-1} = s_{t-1}, \dots, X_0 = s_0) \\ &= \mathbb{P}(X_t = s_t \mid X_{t-1} = s_{t-1}) \times \mathbb{P}(X_{t-1} = s_{t-1}, \dots, X_0 = s_0) \quad (\text{Markov Property}) \\ &= p_{s_{t-1}, s_t} \mathbb{P}(X_{t-1} = s_{t-1} \mid X_{t-2} = s_{t-2}, \dots, X_0 = s_0) \times \mathbb{P}(X_{t-2} = s_{t-2}, \dots, X_0 = s_0) \\ &\quad \vdots \\ &= p_{s_{t-1}, s_t} \times p_{s_{t-2}, s_{t-1}} \times \dots \times p_{s_0, s_1} \times \mathbb{P}(X_0 = s_0) \\ &= p_{s_{t-1}, s_t} \times p_{s_{t-2}, s_{t-1}} \times \dots \times p_{s_0, s_1} \times \pi_{s_0}. \end{aligned}$$

8.9 Worked Example: distribution of X_t and trajectory probabilities

Purpose-flea zooms around the vertices of the transition diagram opposite. Let X_t be Purpose-flea's state at time t ($t = 0, 1, \dots$).



- (a) Find the transition matrix, P .

Answer: $P = \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.4 & 0 & 0.6 \\ 0 & 0.8 & 0.2 \end{pmatrix}$

- (b) Find $\mathbb{P}(X_2 = 3 \mid X_0 = 1)$.

$$\begin{aligned} \mathbb{P}(X_2 = 3 \mid X_0 = 1) &= (P^2)_{13} = \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix} \begin{pmatrix} \cdot & \cdot & 0.2 \\ \cdot & \cdot & 0.6 \\ \cdot & \cdot & 0.2 \end{pmatrix} \\ &= 0.6 \times 0.2 + 0.2 \times 0.6 + 0.2 \times 0.2 \\ &= 0.28. \end{aligned}$$

Note: we only need one element of the matrix P^2 , so don't lose exam time by finding the whole matrix.

- (c) Suppose that Purpose-flea is equally likely to start on any vertex at time 0. Find the probability distribution of X_1 .

From this info, the distribution of X_0 is $\pi^T = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. We need $X_1 \sim \pi^T P$.

$$\pi^T P = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.4 & 0 & 0.6 \\ 0 & 0.8 & 0.2 \end{pmatrix} = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix}.$$

Thus $X_1 \sim (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and therefore X_1 is also equally likely to be 1, 2, or 3.

- (d) Suppose that Purpose-flea begins at vertex 1 at time 0. Find the probability distribution of X_2 .

The distribution of X_0 is now $\pi^T = (1, 0, 0)$. We need $X_2 \sim \pi^T P^2$.

$$\begin{aligned}\pi^T P^2 &= (1 \ 0 \ 0) \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.4 & 0 & 0.6 \\ 0 & 0.8 & 0.2 \end{pmatrix} \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.4 & 0 & 0.6 \\ 0 & 0.8 & 0.2 \end{pmatrix} \\ &= (0.6 \ 0.2 \ 0.2) \begin{pmatrix} 0.6 & 0.2 & 0.2 \\ 0.4 & 0 & 0.6 \\ 0 & 0.8 & 0.2 \end{pmatrix} \\ &= (0.44 \ 0.28 \ 0.28).\end{aligned}$$

Thus $\mathbb{P}(X_2 = 1) = 0.44$, $\mathbb{P}(X_2 = 2) = 0.28$, $\mathbb{P}(X_2 = 3) = 0.28$.

Note that it is quickest to multiply the vector by the matrix first: we don't need to compute P^2 in entirety.

- (e) Suppose that Purpose-flea is equally likely to start on any vertex at time 0. Find the probability of obtaining the trajectory $(3, 2, 1, 1, 3)$.

$$\begin{aligned}\mathbb{P}(3, 2, 1, 1, 3) &= \mathbb{P}(X_0 = 3) \times p_{32} \times p_{21} \times p_{11} \times p_{13} \quad (\text{Section 8.8}) \\ &= \frac{1}{3} \times 0.8 \times 0.4 \times 0.6 \times 0.2 \\ &= 0.0128.\end{aligned}$$

8.10 Class Structure

The state space of a Markov chain can be partitioned into a set of non-overlapping *communicating classes*.

States i and j are in the same communicating class if there is some way of getting from state i to state j , AND there is some way of getting from state j to state i . It needn't be possible to get between i and j in a **single** step, but it must be possible over some number of steps to travel between them both ways.

We write $i \leftrightarrow j$.

Definition: Consider a Markov chain with state space S and transition matrix P , and consider states $i, j \in S$. Then state i communicates with state j if:

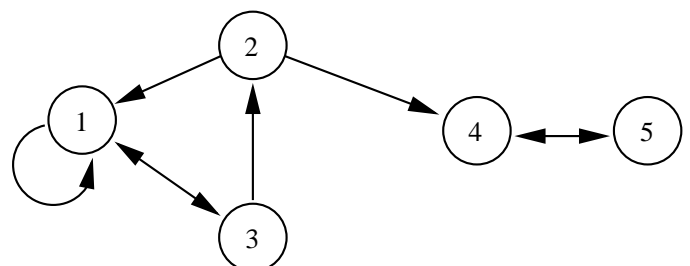
1. there exists some t such that $(P^t)_{ij} > 0$, AND
2. there exists some u such that $(P^u)_{ji} > 0$.

Mathematically, it is easy to show that the communicating relation \leftrightarrow is an equivalence relation, which means that it partitions the sample space S into non-overlapping equivalence classes.

Definition: States i and j are in the same communicating class if $i \leftrightarrow j$: i.e. if each state is accessible from the other.

Every state is a member of *exactly one communicating class*.

Example: Find the communicating classes associated with the transition diagram shown.



Solution:

$\{1, 2, 3\}, \quad \{4, 5\}.$

State 2 leads to state 4, but state 4 does not lead back to state 2, so they are in different communicating classes.

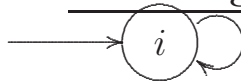
Definition: A communicating class of states is closed if *it is not possible to leave that class*.

That is, the communicating class C is closed if $p_{ij} = 0$ whenever $i \in C$ and $j \notin C$.

Example: In the transition diagram above:

- Class $\{1, 2, 3\}$ is not closed: *it is possible to escape to class $\{4, 5\}$.*
- Class $\{4, 5\}$ is *closed: it is not possible to escape.*

Definition: A state i is said to be absorbing if *the set $\{i\}$ is a closed class*.



Definition: A Markov chain or transition matrix P is said to be irreducible if $i \leftrightarrow j$ for all $i, j \in S$. *That is, the chain is irreducible if the state space S is a single communicating class.*

8.11 Hitting Probabilities

We have been calculating hitting probabilities for Markov chains since Chapter 2, using First-Step Analysis. The hitting probability describes the probability that the Markov chain will *ever* reach some state or set of states.

In this section we show how hitting probabilities can be written in a single vector. We also see a general formula for calculating the hitting probabilities. In general it is easier to continue using our own common sense, but occasionally the formula becomes more necessary.



Vector of hitting probabilities

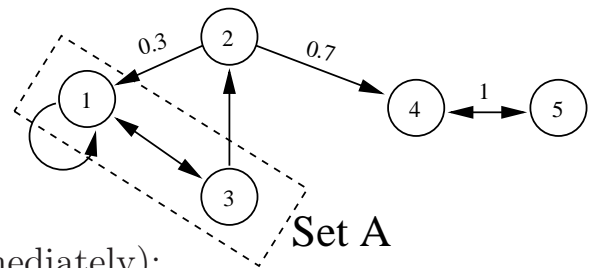
Let A be some subset of the state space S . (A need not be a communicating class: it can be any subset required, including the subset consisting of a single state: e.g. $A = \{4\}$.)

The **hitting probability** from state i to set A is the probability of ever reaching the set A , starting from initial state i . We write this probability as h_{iA} . Thus

$$h_{iA} = \mathbb{P}(X_t \in A \text{ for some } t \geq 0 \mid X_0 = i).$$

Example: Let set $A = \{1, 3\}$ as shown.

The hitting probability for set A is:



- *1 starting from states 1 or 3*
(We are starting in set A , so we hit it immediately);
- *0 starting from states 4 or 5*
(The set $\{4, 5\}$ is a closed class, so we can never escape out to set A);
- *0.3 starting from state 2*
(We could hit A at the first step (probability 0.3), but otherwise we move to state 4 and get stuck in the closed class $\{4, 5\}$ (probability 0.7).)

We can summarize all the information from the example above in a **vector of hitting probabilities**:

$$\mathbf{h}_A = \begin{pmatrix} h_{1A} \\ h_{2A} \\ h_{3A} \\ h_{4A} \\ h_{5A} \end{pmatrix} = \begin{pmatrix} 1 \\ 0.3 \\ 1 \\ 0 \\ 0 \end{pmatrix}.$$

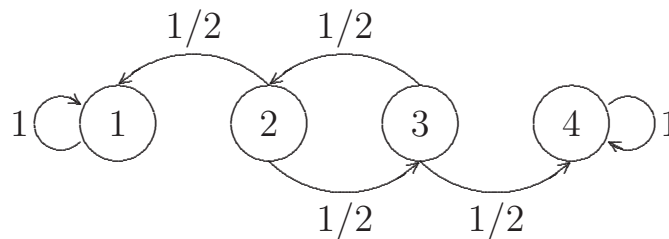
Note: When A is a closed class, the hitting probability h_{iA} is called the **absorption probability**.

In general, if there are N possible states, the vector of hitting probabilities is

$$\mathbf{h}_A = \begin{pmatrix} h_{1A} \\ h_{2A} \\ \vdots \\ h_{NA} \end{pmatrix} = \begin{pmatrix} \mathbb{P}(\text{hit } A \text{ starting from state 1}) \\ \mathbb{P}(\text{hit } A \text{ starting from state 2}) \\ \vdots \\ \mathbb{P}(\text{hit } A \text{ starting from state } N) \end{pmatrix}.$$

Example: finding the hitting probability vector using First-Step Analysis

Suppose $\{X_t : t \geq 0\}$ has the following transition diagram:



Find the vector of hitting probabilities for state 4.

Solution:

Let $h_{i4} = \mathbb{P}(\text{hit state 4, starting from state } i)$. Clearly,

$$h_{14} = 0$$

$$h_{44} = 1$$

Using first-step analysis, we also have:

$$h_{24} = \frac{1}{2}h_{34} + \frac{1}{2} \times 0$$

$$h_{34} = \frac{1}{2} + \frac{1}{2}h_{24}$$

Solving,

$$h_{34} = \frac{1}{2} + \frac{1}{2} \left(\frac{1}{2}h_{34} \right) \Rightarrow h_{34} = \frac{2}{3}. \quad \text{So also, } h_{24} = \frac{1}{2}h_{34} = \frac{1}{3}.$$

So the vector of hitting probabilities is

$$\mathbf{h}_A = \left(0, \frac{1}{3}, \frac{2}{3}, 1 \right).$$

Formula for hitting probabilities

In the previous example, we used our common sense to state that $h_{14} = 0$. While this is easy for a human brain, it is harder to explain a general rule that would describe this ‘common sense’ mathematically, or that could be used to write computer code that will work for all problems.

Although it is usually best to continue to use common sense when solving problems, this section provides a general formula that will *always* work to find a vector of hitting probabilities \mathbf{h}_A .

Theorem 8.11: The vector of hitting probabilities $\mathbf{h}_A = (h_{iA} : i \in S)$ is the minimal non-negative solution to the following equations:

$$h_{iA} = \begin{cases} 1 & \text{for } i \in A, \\ \sum_{j \in S} p_{ij} h_{jA} & \text{for } i \notin A. \end{cases}$$

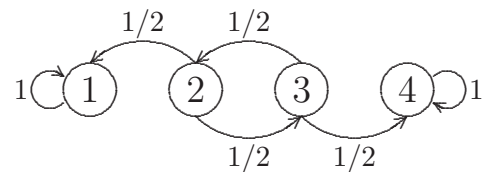
The ‘minimal non-negative solution’ means that:

1. the values $\{h_{iA}\}$ collectively satisfy the equations above;
2. each value h_{iA} is ≥ 0 (non-negative);
3. given any other non-negative solution to the equations above, say $\{g_{iA}\}$ where $g_{iA} \geq 0$ for all i , then $h_{iA} \leq g_{iA}$ for all i (minimal solution).

Example: How would this formula be used to substitute for ‘common sense’ in the previous example?

The equations give:

$$h_{i4} = \begin{cases} 1 & \text{if } i = 4, \\ \sum_{j \in S} p_{ij} h_{j4} & \text{if } i \neq 4. \end{cases}$$



Thus,

$$h_{44} = 1$$

$$h_{14} = h_{14} \quad \text{unspecified! Could be anything!}$$

$$h_{24} = \frac{1}{2}h_{14} + \frac{1}{2}h_{34}$$

$$h_{34} = \frac{1}{2}h_{24} + \frac{1}{2}h_{44} = \frac{1}{2}h_{24} + \frac{1}{2}$$

Because h_{14} could be anything, we have to use the minimal non-negative value, which is $h_{14} = 0$.

(Need to check $h_{14} = 0$ does not force $h_{i4} < 0$ for any other i : OK.)

The other equations can then be solved to give the same answers as before. \square

Proof of Theorem 8.11 (non-examinable):

Consider the equations
$$h_{iA} = \begin{cases} 1 & \text{for } i \in A, \\ \sum_{j \in S} p_{ij} h_{jA} & \text{for } i \notin A. \end{cases} \quad (\star)$$

We need to show that:

- (i) the hitting probabilities $\{h_{iA}\}$ collectively satisfy the equations (\star) ;
- (ii) if $\{g_{iA}\}$ is any other non-negative solution to (\star) , then the hitting probabilities $\{h_{iA}\}$ satisfy $h_{iA} \leq g_{iA}$ for all i (minimal solution).

Proof of (i): Clearly, $h_{iA} = 1$ if $i \in A$ (as the chain hits A immediately).

Suppose that $i \notin A$. Then

$$\begin{aligned} h_{iA} &= \mathbb{P}(X_t \in A \text{ for some } t \geq 1 \mid X_0 = i) \\ &= \sum_{j \in S} \mathbb{P}(X_t \in A \text{ for some } t \geq 1 \mid X_1 = j) \mathbb{P}(X_1 = j \mid X_0 = i) \\ &\hspace{15em} (\text{Partition Rule}) \\ &= \sum_{j \in S} h_{jA} p_{ij} \quad (\text{by definitions}). \end{aligned}$$

Thus the hitting probabilities $\{h_{iA}\}$ must satisfy the equations (\star) .

Proof of (ii): Let $h_{iA}^{(t)} = \mathbb{P}(\text{hit } A \text{ at or before time } t \mid X_0 = i)$.

We use mathematical induction to show that $h_{iA}^{(t)} \leq g_{iA}$ for all t , and therefore $h_{iA} = \lim_{t \rightarrow \infty} h_{iA}^{(t)}$ must also be $\leq g_{iA}$.

Time $t = 0$:
$$h_{iA}^{(0)} = \begin{cases} 1 & \text{if } i \in A, \\ 0 & \text{if } i \notin A. \end{cases}$$

But because g_{iA} is non-negative and satisfies (\star) ,
$$\begin{cases} g_{iA} = 1 & \text{if } i \in A, \\ g_{iA} \geq 0 & \text{for all } i. \end{cases}$$

So $g_{iA} \geq h_{iA}^{(0)}$ for all i .

The inductive hypothesis is true for time $t = 0$.

Time t : Suppose the inductive hypothesis holds for time t , i.e.

$$h_{jA}^{(t)} \leq g_{jA} \quad \text{for all } j.$$

Consider

$$\begin{aligned} h_{iA}^{(t+1)} &= \mathbb{P}(\text{hit } A \text{ by time } t+1 \mid X_0 = i) \\ &= \sum_{j \in S} \mathbb{P}(\text{hit } A \text{ by time } t+1 \mid X_1 = j) \mathbb{P}(X_1 = j \mid X_0 = i) \\ &\hspace{25em} (\text{Partition Rule}) \\ &= \sum_{j \in S} h_{jA}^{(t)} p_{ij} \quad \text{by definitions} \\ &\leq \sum_{j \in S} g_{jA} p_{ij} \quad \text{by inductive hypothesis} \\ &= g_{iA} \quad \text{because } \{g_{iA}\} \text{ satisfies } (\star). \end{aligned}$$

Thus $h_{iA}^{(t+1)} \leq g_{iA}$ for all i , so the inductive hypothesis is proved.

By the Continuity Theorem (Chapter 2), $h_{iA} = \lim_{t \rightarrow \infty} h_{iA}^{(t)}$.

So $h_{iA} \leq g_{iA}$ as required. □

8.12 Expected hitting times

In the previous section we found the **probability** of hitting set A , starting at state i . Now we study **how long** it takes to get from i to A . As before, it is best to solve problems using first-step analysis and common sense. However, a general formula is also available.



Definition: Let A be a subset of the state space S . The **hitting time** of A is the random variable T_A , where

$$T_A = \min\{t \geq 0 : X_t \in A\}.$$

T_A is the time taken before hitting set A *for the first time*.

The hitting time T_A can take values $0, 1, 2, \dots$, and ∞ .

If the chain *never* hits set A , then $T_A = \infty$.

Note: The hitting time is also called the **reaching time**. If A is a closed class, it is also called the **absorption time**.

Definition: The **mean hitting time** for A , starting from state i , is

$$m_{iA} = \mathbb{E}(T_A | X_0 = i).$$

Note: If there is any possibility that the chain *never* reaches A , starting from i , i.e. if the hitting probability $h_{iA} < 1$, then $\mathbb{E}(T_A | X_0 = i) = \infty$.

Calculating the mean hitting times

Theorem 8.12: The vector of expected hitting times $\mathbf{m}_A = (m_{iA} : i \in S)$ is *the minimal non-negative solution to the following equations:*

$$m_{iA} = \begin{cases} 0 & \text{for } i \in A, \\ 1 + \sum_{j \notin A} p_{ij} m_{jA} & \text{for } i \notin A. \end{cases}$$

Proof (sketch):

$$\text{Consider the equations } m_{iA} = \begin{cases} 0 & \text{for } i \in A, \\ 1 + \sum_{j \notin A} p_{ij} m_{jA} & \text{for } i \notin A. \end{cases} \quad (\star).$$

We need to show that:

- (i) the mean hitting times $\{m_{iA}\}$ collectively satisfy the equations (\star) ;
- (ii) if $\{u_{iA}\}$ is any other non-negative solution to (\star) , then the mean hitting times $\{m_{iA}\}$ satisfy $m_{iA} \leq u_{iA}$ for all i (minimal solution).

We will prove point (i) only. A proof of (ii) can be found online at:
<http://www.statslab.cam.ac.uk/~james/Markov/> , Section 1.3.

Proof of (i): Clearly, $m_{iA} = 0$ if $i \in A$ (as the chain hits A immediately).

Suppose that $i \notin A$. Then

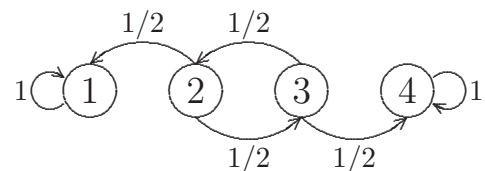
$$\begin{aligned} m_{iA} &= \mathbb{E}(T_A | X_0 = i) \\ &= 1 + \sum_{j \in S} \mathbb{E}(T_A | X_1 = j) \mathbb{P}(X_1 = j | X_0 = i) \\ &\quad \text{(conditional expectation: take 1 step to get to state } j \\ &\quad \text{at time 1, then find } \mathbb{E}(T_A) \text{ from there)} \\ &= 1 + \sum_{j \in S} m_{jA} p_{ij} \quad \text{(by definitions)} \\ &= 1 + \sum_{j \notin A} p_{ij} m_{jA}, \quad \text{because } m_{jA} = 0 \text{ for } j \in A. \end{aligned}$$

Thus the mean hitting times $\{m_{iA}\}$ must satisfy the equations (\star) .

□

Example: Let $\{X_t : t \geq 0\}$ have the same transition diagram as before:

Starting from state 2, find the expected time to absorption.



Solution:

Starting from state $i = 2$, we wish to find the expected time to reach the set $A = \{1, 4\}$ (the set of absorbing states).

Thus we are looking for $m_{iA} = m_{2A}$.

$$\text{Now } m_{iA} = \begin{cases} 0 & \text{if } i \in \{1, 4\}, \\ 1 + \sum_{j \notin A} p_{ij} m_{jA} & \text{if } i \notin \{1, 4\}. \end{cases}$$

Thus,

$$m_{1A} = 0 \quad (\text{because } 1 \in A)$$

$$m_{4A} = 0 \quad (\text{because } 4 \in A)$$

$$m_{2A} = 1 + \frac{1}{2}m_{1A} + \frac{1}{2}m_{3A}$$

$$\Rightarrow m_{2A} = 1 + \frac{1}{2}m_{3A}$$

$$m_{3A} = 1 + \frac{1}{2}m_{2A} + \frac{1}{2}m_{4A}$$

$$= 1 + \frac{1}{2}m_{2A}$$

$$= 1 + \frac{1}{2} \left(1 + \frac{1}{2}m_{3A} \right)$$

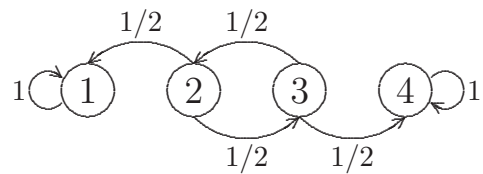
$$\Rightarrow \frac{3}{4}m_{3A} = \frac{3}{2}$$

$$\Rightarrow m_{3A} = 2.$$

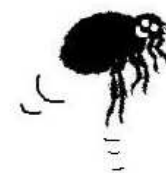
Thus,

$$m_{2A} = 1 + \frac{1}{2}m_{3A} = 2.$$

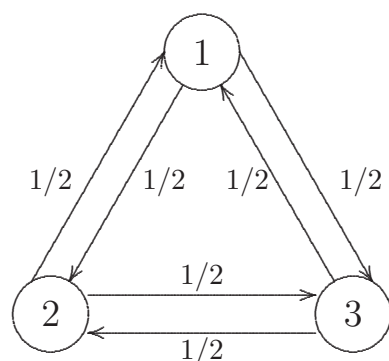
The expected time to absorption is therefore $\mathbb{E}(T_A) = 2$ steps.



Example: Glee-flea hops around on a triangle. At each step he moves to one of the other two vertices at random. What is the expected time taken for Glee-flea to get from vertex 1 to vertex 2?



Solution:



transition matrix, $P = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}.$

We wish to find m_{12} .

$$\text{Now } m_{i2} = \begin{cases} 0 & \text{if } i = 2, \\ 1 + \sum_{j \neq 2} p_{ij} m_{j2} & \text{if } i \neq 2. \end{cases}$$

Thus

$$m_{22} = 0$$

$$m_{12} = 1 + \frac{1}{2}m_{22} + \frac{1}{2}m_{32} = 1 + \frac{1}{2}m_{32}.$$

$$m_{32} = 1 + \frac{1}{2}m_{22} + \frac{1}{2}m_{12}$$

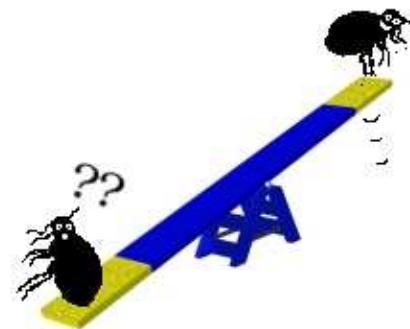
$$= 1 + \frac{1}{2}m_{12}$$

$$= 1 + \frac{1}{2} \left(1 + \frac{1}{2}m_{32} \right)$$

$$\Rightarrow m_{32} = 2.$$

Thus $m_{12} = 1 + \frac{1}{2}m_{32} = 2$ steps.

Chapter 9: Equilibrium



In Chapter 8, we saw that if $\{X_0, X_1, X_2, \dots\}$ is a Markov chain with transition matrix P , then

$$X_t \sim \pi^T \Rightarrow X_{t+1} \sim \pi^T P.$$

This raises the question: is there any distribution π such that $\pi^T P = \pi^T$?

If $\pi^T P = \pi^T$, then

$$\begin{aligned} X_t \sim \pi^T &\Rightarrow X_{t+1} \sim \pi^T P = \pi^T \\ &\Rightarrow X_{t+2} \sim \pi^T P = \pi^T \\ &\Rightarrow X_{t+3} \sim \pi^T P = \pi^T \\ &\Rightarrow \dots \end{aligned}$$

In other words, if $\pi^T P = \pi^T$, and $X_t \sim \pi^T$, then

$$X_t \sim X_{t+1} \sim X_{t+2} \sim X_{t+3} \sim \dots$$

Thus, once a Markov chain has reached a distribution π^T such that $\pi^T P = \pi^T$, *it will stay there*.

If $\pi^T P = \pi^T$, we say that the distribution π^T is an *equilibrium distribution*.

Equilibrium means a *level position*: there is *no more change* in the distribution of X_t as we wander through the Markov chain.

Note: Equilibrium does not mean that the value of X_{t+1} equals the value of X_t . It means that the distribution of X_{t+1} is the same as the distribution of X_t :

e.g. $\mathbb{P}(X_{t+1} = 1) = \mathbb{P}(X_t = 1) = \pi_1;$

$$\mathbb{P}(X_{t+1} = 2) = \mathbb{P}(X_t = 2) = \pi_2, \quad \text{etc.}$$

In this chapter, we will first see how to *calculate* the equilibrium distribution π^T . We will then see the remarkable result that many Markov chains automatically *find their own way* to an equilibrium distribution as the chain wanders through time. This happens for many Markov chains, but not all. We will see the conditions required for the chain to find its way to an equilibrium distribution.

9.1 Equilibrium distribution in pictures

Consider the following 4-state Markov chain:

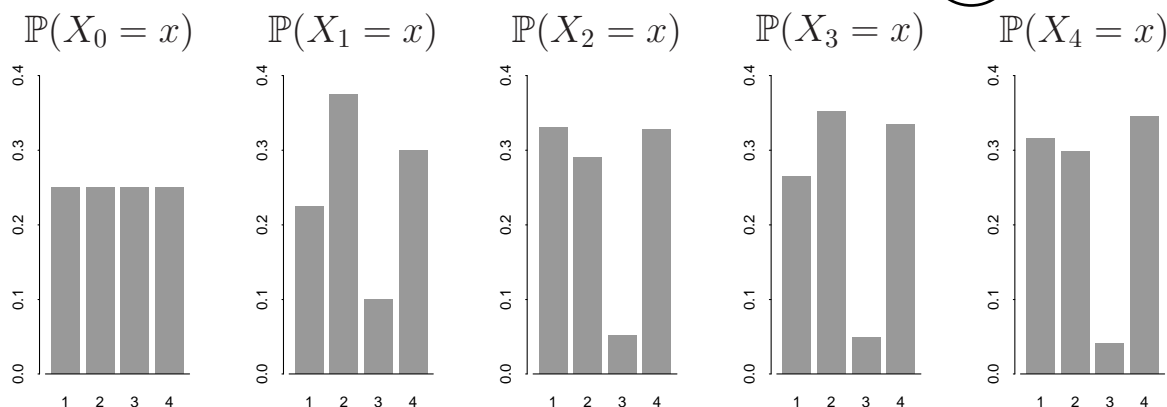
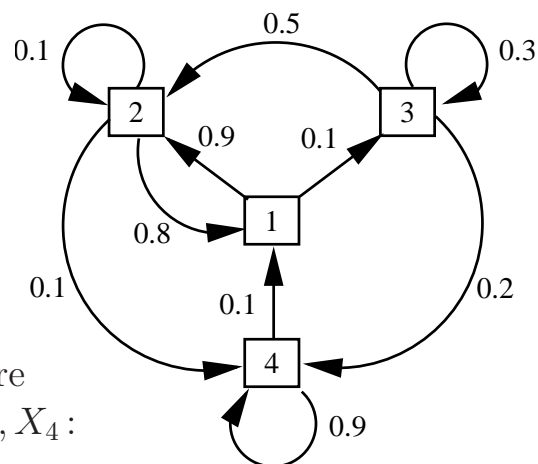
$$P = \begin{pmatrix} 0.0 & 0.9 & 0.1 & 0.0 \\ 0.8 & 0.1 & 0.0 & 0.1 \\ 0.0 & 0.5 & 0.3 & 0.2 \\ 0.1 & 0.0 & 0.0 & 0.9 \end{pmatrix}$$

Suppose we start at time 0 with

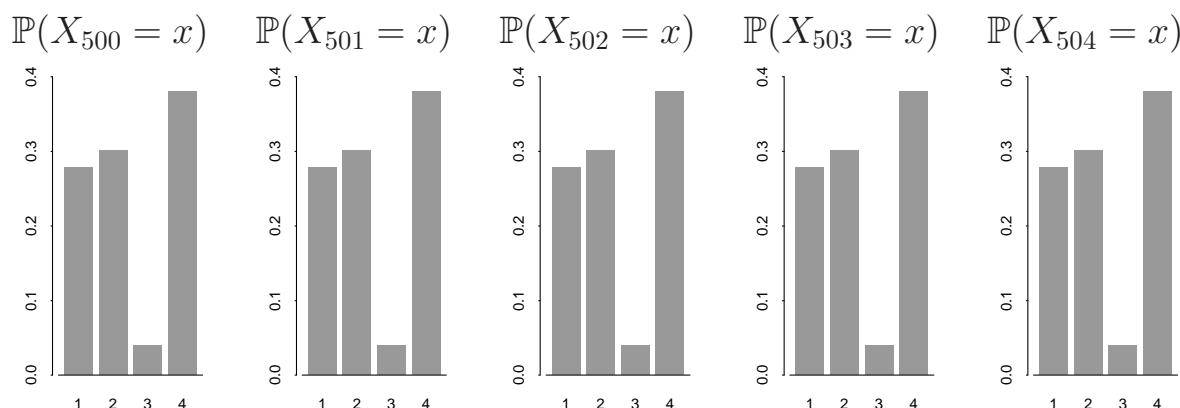
$X_0 \sim (\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$: so the chain is equally

likely to start from any of the four states. Here

are pictures of the distributions of X_0, X_1, \dots, X_4 :



The distribution starts off level, but quickly changes: for example the chain is least likely to be found in state 3. The distribution of X_t changes between each $t = 0, 1, 2, 3, 4$. Now look at the distribution of X_t 500 steps into the future:



The distribution has reached a steady state: it **does not change** between $t = 500, 501, \dots, 504$. *The chain has reached equilibrium of its own accord.*

9.2 Calculating equilibrium distributions

Definition: Let $\{X_0, X_1, \dots\}$ be a Markov chain with transition matrix P and state space S , where $|S| = N$ (possibly infinite). Let π^T be a row vector denoting a probability distribution on S : so each element π_i denotes the probability of being in state i , and $\sum_{i=1}^N \pi_i = 1$, where $\pi_i \geq 0$ for all $i = 1, \dots, N$. The probability distribution π^T is an **equilibrium** distribution for the Markov chain if $\pi^T P = \pi^T$.

That is, π^T is an equilibrium distribution if

$$(\pi^T P)_j = \sum_{i=1}^N \pi_i p_{ij} = \pi_j \quad \text{for all } j = 1, \dots, N.$$

By the argument given on page 174, we have the following Theorem:

Theorem 9.2: Let $\{X_0, X_1, \dots\}$ be a Markov chain with transition matrix P . Suppose that π^T is an equilibrium distribution for the chain. If $X_t \sim \pi^T$ for any t , then $X_{t+r} \sim \pi^T$ for all $r \geq 0$. \square

Once a chain has hit an equilibrium distribution, *it stays there for ever*.

Note: There are several other names for an equilibrium distribution. If π^T is an equilibrium distribution, it is also called:

- **invariant:** *it doesn't change:* $\pi^T P = \pi^T$;
- **stationary:** *the chain 'stops' here.*

Stationarity: the Chain Station



a BUS station is where a BUS stops

a train station is where a train stops

a **workstation** is where ... ???



a stationary distribution is where a Markov chain stops

9.3 Finding an equilibrium distribution

Vector π^T is an equilibrium distribution for P if:

1. $\pi^T P = \pi^T$;
2. $\sum_{i=1}^N \pi_i = 1$;
3. $\pi_i \geq 0$ for all i .

Conditions 2 and 3 ensure that π^T is a *genuine probability distribution*.

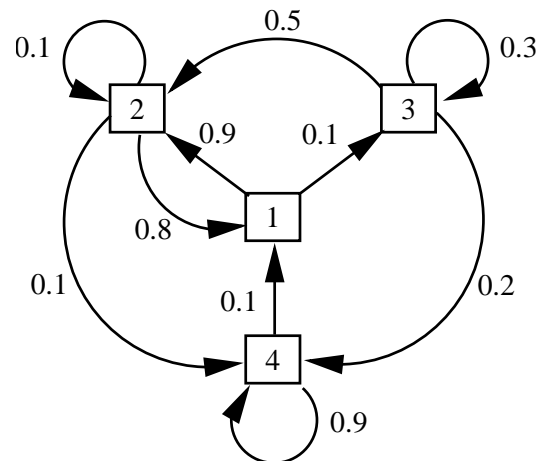
Condition 1 means that π is a row eigenvector of P .

Solving $\pi^T P = \pi^T$ by itself will just specify π up to a *scalar multiple*.

We need to include Condition 2 to scale π to a genuine probability distribution, and then check with Condition 3 that the scaled distribution is valid.

Example: Find an equilibrium distribution for the Markov chain below.

$$P = \begin{pmatrix} 0.0 & 0.9 & 0.1 & 0.0 \\ 0.8 & 0.1 & 0.0 & 0.1 \\ 0.0 & 0.5 & 0.3 & 0.2 \\ 0.1 & 0.0 & 0.0 & 0.9 \end{pmatrix}$$



Solution:

Let $\pi^T = (\pi_1, \pi_2, \pi_3, \pi_4)$.

The equations are $\pi^T P = \pi^T$ and $\pi_1 + \pi_2 + \pi_3 + \pi_4 = 1$.

$$\pi^T P = \pi^T \quad \Rightarrow \quad (\pi_1 \ \pi_2 \ \pi_3 \ \pi_4) \begin{pmatrix} 0.0 & 0.9 & 0.1 & 0.0 \\ 0.8 & 0.1 & 0.0 & 0.1 \\ 0.0 & 0.5 & 0.3 & 0.2 \\ 0.1 & 0.0 & 0.0 & 0.9 \end{pmatrix} = (\pi_1 \ \pi_2 \ \pi_3 \ \pi_4)$$

$$.8\pi_2 + .1\pi_4 = \pi_1 \quad (1)$$

$$.9\pi_1 + .1\pi_2 + .5\pi_3 = \pi_2 \quad (2)$$

$$.1\pi_1 + .3\pi_3 = \pi_3 \quad (3)$$

$$.1\pi_2 + .2\pi_3 + .9\pi_4 = \pi_4 \quad (4)$$

$$\text{Also} \quad \pi_1 + \pi_2 + \pi_3 + \pi_4 = 1. \quad (5)$$

$$(3) \Rightarrow \pi_1 = 7\pi_3$$

$$\begin{aligned} \text{Substitute in (2)} \Rightarrow .9(7\pi_3) + .5\pi_3 &= .9\pi_2 \\ \Rightarrow \pi_2 &= \frac{68}{9}\pi_3 \end{aligned}$$

$$\begin{aligned} \text{Substitute in (1)} \Rightarrow .8\left(\frac{68}{9}\pi_3\right) + .1\pi_4 &= 7\pi_3 \\ \Rightarrow \pi_4 &= \frac{86}{9}\pi_3 \end{aligned}$$

$$\begin{aligned} \text{Substitute all in (5)} \Rightarrow \pi_3 \left(7 + \frac{68}{9} + 1 + \frac{86}{9}\right) &= 1 \\ \Rightarrow \pi_3 &= \frac{9}{226} \end{aligned}$$

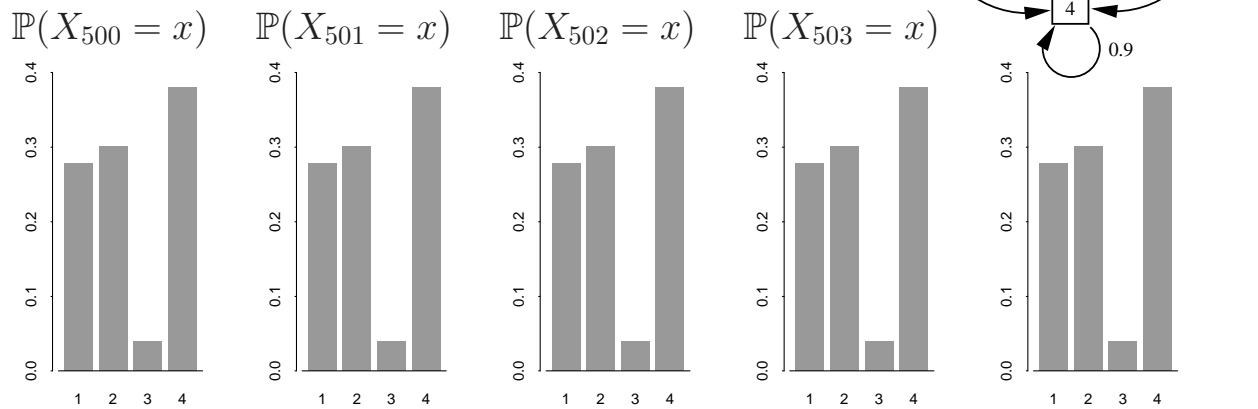
Overall:

$$\begin{aligned} \boldsymbol{\pi}^T &= \left(\frac{63}{226}, \frac{68}{226}, \frac{9}{226}, \frac{86}{226} \right) \\ &= (0.28, 0.30, 0.04, 0.38). \end{aligned}$$

This is the distribution the chain converged to in Section 9.1.

9.4 Long-term behaviour

In Section 9.1, we saw an example where the Markov chain wandered of its own accord into its equilibrium distribution:



This will always happen for this Markov chain. In fact, the distribution it converges to (found above) does not depend upon the starting conditions: *for ANY value of X_0 , we will always have $X_t \sim (0.28, 0.30, 0.04, 0.38)$ as $t \rightarrow \infty$.*

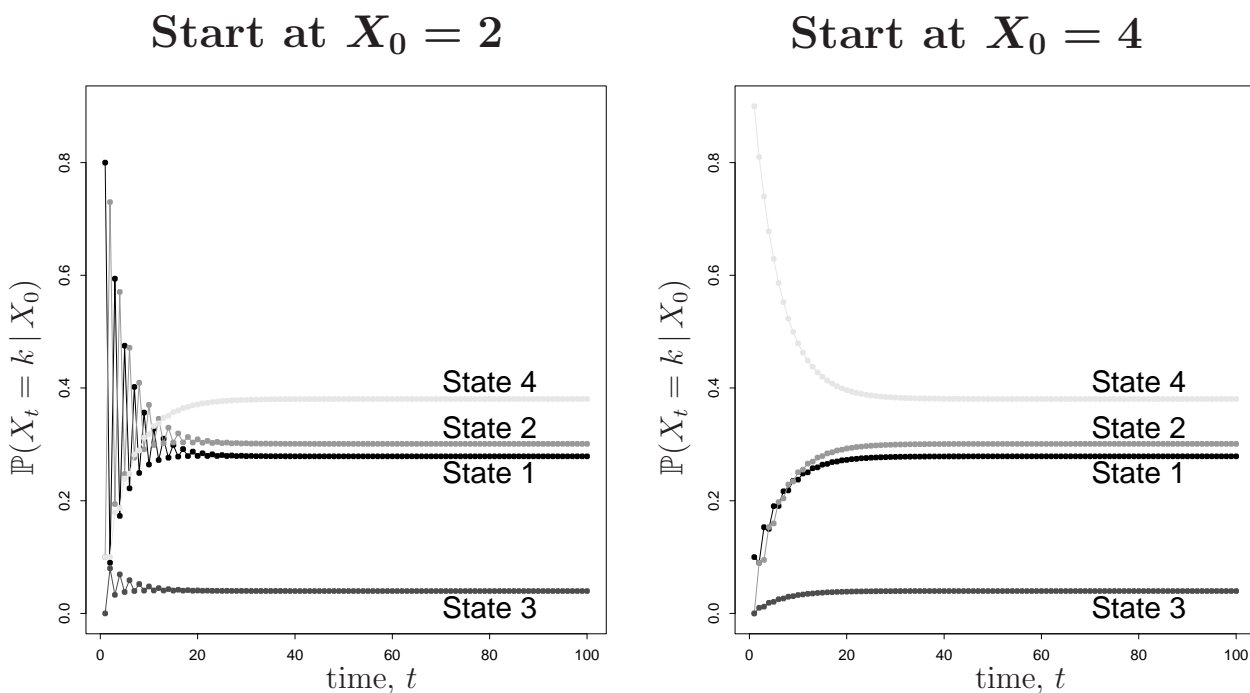
What is happening here is that *each row of the transition matrix P^t converges to the equilibrium distribution $(0.28, 0.30, 0.04, 0.38)$ as $t \rightarrow \infty$:*

$$P = \begin{pmatrix} 0.0 & 0.9 & 0.1 & 0.0 \\ 0.8 & 0.1 & 0.0 & 0.1 \\ 0.0 & 0.5 & 0.3 & 0.2 \\ 0.1 & 0.0 & 0.0 & 0.9 \end{pmatrix} \Rightarrow P^t \rightarrow \begin{pmatrix} 0.28 & 0.30 & 0.04 & 0.38 \\ 0.28 & 0.30 & 0.04 & 0.38 \\ 0.28 & 0.30 & 0.04 & 0.38 \\ 0.28 & 0.30 & 0.04 & 0.38 \end{pmatrix} \text{ as } t \rightarrow \infty.$$

(If you have a calculator that can handle matrices, try finding P^t for $t = 20$ and $t = 30$: you will find the matrix is already converging as above.)

This convergence of P^t means that *for large t , no matter WHICH state we start in, we always have probability*

- about **0.28** of being in State **1** after t steps;
- about **0.30** of being in State **2** after t steps;
- about **0.04** of being in State **3** after t steps;
- about **0.38** of being in State **4** after t steps.



The **left graph** shows the probability of getting from state 2 to state k in t steps, as t changes: $(P^t)_{2,k}$ for $k = 1, 2, 3, 4$.

The **right graph** shows the probability of getting from state 4 to state k in t steps, as t changes: $(P^t)_{4,k}$ for $k = 1, 2, 3, 4$.

The *initial behaviour* differs greatly for the different start states.

The *long-term behaviour* (large t) is the same for both start states.

However, this does not always happen. Consider the two-state chain below:



As t gets large, P^t *does not converge*:

$$P^{500} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad P^{501} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad P^{502} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad P^{503} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \dots$$

For this Markov chain, we *never ‘forget’ the initial start state*.

General formula for P^t

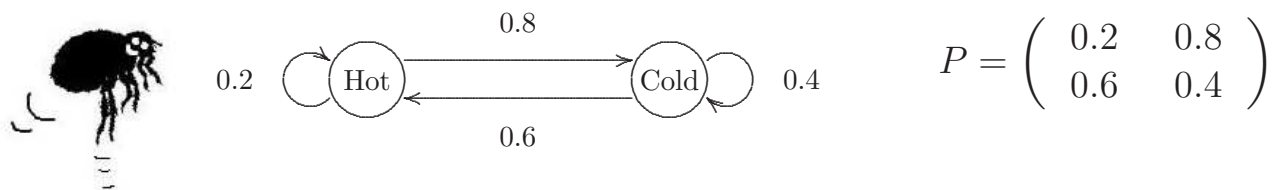
We have seen that we are interested in whether P^t converges to *a fixed matrix with all rows equal* as $t \rightarrow \infty$.

If it does, then the Markov chain will *reach an equilibrium distribution that does not depend upon the starting conditions*.

The equilibrium distribution is then given by *any row of the converged P^t* .

It can be shown that a general formula is available for P^t for any t , based on the eigenvalues of P . Producing this formula is beyond the scope of this course, but if you are given the formula, you should be able to recognise whether P^t is going to converge to a fixed matrix with all rows the same.

Example 1:



We can show that the general solution for P^t is:

$$P^t = \frac{1}{7} \left\{ \begin{pmatrix} 3 & 4 \\ 3 & 4 \end{pmatrix} + \begin{pmatrix} 4 & -4 \\ -3 & 3 \end{pmatrix} (-0.4)^t \right\}$$

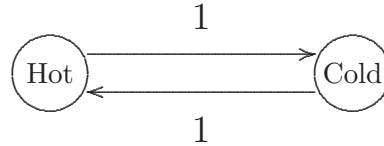
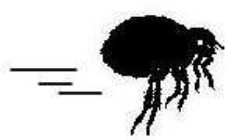
As $t \rightarrow \infty$, $(-0.4)^t \rightarrow 0$, so

$$P^t \rightarrow \frac{1}{7} \begin{pmatrix} 3 & 4 \\ 3 & 4 \end{pmatrix} = \begin{pmatrix} \frac{3}{7} & \frac{4}{7} \\ \frac{3}{7} & \frac{4}{7} \end{pmatrix}$$

This Markov chain will therefore converge to the equilibrium distribution $\pi^T = (\frac{3}{7}, \frac{4}{7})$ as $t \rightarrow \infty$, regardless of whether the flea starts in state 1 or state 2.

Exercise: Verify that $\pi^T = (\frac{3}{7}, \frac{4}{7})$ is the same as the result you obtain from solving the equilibrium equations: $\pi^T P = \pi^T$ and $\pi_1 + \pi_2 = 1$.

Example 2: Purposeflea knows exactly what he is doing, so his probabilities are all 1:



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

We can show that the general solution for P^t is:

$$P^t = \frac{1}{2} \left\{ \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} + \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} (-1)^t \right\}$$

As $t \rightarrow \infty$, $(-1)^t$ does not converge to 0, so

$$P^t = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \text{if } t \text{ is even,}$$

$$P^t = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \text{if } t \text{ is odd,}$$

for all t .

In this example, P^t never converges to a matrix with both rows identical as t gets large. The chain never ‘forgets’ its starting conditions as $t \rightarrow \infty$.

Exercise: Verify that this Markov chain *does* have an equilibrium distribution, $\pi^T = (\frac{1}{2}, \frac{1}{2})$. However, the chain does not *converge* to this distribution as $t \rightarrow \infty$.

These examples show that some Markov chains forget their starting conditions in the long term, and ensure that X_t will have the same distribution as $t \rightarrow \infty$ regardless of where we started at X_0 . However, for other Markov chains, the initial conditions are never forgotten. In the next sections we look for general criteria that will ensure the chain converges.

Target Result:

- If a Markov chain is *irreducible* and *aperiodic*, and if an equilibrium distribution π^T *exists*, then the chain converges to this distribution as $t \rightarrow \infty$, regardless of the initial starting states.

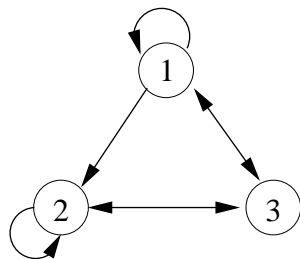
To make sense of this, we need to revise the concept of *irreducibility*, and introduce the idea of *aperiodicity*.

9.5 Irreducibility

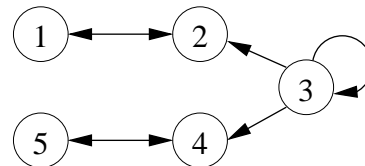
Recall from Chapter 8:

Definition: A Markov chain or transition matrix P is said to be irreducible if $i \leftrightarrow j$ for all $i, j \in S$. That is, the chain is irreducible if the state space S is a single communicating class.

An irreducible Markov chain consists of a single class.



Irreducible



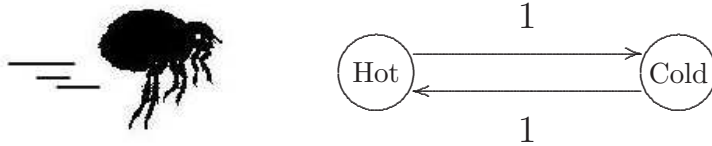
Not irreducible

Irreducibility of a Markov chain is important for convergence to equilibrium as $t \rightarrow \infty$, because *we want the convergence to be independent of start state*.

This can happen if the chain is irreducible. When the chain is not irreducible, different start states might cause the chain to get stuck in different closed classes. In the example above, a start state of $X_0 = 1$ means that the chain is restricted to states 1 and 2 as $t \rightarrow \infty$, whereas a start state of $X_0 = 4$ means that the chain is restricted to states 4 and 5 as $t \rightarrow \infty$. A single convergence that ‘forgets’ the initial state is therefore not possible.

9.6 Periodicity

Consider the Markov chain with transition matrix $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$.



Suppose that $X_0 = 1$.

Then $X_t = 1$ for all even values of t , and $X_t = 2$ for all odd values of t .

This sort of behaviour is called **periodicity**: *the Markov chain can only return to a state at particular values of t .*

Clearly, periodicity of the chain will interfere with convergence to an equilibrium distribution as $t \rightarrow \infty$. For example,

$$\mathbb{P}(X_t = 1 \mid X_0 = 1) = \begin{cases} 1 & \text{for even values of } t, \\ 0 & \text{for odd values of } t. \end{cases}$$

Therefore, the probability can not converge to any single value as $t \rightarrow \infty$.

Period of state i

To formalize the notion of periodicity, we define the **period** of a state i . Intuitively, *the period is defined so that the time taken to get from state i back to state i again is always a multiple of the period.*

In the example above, the chain can return to state 1 after 2 steps, 4 steps, 6 steps, 8 steps, ...

The period of state 1 is therefore 2.

In general, the chain can return from state i back to state i again in t steps if $(P^t)_{ii} > 0$. This prompts the following definition.

Definition: The **period** $d(i)$ of a state i is

$$d(i) = \gcd\{t : (P^t)_{ii} > 0\},$$

the greatest common divisor of the times at which return is possible.

Definition: The state i is said to be **periodic** if $d(i) > 1$.

For a periodic state i , $(P^t)_{ii} = 0$ if t is not a multiple of $d(i)$.

Definition: The state i is said to be **aperiodic** if $d(i) = 1$.

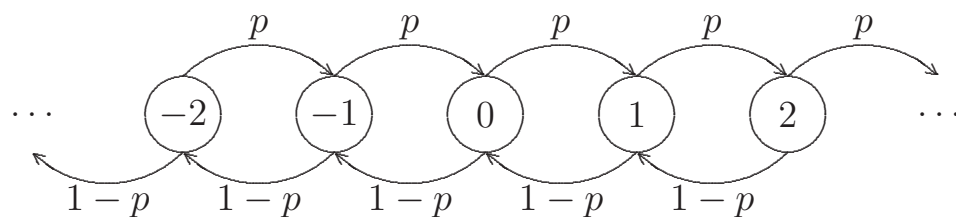
If state i is aperiodic, it means that *return to state i is not limited only to regularly repeating times*.

For convergence to equilibrium as $t \rightarrow \infty$, we will be interested only in ***aperiodic states***.

The following examples show how to calculate the period for both aperiodic and periodic states.

Examples: Find the periods of the given states in the following Markov chains, and state whether or not the chain is irreducible.

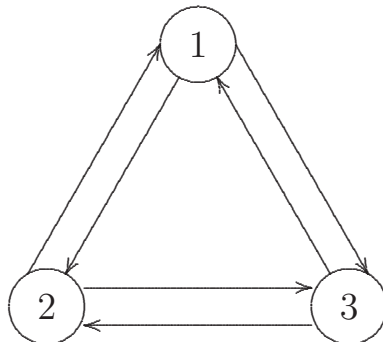
1. The simple random walk.



$$d(0) = \gcd\{2, 4, 6, \dots\} = 2.$$

Chain is irreducible.

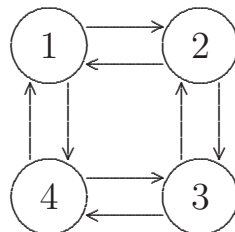
2.



$$d(1) = \gcd\{2, 3, 4, \dots\} = 1.$$

Chain is irreducible.

3.



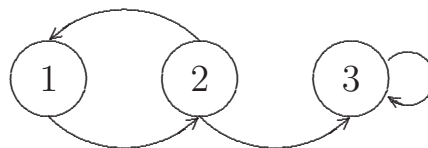
$$d(1) = \gcd\{2, 4, 6, \dots\} = 2.$$

Chain is irreducible.

4.

$$d(1) = \gcd\{2, 4, 6, \dots\} = 2.$$

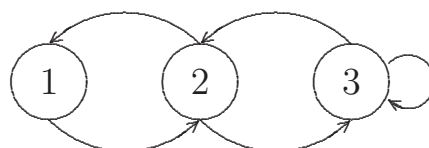
Chain is NOT irreducible (i.e. Reducible).



5.

$$d(1) = \gcd\{2, 4, 5, 6, \dots\} = 1.$$

Chain is irreducible.



9.7 Convergence to Equilibrium

We now draw together the threads of the previous sections with the following results.

Fact: If $i \leftrightarrow j$, then i and j have the same period. (Proof omitted.)

This leads immediately to the following result:

If a Markov chain is *irreducible* and has *one* aperiodic state, then *all* states are aperiodic.

We can therefore talk about an irreducible, aperiodic chain, meaning that *all states are aperiodic*.

Theorem 9.7: Let $\{X_0, X_1, \dots\}$ be an irreducible and aperiodic Markov chain with transition matrix P . Suppose that there *exists* an equilibrium distribution π^T . Then, from *any* starting state i , and for any end state j ,

$$\mathbb{P}(X_t = j \mid X_0 = i) \rightarrow \pi_j \quad \text{as } t \rightarrow \infty.$$

In particular,

$$(P^t)_{ij} \rightarrow \pi_j \quad \text{as } t \rightarrow \infty, \text{ for } \underline{\text{all}} \ i \text{ and } j,$$

so P^t converges to a matrix with all rows identical and equal to π^T . □

For an irreducible, aperiodic Markov chain,
with finite or infinite state space,
the existence of an equilibrium distribution π^T ensures
that the Markov chain will converge to π^T as $t \rightarrow \infty$.

Note: If the state space is infinite, it is not guaranteed that an equilibrium distribution π^T exists. See Example 3 below.

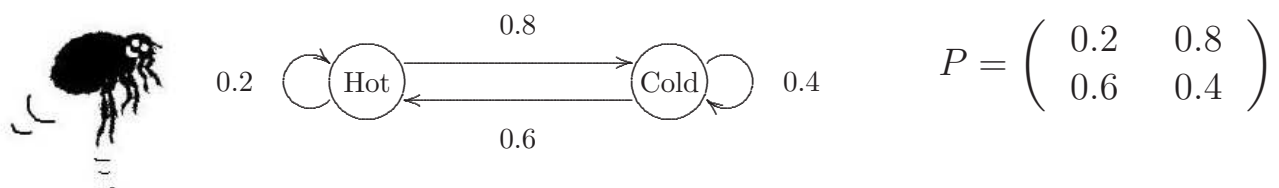
Note: If the chain converges to an equilibrium distribution π^T as $t \rightarrow \infty$, then *the long-run proportion of time spent in state k is π_k* .

9.8 Examples

A typical exam question gives you a Markov chain on a finite state space and asks if it converges to an equilibrium distribution as $t \rightarrow \infty$. An equilibrium distribution will always exist for a finite state space. You need to check whether the chain is irreducible and aperiodic. If so, it will converge to equilibrium. If the chain is irreducible but *periodic*, it cannot converge to an equilibrium distribution that is independent of start state. If the chain is *reducible*, it may or may not converge.

The first two examples are the same as the ones given in Section 9.4.

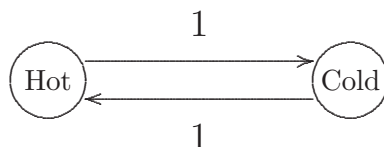
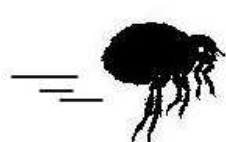
Example 1: State whether the Markov chain below converges to an equilibrium distribution as $t \rightarrow \infty$.



The chain is irreducible and aperiodic, and an equilibrium distribution will exist for a finite state space. So the chain does converge.

(From Section 9.4, the chain converges to $\pi^T = (\frac{3}{7}, \frac{4}{7})$ as $t \rightarrow \infty$.)

Example 2: State whether the Markov chain below converges to an equilibrium distribution as $t \rightarrow \infty$.



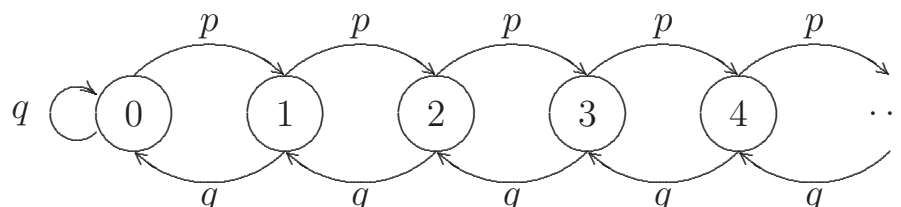
$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

The chain is irreducible, but it is NOT aperiodic: period = 2.

Thus the chain does NOT converge to an equilibrium distribution as $t \rightarrow \infty$.

It is important to check for aperiodicity, because the existence of an equilibrium distribution does NOT ensure convergence to this distribution if the matrix is not aperiodic.

Example 3: Random walk with retaining barrier at 0.



Find whether the chain converges to equilibrium as $t \rightarrow \infty$, and if so, find the equilibrium distribution.

The chain is irreducible and aperiodic, so if an equilibrium distribution exists, then the chain will converge to this distribution as $t \rightarrow \infty$.

However, the chain has an infinite state space, so we cannot guarantee that an equilibrium distribution exists.

Try to solve the equilibrium equations:

$$\pi^T P = \pi^T \text{ and } \sum_{i=0}^{\infty} \pi_i = 1.$$

$$P = \begin{pmatrix} q & p & 0 & 0 & \dots \\ q & 0 & p & 0 & \dots \\ 0 & q & 0 & p & \dots \\ \vdots & & & & \end{pmatrix} \quad \begin{array}{lcl} q\pi_0 + q\pi_1 & = & \pi_0 \\ p\pi_0 + q\pi_2 & = & \pi_1 \\ p\pi_1 + q\pi_3 & = & \pi_2 \\ \vdots & & \\ p\pi_{k-1} + q\pi_{k+1} & = & \pi_k \end{array} \quad \begin{array}{l} (\star) \\ \\ \\ \text{for } k = 1, 2, \dots \end{array}$$

From (\star) , we have $p\pi_0 = q\pi_1$,

$$\text{so } \pi_1 = \frac{p}{q}\pi_0$$

$$\Rightarrow \pi_2 = \frac{1}{q}(\pi_1 - p\pi_0) = \frac{1}{q}\left(\frac{p}{q}\pi_0 - p\pi_0\right) = \frac{p}{q}\left(\frac{1-q}{q}\right)\pi_0 = \left(\frac{p}{q}\right)^2\pi_0.$$

We suspect that $\pi_k = \left(\frac{p}{q}\right)^k \pi_0$. Prove by induction.

The hypothesis is true for $k = 0, 1, 2$. Suppose that $\pi_k = \left(\frac{p}{q}\right)^k \pi_0$. Then

$$\begin{aligned} \pi_{k+1} &= \frac{1}{q}(\pi_k - p\pi_{k-1}) \\ &= \frac{1}{q}\left\{\left(\frac{p}{q}\right)^k \pi_0 - p\left(\frac{p}{q}\right)^{k-1} \pi_0\right\} \\ &= \frac{p^k}{q^k}\left(\frac{1}{q} - 1\right)\pi_0 \\ &= \left(\frac{p}{q}\right)^{k+1} \pi_0. \end{aligned}$$

The inductive hypothesis holds, so $\pi_k = \left(\frac{p}{q}\right)^k \pi_0$ for all $k \geq 0$.

We now need $\sum_{i=0}^{\infty} \pi_i = 1$, i.e. $\pi_0 \sum_{k=0}^{\infty} \left(\frac{p}{q}\right)^k = 1$.

The sum is a Geometric series, and converges only for $\left|\frac{p}{q}\right| < 1$. Thus when $p < q$, we have

$$\pi_0 \left(\frac{1}{1 - \frac{p}{q}}\right) = 1 \Rightarrow \pi_0 = 1 - \frac{p}{q}.$$

If $p \geq q$, there is no equilibrium distribution.

Solution:

If $p < q$, the chain converges to an equilibrium distribution π , where $\pi_k = \left(1 - \frac{p}{q}\right) \left(\frac{p}{q}\right)^k$ for $k = 0, 1, \dots$

If $p \geq q$, the chain does not converge to an equilibrium distribution as $t \rightarrow \infty$.

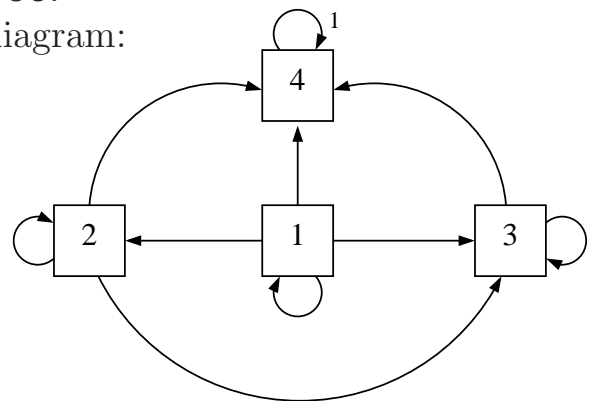
Example 4: Sketch of Exam Question 2006.

Consider a Markov chain with transition diagram:

- (a) Identify all communicating classes.
For each class, state whether or not it is closed.

Classes are:

$\{1\}, \{2\}, \{3\}$ (each not closed);
 $\{4\}$ (closed).



- (b) State whether the Markov chain is irreducible, and whether or not all states are aperiodic.

Not irreducible: there are 4 classes.

All states are aperiodic.

- (c) The equilibrium distribution is $\pi^T = (0, 0, 0, 1)$. Does the Markov chain converge to this distribution as $t \rightarrow \infty$, regardless of its start state?

Yes, it clearly will converge to $\pi^T = (0, 0, 0, 1)$, despite failure of irreducibility.

Note: Equilibrium results also exist for chains that are *not* aperiodic. Also, states can be classified as **transient** (return to the state is not certain), **null recurrent** (return to the state is certain, but the expected return time is infinite), and **positive recurrent** (return to the state is certain, and the expected return time is finite). For each type of state, the long-term behaviour is known:

- If the state k is **transient** or **null-recurrent**,

$$\mathbb{P}(X_t = k \mid X_0 = k) = (P^t)_{kk} \rightarrow 0 \text{ as } t \rightarrow \infty.$$

- If the state is **positive recurrent**, then

$$\mathbb{P}(X_t = k \mid X_0 = k) = (P^t)_{kk} \rightarrow \pi_k \text{ as } t \rightarrow \infty, \text{ where } \pi_k > 0.$$

The expected return time for the state is $1/\pi_k$.

A detailed treatment is available at
<http://www.statslab.cam.ac.uk/~james/Markov/>.

9.9 Special Process: the Two-Armed Bandit

A well-known problem in probability is called the **two-armed bandit** problem. The name is a reference to a type of gambling machine called the two-armed bandit. The two arms of the two-armed bandit offer different rewards, and the gambler has to decide which arm to play without knowing which is the better arm.



One-armed bandit

A similar problem arises when doctors are experimenting with two different treatments, without knowing which one is better. Call the treatments A and B . One of them is likely to be better, but we don't know which one. A series of patients will each be given one of the treatments. We aim to find a strategy that ensures that as many as possible of the patients are given the *better* treatment — though we don't know which one this is.

Suppose that, for any patient, treatment A has $\mathbb{P}(\text{success}) = \alpha$, and treatment B has $\mathbb{P}(\text{success}) = \beta$, and all patients are independent. Assume that $0 < \alpha < 1$ and $0 < \beta < 1$.

First let's look at a simple strategy the doctors might use:

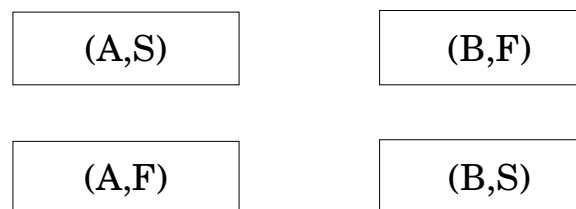
- The **random strategy** for allocating patients to treatments A and B is to choose from the two treatments at random, each with probability 0.5, for each patient.
- Let p_R be the overall probability of **success** for each patient with the random strategy. Show that $p_R = \frac{1}{2}(\alpha + \beta)$.

The **two-armed bandit strategy** is more clever. For the first patient, we choose treatment A or B at random (probability 0.5 each). If patient n is given treatment A and it is *successful*, then we use treatment A again for patient $n+1$, for all $n = 1, 2, 3, \dots$. If A is a failure for patient n , we switch to treatment B for patient $n+1$. A similar rule is applied if patient n is given treatment B : if it is successful, we keep B for patient $n+1$; if it fails, we switch to A for patient $n+1$.

Define the **two-armed bandit process** to be a Markov chain with state space $\{(A, S), (A, F), (B, S), (B, F)\}$, where (A, S) means that patient n is given treatment A and it is successful, and so on.

Transition diagram:

Exercise: Draw on the missing arrows and find their probabilities in terms of α and β .



Transition matrix:

$$\begin{matrix} AS \\ AF \\ BS \\ BF \end{matrix} \begin{pmatrix} AS & AF & BS & BF \\ & & & \end{pmatrix}$$

Probability of success under the two-armed bandit strategy

Define p_T to be the long-run probability of **success** using the two-armed bandit strategy.

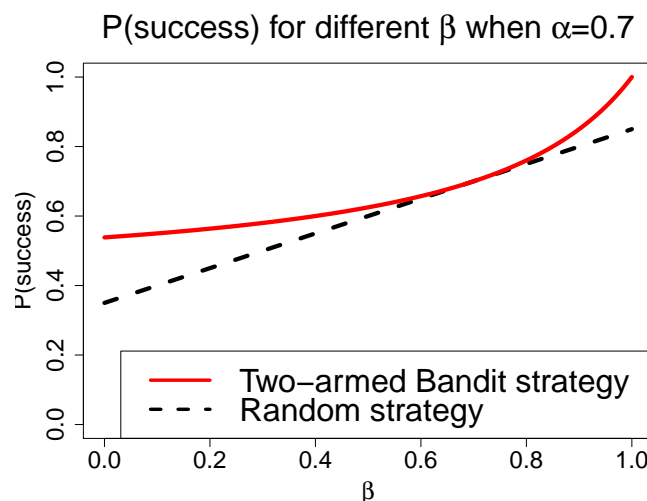
Exercise: Find the equilibrium distribution π for the two-armed bandit process. Hence show that the long-run probability of success for each patient under this strategy is:

$$p_T = \frac{\alpha + \beta - 2\alpha\beta}{2 - \alpha - \beta}.$$

Which strategy is better?

Exercise: Prove that $p_T - p_R \geq 0$ always, regardless of the values of α and β .

This proves that the two-armed bandit strategy is always better than, or equal to, the random strategy. It shows that we have been able to construct a strategy that gives all patients an increased chance of success, even though we don't know which treatment is better!



The graph shows the probability of success under the two different strategies, for $\alpha = 0.7$ and for $0 \leq \beta \leq 1$. Notice how $p_T \geq p_R$ for all possible values of β .