
Methods of Monte Carlo Simulation II

Ulm University
Institute of Stochastics

Lecture Notes
Dr. Tim Brereton
Summer Term 2014

ULM, 2014

Contents

1	Some Simple Stochastic Processes	7
1.1	Stochastic Processes	7
1.2	Random Walks	7
1.2.1	Bernoulli Processes	7
1.2.2	Random Walks	10
1.2.3	Probabilities of Random Walks	13
1.2.4	Distribution of X_n	13
1.2.5	First Passage Time	14
2	Estimators	17
2.1	Bias, Variance, the Central Limit Theorem and Mean Square Error	19
2.2	Non-Asymptotic Error Bounds	22
2.3	Big O and Little o Notation	23
3	Markov Chains	25
3.1	Simulating Markov Chains	28
3.1.1	Drawing from a Discrete Uniform Distribution	28
3.1.2	Drawing From A Discrete Distribution on a Small State Space	28
3.1.3	Simulating a Markov Chain	28
3.2	Communication	29
3.3	The Strong Markov Property	30
3.4	Recurrence and Transience	31
3.4.1	Recurrence of Random Walks	33
3.5	Invariant Distributions	34
3.6	Limiting Distribution	36
3.7	Reversibility	37
4	The Poisson Process	39
4.1	Point Processes on $[0, \infty)$	39

4.2	Poisson Process	41
4.2.1	Order Statistics and the Distribution of Arrival Times	44
4.2.2	Distribution of Arrival Times	45
4.3	Simulating Poisson Processes	46
4.3.1	Using the Infinitesimal Definition to Simulate Approx- imately	46
4.3.2	Simulating the Arrival Times	47
4.3.3	Simulating the Inter-Arrival Times	48
4.4	Inhomogenous Poisson Processes	48
4.5	Simulating an Inhomogenous Poisson Process	49
4.5.1	Acceptance-Rejection	49
4.5.2	Infinitesimal Approach (Approximate)	50
4.6	Compound Poisson Processes	51
5	Continuous Time Markov Chains	53
5.1	Transition Function	53
5.2	Infinitesimal Generator	54
5.3	Continuous Time Markov Chains	54
5.4	The Jump Chain and Holding Times	55
5.5	Examples of Continuous Time Markov Chains	56
5.5.1	Poisson Process	56
5.5.2	Birth-Death Process	56
5.6	Simulating Continuous Time Markov Chains	56
5.7	The Relationship Between P and Q in the Finite Case	58
5.8	Irreducibility, Recurrence and Positive Recurrence	59
5.9	Invariant Measures and Stationary Distribution	61
5.9.1	Reversibility and Detailed Balance	62
6	Gaussian Processes	63
6.1	The Multivariate Normal Distribution	63
6.1.1	Symmetric Positive Definite and Semi-Positive Definite Matrices	65
6.1.2	Densities of Multivariate Normals	66
6.1.3	Simulating Multivariate Normals	67
6.2	Simulating a Gaussian Processes Version 1	68
6.3	Stationary and Weak Stationary Gaussian Processes	71
6.4	Finite Dimensional Distributions	71
6.5	Marginal and Conditional Multivariate Normal Distributions	72
6.6	Interpolating Gaussian Processes	73
6.7	Markovian Gaussian Processes	75
6.8	Brownian Motion	77

7	Random Fields	81
7.1	Gaussian Random Fields	81
7.2	Markov Random Fields	85
7.3	Gaussian Random Markov Fields	86

Chapter 1

Some Simple Stochastic Processes

1.1 Stochastic Processes

To begin, we need to define the basic objects we will be learning about.

Definition 1.1.1 (Stochastic Process). A stochastic process is a set of random variables $\{X_i\}_{i \in I}$, taking values in a *state space* \mathcal{X} , with index set $I \subset \mathbb{R}$.

In general, i represents a point in time. However, it could represent a point in 1D space as well.

We will say a process is *discrete time* if I is discrete. For example, I could be \mathbb{N} or $\{1, 2, 3, 10, 20\}$. Normally, when we talk about discrete time stochastic processes, we will use the index n (e.g., $\{X_n\}_{n \in \mathbb{N}}$).

We will say a process is *continuous time* if I is an interval. For example $[0, \infty)$ or $[1, 2]$. Normally, when we talk about continuous time stochastic processes, we will use the index t (e.g. $\{X_t\}_{t \in [0, T]}$).

1.2 Random Walks

1.2.1 Bernoulli Processes

One of the simplest stochastic processes is a *random walk*. However, even though the random walk is very simple, it has a number of properties that will be important when we think about more complicated processes. To define a random walk, we begin with an even simpler process called a *Bernoulli process*. A Bernoulli process is a discrete time process taking values in the state space $\{0, 1\}$.

Definition 1.2.1 (Bernoulli Process). A Bernoulli process with parameter $p \in [0, 1]$ is a sequence of *independent and identically distributed* (i.i.d.) random variables, $\{Y_n\}_{n \geq 1}$, such that

$$\mathbb{P}(Y_1 = 1) = 1 - \mathbb{P}(Y_1 = 0) = p. \quad (1.1)$$

A variable with the *probability mass function* (pmf) described by (1.1) is called a Bernoulli random variable with distribution $\text{Ber}(p)$.

In the case where $p = 1/2$, you can think of a Bernoulli process as a sequence of fair coin tosses. In the case where $p \neq 1/2$, you can think of a Bernoulli process as a sequence of tosses of an unfair coin.

Note that if $U \sim \text{U}(0, 1)$, then $\mathbb{P}(U \leq p) = p$ for $p \in [0, 1]$. This means if we define the random variable $Y = \mathbb{I}(U \leq p)$, where $\mathbb{I}(\cdot)$ is the indicator function and $U \sim \text{U}(0, 1)$, we have

$$\mathbb{P}(Y = 1) = 1 - \mathbb{P}(Y = 0) = p.$$

In Matlab, we can make these variables as follows.

Listing 1.1: Generating a Bernoulli Random Variable

```
1 Y = (rand <= p)
```

If we want to make a realization of the first n steps of a Bernoulli process, we simply make n such variables. One way to do this is the following.

Listing 1.2: Generating a Bernoulli Process 1

```
1 n = 20; p = 0.6; Y = zeros(N,1);
2
3 for i = 1:n
4     Y(i) = (rand <= p);
5 end
```

We can do this more quickly in Matlab though.

Listing 1.3: Generating a Bernoulli Process 2

```
1 n = 20; p = 0.6;
2
3 Y = (rand(n,1) <= p);
```

It is important to be able to visualize stochastic objects. One way to represent / visualize a Bernoulli process is to put $n = 1, 2, \dots$ on the x -axis and the values of Y_n on the y -axis. I draw a line (almost of length 1) at each value of Y_n , as this is easier to see than dots. Figure 1.2.1 shows a realization of the first 20 steps of a Bernoulli process.

Listing 1.4: Generating and Visualizing a Bernoulli Process

```

1 n = 20; p = 0.6;
2
3 Y = (rand(n,1) <= p);
4
5 clf;
6 axis([1 n+1 -0.5 1.5]);
7 for i = 1:n
8     line([i, i+.9],[Y(i) Y(i)]);
9 end

```

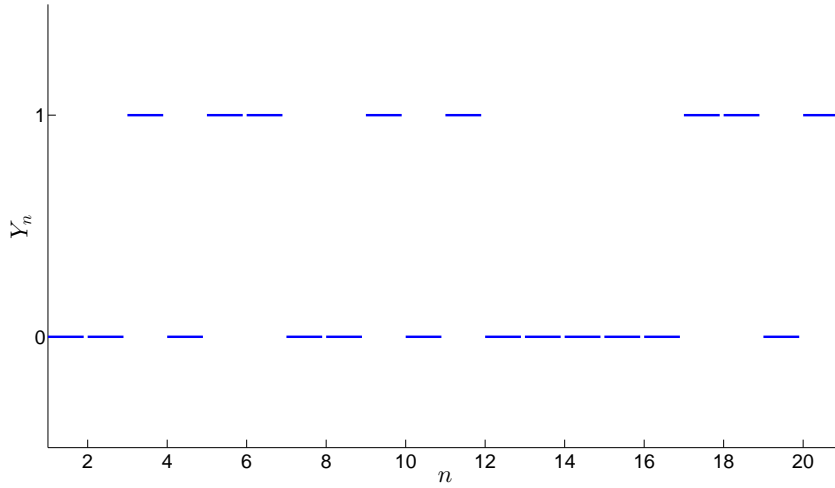


Figure 1.2.1: A realization of the first 20 steps of a Bernoulli process.

Now, it is easy to write out the probability of seeing a particular realization of n steps of a Bernoulli process. This is given by

$$\mathbb{P}(Y_n = y_n, Y_{n-1} = y_{n-1}, \dots, Y_1 = y_1) = p^{N_U} (1-p)^{n-N_U}.$$

where N_U is the number of times the Bernoulli process takes the value 1. More technically, we define N_U as

$$N_U = \#\{0 \leq i \leq n : y_i = 1\},$$

where $\#$ denotes the cardinality of the set.

1.2.2 Random Walks

Now that we can generate Bernoulli processes, we are ready to consider our main object of interest in this chapter: the random walk. The random walk is a discrete time random process taking values in the state space \mathbb{Z} .

Definition 1.2.2 (Random Walk). Given a Bernoulli process $\{Y_n\}_{n \geq 1}$ with parameter p , we define a random walk $\{X_n\}_{n \geq 0}$ with parameter p and initial condition $X_0 = x_0$ by

$$X_{n+1} = X_n + (2Y_{n+1} - 1).$$

Note that $(2Y_n - 1)$ is 1 if $Y_n = 1$ and -1 if $Y_n = 0$. So, essentially, at each step of the process, it goes up or down by 1.

There are lots of ways to think about random walks. One way is in terms of gambling (which is a classical setting for probability). Think of a game which is free to enter. The person running the game flips a coin. With probability p , the coin shows heads and you win €1. With probability $1 - p$, you lose €1. If you start with € x_0 , and assuming you can go into debt, X_n is the random variable describing how much money you have after n games.

Given that we know how to generate a Bernoulli process, it is straightforward to generate a random walk.

Listing 1.5: Generating a Random Walk 1

```

1 n = 100; X = zeros(n,1);
2 p = 0.5; X_0 = 0;
3
4 Y = (rand(n,1) <= p);
5 X(1) = X_0 + 2*Y(1) - 1;
6 for i = 2:n
7     X(i) = X(i-1) + 2*Y(i) - 1;
8 end

```

A more compact way is as follows.

Listing 1.6: Generating a Random Walk 2

```

1 n = 100; X = zeros(n,1);
2 p = 0.5; X_0 = 0;
3
4 Y = (rand(n,1) <= p);
5 X = X_0 + cumsum((2*Y - 1));

```

It is good to be able to plot these. Here, we have at least two options. One is to draw a line for each value of X_n , as we did for the Bernoulli process.

Because the value of $\{X_n\}_{n \geq 0}$ changes at every value of n , we can draw the lines to be length 1. This approach is nice, because it reinforces the idea that $\{X_n\}_{n \geq 0}$ jumps at each value of n (this is made obvious by the gaps between the lines).

Listing 1.7: Plotting a Random Walk 1

```

1 clf
2 axis([1 n+1 min(X)-1 max(X)+1]);
3 for i = 1:n
4     line([i,i+1],[X(i) X(i)], 'LineWidth',3);
5 end

```

For larger values of n , it is help to draw vertical lines as well as horizontal lines (otherwise, things get a bit hard to see). We can do this as follows.

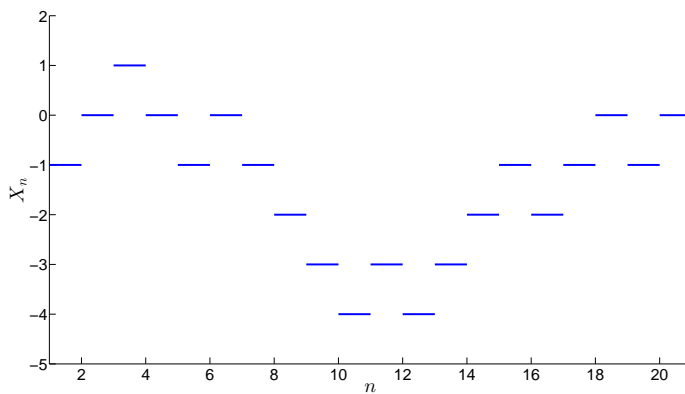
Listing 1.8: Plotting a Random Walk 2

```

1 clf
2 stairs(X);
3 axis([1 n+1 min(X)-1 max(X)+1]);

```

Below are a number of realizations of random walks. Figures 1.2.2 and 1.2.3 are generated using lines. Figures 1.2.4 and 1.2.5 are generated using the ‘stairs’ function.

Figure 1.2.2: A realization of the first 20 steps of a random walk with $p = 0.5$.

Note how, as the number of steps simulated gets bigger, the sample paths look wilder and wilder.

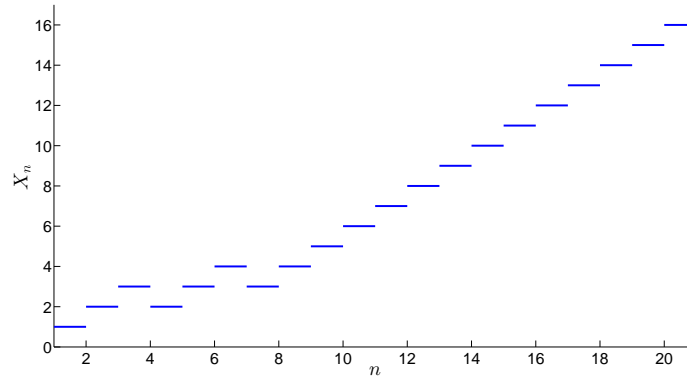


Figure 1.2.3: A realization of the first 20 steps of a random walk with $p = 0.9$.

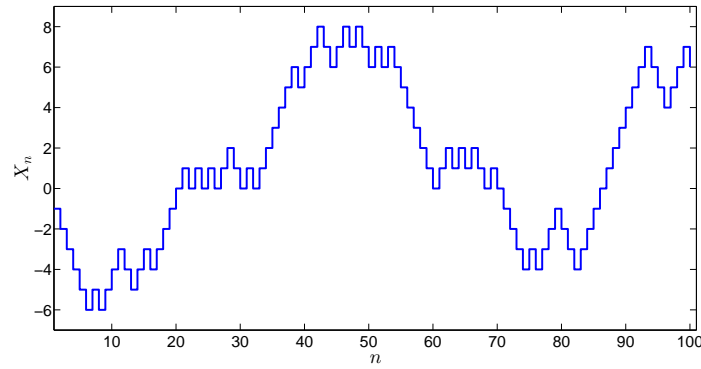


Figure 1.2.4: A realization of the first 100 steps of a random walk with $p = 0.5$.

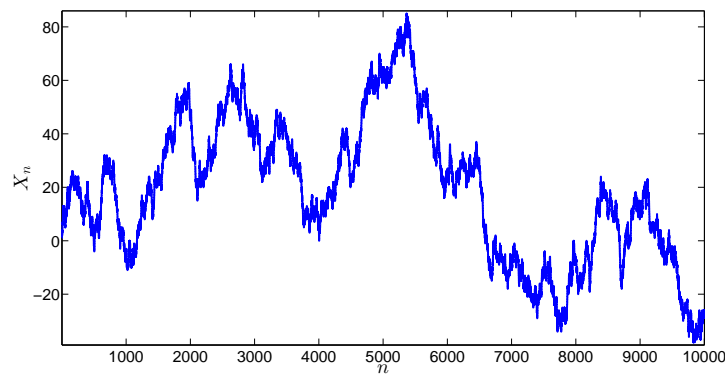


Figure 1.2.5: A realization of the first 10000 steps of a random walk with $p = 0.5$.

1.2.3 Probabilities of Random Walks

A random walk is a special kind of process, called a *Markov process*. Discrete time discrete state space Markov processes (called *Markov chains*) are defined by the following property.

Definition 1.2.3 (Markov Property: Discrete Time Discrete State Space Version). We say a discrete time, discrete state space stochastic process, $\{X_n\}_{n \geq 0}$, has the Markov property if

$$\mathbb{P}(X_{n+1} = x_{n+1} \mid X_n = x_n, \dots, X_0 = x_0) = \mathbb{P}(X_{n+1} = x_{n+1} \mid X_n = x_n),$$

for all $n \geq 0$ and $x_0, \dots, x_n \in \mathcal{X}$.

If we think about how a random walk can move after n steps (assuming we know the values of the first n steps), we have

$$\mathbb{P}(X_{n+1} = x_{n+1} \mid X_n = x_n, \dots, X_0 = x_0) = \begin{cases} p & \text{if } x_{n+1} = x_n + 1, \\ 1 - p & \text{if } x_{n+1} = x_n - 1, \\ 0 & \text{otherwise.} \end{cases}$$

Note that these probabilities do not depend on any value other than X_n . That is, it is always the case that $\mathbb{P}(X_{n+1} = x_{n+1} \mid X_n = x_n, \dots, X_0 = x_0) = \mathbb{P}(X_{n+1} = x_{n+1} \mid X_n = x_n)$, so $\{X_n\}_{n \geq 0}$ is a Markov chain.

Using this result, we can assign probabilities to various paths of the random walk. That is, we can write

$$\begin{aligned} & \mathbb{P}(X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1) \\ &= \mathbb{P}(X_n = x_n \mid X_{n-1} = x_{n-1}, \dots, X_1 = x_1) \cdots \mathbb{P}(X_2 = x_2 \mid X_1 = x_1) \mathbb{P}(X_1 = x_1) \\ &= \mathbb{P}(X_n = x_n \mid X_{n-1} = x_{n-1}) \cdots \mathbb{P}(X_2 = x_2 \mid X_1 = x_1) \mathbb{P}(X_1 = x_1) \end{aligned}$$

Now, each of these probabilities is either p or $1 - p$. So, we can write

$$\mathbb{P}(X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1) = p^{N_U} (1 - p)^{n - N_U}.$$

where N_U is the number of times the random walk goes up. More precisely,

$$N_U = \#\{1 \leq i \leq n : x_i - x_{i-1} = 1\}$$

1.2.4 Distribution of X_n

One thing we might be interested in is the distribution of X_n , given we only know the initial value of the random walk, x_0 . We can make this a bit

simpler by thinking about the distribution of $X_n - x_0$ (this always starts at 0). Notice that X_n cannot be more than n steps from x_0 . That is, $|X_n - x_0| \leq n$. Also, if x_0 is even then X_n must be even if n is even and odd if n is odd. If x_0 is odd, it is the other way around.

Now, $X_n - x_0 = n$ can only happen if all n of the steps are up. There is only one possible path for which this is true. This path has probability p^n . Likewise, $X_n - x_0 = n - 2$ can only happen if all but one of the steps are up. Any path for which this is true has probability $p^{n-1}(1-p)^1$. However, there is more than one path that has $n - 1$ up steps and 1 down step (the first step could be the down one or the second step could be the down one, etc.). In fact, there are $\binom{n}{n-1}$ possible paths. Continuing with this pattern, we have

$$\begin{aligned} \mathbb{P}(X_n - x_0 = x) &= \begin{cases} \binom{n}{(n+x)/2} p^{(n+x)/2} (1-p)^{(n-x)/2} & \text{if } |x| \leq n \text{ and } x \text{ have the same parity as } n \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Having the same parity means both n and x are even or both n and x are odd.

1.2.5 First Passage Time

A quantity we are often very interested in is the *first passage time* of $\{X_n\}_{n \geq 0}$ into the set A .

Definition 1.2.4 (First Passage Time: Discrete Time Version). Given a stochastic process $\{X_n\}_{n \geq 0}$ we define the first passage time of $\{X_n\}_{n \geq 0}$ into the set A by

$$\tau_A = \begin{cases} \inf\{n \geq 0 : X_n \in A\} & \text{if } X_n \in A \text{ for some } n \geq 0 \\ \infty & \text{if } X_n \notin A \text{ for all } n \geq 0 \end{cases}.$$

Note that a first passage time is a random variable.

In general, for random walks, we want to calculate the first time the process hits a certain level (say 10). In this case, $A = \{10, 11, \dots\}$. In a random walk context, a first passage time could correspond to the first time a queue gets too long, or (in a very simplistic model) the first time a stock price hits a certain level.

It is possible to work out the distribution of the first passage time using a beautiful and simple piece of mathematics called the *reflection principle*. However, I will just give you the result.

Lemma 1.2.5 (First Passage Time Distribution for a Random Walk). The distribution of the first passage time for a random walk, $\{X_n\}_{n \geq 0}$ with $x_0 = 0$, where $a \neq 0$ and $A = \{a, a + 1, \dots\}$ if $a > 0$ or $A = \{a, a - 1, a - 2, \dots\}$ if $a < 0$, is given by

$$\mathbb{P}(\tau_A = n) = \begin{cases} \frac{|a|}{n} \mathbb{P}(X_n = a) & \text{if } a \text{ and } n \text{ have the same parity and } n \geq |a| \\ 0 & \text{otherwise} \end{cases}.$$

Working out first passage time probabilities for a random walk hitting a with $x_0 \neq 0$ is the same as working out first passage time probabilities for a random walk starting from 0 hitting $a - x_0$.

Example 1.2.6 (Probability a random walk returns to 0 in 4 steps). We can make everything a bit clearer with an example. What is the probability a random walker returns to zero in 4 steps? Our formula for first passage time probabilities only helps with levels other than 0. However, we can make it work by realising that, at the first step, the random walk must step up to 1 (with probability p) or step down to -1 (with probability $1 - p$). Now, the probability a random walk starting at 1 first hits 0 in 3 steps is the same as the probability a random walk starting at 0 first hits -1 in 3 steps and the probability a random walk starting at -1 first hits 0 in 3 steps is the same as the probability a random walk starting at 0 first hits 1 in 3 steps. So our answer is of the form

$$\begin{aligned} & p \mathbb{P}(\tau_{\{-1, -2, \dots\}} = 3) + (1 - p) \mathbb{P}(\tau_{\{1, 2, \dots\}} = 3) \\ &= p \frac{1}{3} \mathbb{P}(X_3 = -1) + (1 - p) \frac{1}{3} \mathbb{P}(X_3 = 1) \\ &= p \frac{1}{3} \binom{3}{1} p^1 (1 - p)^2 + (1 - p) \frac{1}{3} \binom{3}{2} p^2 (1 - p)^1 \\ &= p^2 (1 - p)^2 + (1 - p)^2 p^2 = 2p^2 (1 - p)^2 \end{aligned}$$

Chapter 2

Estimators

When we generate stochastic processes, we are usually interested in calculating actual numbers. These numbers are almost always in the form of an expectation. Sometimes we just want to know $\mathbb{E}X_n$ for some value of n but usually it is something slightly more complicated. Often we are interested in a probability, for example $\mathbb{P}(X_5 > 10)$ or $\mathbb{P}(\tau_A < 10)$. A probability can be written as an expectation using the indicator function. For example $\mathbb{P}(X_5 > 3) = \mathbb{E}\mathbb{I}(X_5 > 3)$ and $\mathbb{P}(\tau_A < 10) = \mathbb{E}\mathbb{I}(\tau_A < 10)$.

In general, we will not know how to calculate the expectations we are interested in explicitly. In the first few weeks we will practice by calculating expectations of things we know exactly, but soon the objects we are considering will become too complicated for us to do this. We will need to estimate things.

Monte Carlo methods are based on one very important result in probability theory: the *strong law of large numbers*. In simple terms, this says we can estimate an expected value using a sample (or empirical) average.

Theorem 2.0.7 (Strong Law of Large Numbers).

Let $\{X^{(i)}\}_{i \geq 0}$ be a sequence of i.i.d. random variables with values in \mathbb{R} such that $\mathbb{E}|X^{(1)}| \leq \infty$. Let S_N , $N > 0$, be the sample average defined by

$$S_N = \frac{1}{N} \sum_{i=1}^N X^{(i)}.$$

Then,

$$\lim_{N \rightarrow \infty} S_N = \mathbb{E}X^{(1)} \text{ almost surely.}$$

This means we can estimate $\ell = \mathbb{E}X$ by simulating a sequence of random variables with the same distribution as X and taking their average value.

That is, if we can produce a sequence of i.i.d. random variables $\{X^{(i)}\}_{i=1}^N$ with the same distribution as X , then the *estimator*

$$\widehat{\ell} = \frac{1}{N} \sum_{i=1}^N X^{(i)},$$

is approximately equal to ℓ for large enough N .

Example 2.0.8 (Estimating $\mathbb{P}(X_5 > 3)$ for a random walk). To estimate $\ell = \mathbb{P}(X_5 > 3) = \mathbb{E}\mathbb{I}(X_5 > 3)$ we can simulate N random walks until the 5th step. We then look at the N values at the 5th step, $X_5^{(1)}, \dots, X_5^{(N)}$, and check how many of them have a value bigger than 3. Our estimator is

$$\widehat{\ell} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(X_5^{(i)} > 3).$$

We can easily do this in Matlab.

Listing 2.1: Estimating $\mathbb{P}(X_5 > 3)$

```

1 N = 10^3;
2 n = 5; p = 0.5; X_0 = 0;
3
4 X = zeros(N,5);
5
6 Y = (rand(N,n) <= p);
7 X = X_0 + cumsum((2*Y - 1),2);
8 X_5 = X(:,5);
9
10 ell_est = mean(X_5 > 3)
```

In this simple example, we can calculate the probability exactly. It is only possible for X_5 to be greater than 3 if all the steps of the random walk are up. That is, $\ell = \mathbb{P}(X_5 > 3) = p^5$. In the case where $p = 0.5$, $\ell = 0.03125$. We can look at example output from our estimator: for $N = 10^1$, we have $\widehat{\ell} = 0$; for $N = 10^2$, we have $\widehat{\ell} = 0.04$; for $N = 10^3$, we have $\widehat{\ell} = 0.026$; for $N = 10^4$, we have $\widehat{\ell} = 0.0343$; and so on. On average, our estimator will become more accurate each time we increase N .

Example 2.0.9 (Estimating $\mathbb{E}\tau_{\{4,5,\dots\}}$). We can estimate expected hitting times as well. For example, we might want to estimate $\mathbb{E}\tau_{\{4,5,\dots\}}$ for a random walk with $x_0 = 0$. We should do this for a random walk with $p > (1 - p)$, so we can be certain the random walk hits 4. In Matlab, we would do the following.

Listing 2.2: Estimating $\mathbb{E}\tau_{\{4,5,\dots\}}$

```

1 N = 10^4;
2 n = 5; p = 0.8; X_0 = 0;
3
4 tau_4 = zeros(N,1);
5
6 for i = 1:N
7     X = X_0; n = 0;
8     while(X ~= 4)
9         X = X + 2*(rand <= p) - 1;
10        n = n + 1;
11    end
12    tau_4(i) = n;
13 end
14
15 est_expec_tau_4 = mean(tau_4)

```

Notice that we have to use a ‘while’ loop instead of a ‘for’ loop, as we do not know how long we will have to simulate $\{X_n\}_{n \geq 0}$ for.

2.1 Bias, Variance, the Central Limit Theorem and Mean Square Error

The strong law of large numbers tells that, if we use an infinite sample, a Monte Carlo estimator will give us the right answer (with probability 1). This is not really a practical result. In practice, we can only run a computer for a finite amount of time. What we are really interested in are the properties of our estimator for a fixed N .

It is comforting if our estimator on average gives the right answer. The *bias* of an estimator is the amount by which, on average, it deviates from the value it is supposed to estimate.

Definition 2.1.1 (Bias). The bias of an estimator, $\widehat{\ell}$, for the value ℓ is given by

$$\text{Bias}(\widehat{\ell}) = \mathbb{E}[\widehat{\ell} - \ell] = \mathbb{E}\widehat{\ell} - \ell.$$

An estimator without bias is called *unbiased*.

Definition 2.1.2 (Unbiased Estimator). We say an estimator, $\widehat{\ell}$, of ℓ is unbiased if $\text{Bias}(\widehat{\ell}) = 0$ or, equivalently,

$$\mathbb{E}\widehat{\ell} = \ell.$$

Many Monte Carlo estimators are unbiased. It is usually pretty easy to prove this.

Example 2.1.3 (Our estimator of $\mathbb{P}(X_5 > 3)$ is unbiased). Our estimator of $\ell = \mathbb{P}(X_5 > 3)$ is

$$\widehat{\ell} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(X_5^{(i)} > 3),$$

where the $\{X_5^i\}_{i=1}^N$ are i.i.d. Now,

$$\mathbb{E}\widehat{\ell} = \mathbb{E} \frac{1}{N} \sum_{i=1}^N \mathbb{I}(X_5^{(i)} > 3) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}\mathbb{I}(X_5^{(i)} > 3)$$

(we can justify exchanging expectation and summation using the fact that the sum is finite) and, because the $\{X_5^{(i)}\}_{i=1}^N$ are i.i.d. with the same distribution as X_5 ,

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E}\mathbb{I}(X_5^{(i)} > 3) = \frac{1}{N} N \mathbb{E}\mathbb{I}(X_5^{(1)} > 3) = \mathbb{E}\mathbb{I}(X_5 > 3) = \mathbb{P}(X_5 > 3) = \ell.$$

Because Monte Carlo estimators are random (they are averages of random variables), their errors (deviations from the true value) are random as well. This means we should use probabilistic concepts to describe these errors. For unbiased estimators, variance / standard deviation is an obvious choice of error measure. This is because variance is a measure of the deviation of a random variable from the mean of its distribution (which, for unbiased estimators, is the value which is being estimated). For an estimator of the form $\widehat{\ell} = \frac{1}{N} \sum_{i=1}^N X^{(i)}$, with the $\{X^{(i)}\}_{i=1}^N$ i.i.d., the variance is given by

$$\text{Var}(\widehat{\ell}) = \text{Var}\left(\frac{1}{N} \sum_{i=1}^N X^{(i)}\right) = \frac{1}{N^2} N \text{Var}(X^{(1)}) = \frac{1}{N} \text{Var}(X^{(1)}). \quad (2.1)$$

It should be clear that, as N gets larger, the variance of Monte Carlo estimators gets smaller. The standard deviation of the estimator is given by

$$\text{Std}(\widehat{\ell}) = \sqrt{\text{Var}(\widehat{\ell})} = \frac{\sqrt{\text{Var}(X^{(1)})}}{\sqrt{N}}.$$

Example 2.1.4 (The variance of our estimator of $\mathbb{P}(X_5 > 3)$). The variance of

$$\widehat{\ell} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(X_5^{(i)} > 3),$$

is given by

$$\text{Var}(\widehat{\ell}) = \frac{1}{N} \text{Var}(\mathbb{I}(X_5 > 3)).$$

Observe that $\mathbb{I}(X_5 > 3)$ takes the value 1 with probability $\ell = \mathbb{P}(X_5 > 3)$ and the value 0 with probability $1 - \ell$. That is, it is a Bernoulli random variable with parameter $p = \ell$. Now, the variance of a $\text{Ber}(p)$ random variable is given by $p(1 - p)$. So, $\text{Var}(\mathbb{I}(X_5 > 3)) = \ell(1 - \ell)$. In example 2.0.8 we observed that $\ell = p^5$. So,

$$\text{Var}(\widehat{\ell}) = \frac{\ell(1 - \ell)}{N} = \frac{p^5(1 - p^5)}{N}.$$

In order to calculate variances of the form (2.1), we need to know $\text{Var}(X^{(1)})$. Unfortunately, we usually do not know this. In the example above, calculating $\text{Var}(\mathbb{I}(X_5 > 3))$ required us to know ℓ , which was the very quantity we were trying to estimate. In practice, we usually estimate the variance (or, more meaningfully, standard deviation) instead. This is easy to do in Matlab.

Listing 2.3: Estimating the mean and variance of our estimator of $\mathbb{P}(X_5 > 3)$

```

1 N = 10^4;
2 n = 5; p = 0.5; X_0 = 0;
3
4 X = zeros(N,5);
5
6 Y = (rand(N,n) <= p);
7 X = X_0 + cumsum((2*Y - 1),2);
8 X_5 = X(:,5);
9
10 ell_est = mean(X_5 > 3)
11 var_est = var(X_5 > 3) / N
12 std_est = std(X_5 > 3) / sqrt(N)

```

Knowing the variance / standard deviation is very useful, because it allows us to make confidence intervals for estimators. This is because of an even more famous result in probability theory than the strong law of large numbers: the central limit theorem.

Theorem 2.1.5 (Central Limit Theorem). Let $\{X^{(i)}\}_{i=1}^N$ be a sequence of i.i.d. random values taking values in \mathbb{R} such that $\mathbb{E}(X^{(1)}) = \mu$ and $\text{Var}(X^{(1)}) = \sigma^2$, with $0 < \sigma^2 < \infty$. Then, the random variables

$$Z_n = \frac{\sum_{i=1}^N X^{(i)} - N\mu}{\sigma\sqrt{N}}$$

converge in distribution to a random variable $Z \sim \mathcal{N}(0, 1)$ as $N \rightarrow \infty$.

This is important because it implies that, for large enough N ,

$$\frac{1}{N} \sum_{i=1}^N X^{(i)} - \mathbb{E}X^{(1)}$$

is approximately Normally distributed with mean 0 and variance $\text{Var}(X^{(1)})/N$. As a result, we can make *confidence intervals* for our estimators. For example, a 95% confidence interval would be of the form

$$\left(\hat{\ell} - 1.96 \frac{\sqrt{\text{Var}(X^{(1)})}}{\sqrt{N}}, \hat{\ell} + 1.96 \frac{\sqrt{\text{Var}(X^{(1)})}}{\sqrt{N}} \right).$$

Sometimes, for whatever reason, we are unable to find an unbiased estimator, or our unbiased estimator is not very good. For biased estimators (and, arguably, for unbiased estimators) a good measure of the error is *mean squared error*.

Definition 2.1.6 (Mean Squared Error). The mean square error of the estimator, $\hat{\ell}$, of ℓ is given by

$$\text{MSE}(\hat{\ell}) = \mathbb{E} \left(\hat{\ell} - \ell \right)^2.$$

The mean square error can be decomposed into variance and bias terms.

Lemma 2.1.7. It holds true that $\text{MSE}(\hat{\ell}) = \text{Var}(\hat{\ell}) + \text{Bias}(\hat{\ell})^2$.

Proof. We have

$$\begin{aligned} \text{MSE}(\hat{\ell}) &= \mathbb{E} \left(\hat{\ell} - \ell \right)^2 = \mathbb{E} \hat{\ell}^2 - 2\ell \mathbb{E} \hat{\ell} + \ell^2 \\ &= \mathbb{E} \hat{\ell}^2 + \ell^2 - 2\ell \mathbb{E} \hat{\ell} + \left(\mathbb{E} \hat{\ell} \right)^2 - \left(\mathbb{E} \hat{\ell} \right)^2 \\ &= \mathbb{E} \hat{\ell}^2 - \left(\mathbb{E} \hat{\ell} \right)^2 + (\ell - \mathbb{E} \hat{\ell})^2 \\ &= \text{Var}(\hat{\ell}) + \text{Bias}(\hat{\ell})^2 \end{aligned}$$

□

2.2 Non-Asymptotic Error Bounds

Another way to measure the error of an estimator is by $|\ell - \hat{\ell}|$, the distance between the estimator and the value it is trying to estimate. Pretty obviously, we want this to be as small as possible. We can get bounds on this error using some famous inequalities from probability theory.

Theorem 2.2.1 (Markov's inequality). Given a random variable X taking values in \mathbb{R} , a function $g : \mathbb{R} \rightarrow [0, \infty)$ (that is, a function that never returns negative values), and $a > 0$, we have

$$\mathbb{P}(g(X) \geq a) \leq \frac{\mathbb{E}g(X)}{a}.$$

Proof. It is clear that $g(X) \geq a\mathbb{I}(g(X) \geq a)$, so $\mathbb{E}g(X) \geq \mathbb{E}a\mathbb{I}(g(X) \geq a) = a\mathbb{E}\mathbb{I}(g(X) \geq a) = a\mathbb{P}(g(X) \geq a)$. \square

If we set $g(x) = (x - \mathbb{E}X)^2$ and $a = \epsilon^2$, where $\epsilon > 0$, we have

$$\mathbb{P}((X - \mathbb{E}X)^2 \geq \epsilon^2) \leq \frac{(X - \mathbb{E}X)^2}{\epsilon^2} \Rightarrow \mathbb{P}(|X - \mathbb{E}X| \geq \epsilon) \leq \frac{\text{Var}(X)}{\epsilon^2}.$$

This is Chebyshev's inequality.

Theorem 2.2.2 (Chebyshev's inequality). Given a random variable X taking values in \mathbb{R} with $\text{Var}(X) < \infty$ and $\epsilon > 0$, we have

$$\mathbb{P}(|X - \mathbb{E}X| \geq \epsilon) \leq \frac{\text{Var}(X)}{\epsilon^2}.$$

Example 2.2.3 (Error Bounds on Probability Estimators). Consider the standard estimator of $\ell = \mathbb{P}(X > \gamma)$,

$$\widehat{\ell} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(X_i > \gamma).$$

We know the variance of this estimator is $N^{-1}\mathbb{P}(X > \gamma)(1 - \mathbb{P}(X > \gamma))$. So, we have the error bound

$$\mathbb{P}(|\widehat{\ell} - \ell| > \epsilon) \leq \frac{\mathbb{P}(X > \gamma)(1 - \mathbb{P}(X > \gamma))}{N\epsilon^2}.$$

2.3 Big O and Little o Notation

Although it is not always sensible to concentrate on how things behave asymptotically (for example, how an estimator behaves as $N \rightarrow \infty$), it is often difficult to get meaningful non-asymptotic results. We will use two types of asymptotic notation in this course. The first is what is called big O notation (or, sometimes, Landau notation).

Definition 2.3.1 (Big O). We say $f(x) = O(g(x))$, or f is of order $g(x)$, if there exists $C > 0$ and $x_0 > 0$ such that

$$|f(x)| \leq Cg(x)$$

for all $x \geq x_0$ as $x \rightarrow \infty$ (or, sometimes, as $x \rightarrow 0$).

Example 2.3.2. The quadratic $x^2 + 3x + 1$ is $O(x^2)$, as, for $x \geq x_0 = 1$, we have $x^2 + 3x + 1 \leq x^2 + 3x^2 + x^2 = 5x^2$, so $x^2 + 3x + 1 \leq Cx^2$ where $C = 5$.

If we can break a function into a sum of other functions (e.g., $f(x) = x^2 + 3x = f_1(x) + f_2(x)$, where $f_1(x) = x^2$ and $f_2(x) = 3x$), then the order of f is the order of the component function with the biggest order. In our example, $f_1(x) = O(x^2)$ and $f_2(x) = O(x)$, so $f(x) = O(x^2)$.

We use big O notation to describe the behavior of algorithms. If we measure the dimension of a problem by n — for example, the number of items in a list that we have to sort — then the work done by the algorithm will usually be a function of n . Usually, we prefer algorithms with smaller growth rates for the work they have to do. For example, if algorithm 1 accomplishes a task with $O(n)$ work and another algorithm needs $O(n^2)$ work, then we will tend to prefer algorithm 1. It is worth noting, however, that if the work done by algorithm 1 is $f_1(n) = 10^6 n$ and the work done by algorithm 2 is $f_2(n) = n^2$, then the second algorithm is actually better if $n < 10^6$, even though its order is worse.

It is worth noting that $x^2 = O(x^2)$, but it is also true that $x^2 = O(x^3)$, so the equals sign is something of an abuse of notation.

Example 2.3.3 (The Standard Deviations of the Monte Carlo Estimator). The standard deviations of the standard Monte Carlo estimator

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N X^{(i)},$$

is

$$\frac{1}{\sqrt{N}} \sqrt{\text{Var}(X^{(1)})} = O(N^{-1/2}).$$

Definition 2.3.4 (Small o Notation). We say $f(x) = o(g(x))$ if

$$\frac{f(x)}{g(x)} \rightarrow 0$$

as $x \rightarrow \infty$ (or, sometimes, as $x \rightarrow 0$).

Basically, $f(x) = o(g(x))$ means that $f(x)$ is growing more slowly than $g(x)$ as x gets large (or small).

Chapter 3

Markov Chains

The following material is closely based on the book “Markov Chains” by James Norris ([1]). This is really a wonderful book and well worth buying if you are interested in Markov chains.

As stated above, random walks are examples of Markov Chains, discrete time discrete state space stochastic processes with the Markov property. While random walks can only move up or down by 1 at each step, there is no such restriction on Markov chains in general.

On a finite state space (i.e., $|\mathcal{X}| < \infty$) a Markov chain can be represented by a transition matrix, P . The element in the i th row and j th column, $P_{i,j}$, describes the probability of going from state i to state j in one step. That is $P_{i,j} = \mathbb{P}(X_1 = j \mid X_0 = i)$. We will always work with homogenous Markov chains (that is, the transition probabilities will never depend on n), so we have that $P_{i,j} = \mathbb{P}(X_1 = j \mid X_0 = i) = \mathbb{P}(X_{n+1} = j \mid X_n = i)$ for all $n \geq 0$.

Example 3.0.5. Consider the Markov chain with the following graphical representation (the nodes are the states and the arrows represent possible transitions, with the probabilities attached).

PICTURE HERE

We can write this in matrix form as

$$P = \begin{bmatrix} 0 & 1/3 & 1/3 & 1/3 \\ 0 & 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 & 1/2 \\ 1/2 & 0 & 1/2 & 0 \end{bmatrix}$$

The convention when working with Markov chains is to describe probability distributions using row vectors. So, for example, $\boldsymbol{\mu} = (1/4, 1/4, 1/4, 1/4)$ would be a possible probability distribution for 4 states.

The initial state of the Markov chain, x_0 , will be given by a probability distribution. I will try to use λ for this. So, for example, $\mathbb{P}(X_0 = 1) = \lambda_1$. When the Markov chain starts at a fixed state, x_0 , this will be represented by the distribution δ_i which is 0 for all states except the i th one, where it is 1. Because the Markov property tells us that it is sufficient to know the current state in order to calculate the probability of the next state, the transition probability matrix P and the initial distribution λ fully specify the Markov chain. We will call a Markov chain a (P, λ) Markov chain if it has transition probability matrix P and the initial distribution λ .

The path probabilities of a Markov chain are straightforward to calculate. We have

$$\begin{aligned} & \mathbb{P}(X_n = x_n, X_{n-1} = x_{n-1}, \dots, X_1 = x_1, X_0 = x_0) \\ &= \mathbb{P}(X_n = x_n | X_{n-1} = x_{n-1}) \cdots \mathbb{P}(X_1 = x_1 | X_0 = x_0) \mathbb{P}(X_0 = x_0) \\ &= P_{x_{n-1}, x_n} P_{x_{n-2}, x_{n-1}} \cdots P_{x_0, x_1} \lambda_{x_0} \end{aligned}$$

We can also work out the distribution of X_n quite easily. Labeling the states from $1, \dots, L$, the probability that we are in state 1 in the first step is given by

$$\mathbb{P}(X_1 = 1) = \lambda_1 P_{1,1} + \lambda_2 P_{2,1} + \cdots + \lambda_L P_{L,1}.$$

Likewise,

$$\mathbb{P}(X_1 = 2) = \lambda_1 P_{1,2} + \lambda_2 P_{2,2} + \cdots + \lambda_L P_{L,2}.$$

and,

$$\mathbb{P}(X_1 = L) = \lambda_1 P_{1,L} + \lambda_2 P_{2,L} + \cdots + \lambda_L P_{L,L}.$$

It is easy to see, in fact, that the distribution

$$(\mathbb{P}(X_1 = 1), \mathbb{P}(X_1 = 2), \dots, \mathbb{P}(X_1 = L))$$

is given by λP . If we consider $\mathbb{P}(X_2 = x_2)$, we have

$$\mathbb{P}(X_2 = 1) = \mathbb{P}(X_1 = 1)P_{1,1} + \mathbb{P}(X_1 = 2)P_{2,1} + \cdots + \mathbb{P}(X_1 = L)P_{L,1},$$

and so on. This implies that the distribution of X_2 is given by λP^2 . If we keep on going, we have the general result that, for a (P, λ) Markov chain,

$$\mathbb{P}(X_n = j) = (\lambda P^n)_j.$$

We call P^n the *n-step transition matrix*.

Example 3.0.6 (Calculating the Distribution of X_n). Consider the Markov chain from example 3.0.5, with transition matrix

$$P = \begin{bmatrix} 0 & 1/3 & 1/3 & 1/3 \\ 0 & 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 & 1/2 \\ 1/2 & 0 & 1/2 & 0 \end{bmatrix}.$$

Given an initial distribution, $\lambda = (1/4, 1/4, 1/4, 1/4)$, we can calculate the distribution of X_n in Matlab as follows.

Listing 3.1: Finding the Distribution of X_n

```

1 n = 2;
2 P = [0 1/3 1/3 1/3; 0 1/2 1/2 0; 1/2 0 0 1/2; 1/2 0 1/2 0];
3 lambda = [1/4 1/4 1/4 1/4];
4
5 X_n_dist = lambda * P^n

```

For $n = 1$, we get

$$\lambda P = (0.2500, 0.2083, 0.3333, 0.2083).$$

For $n = 2$, we have

$$\lambda P^2 = (0.2500, 0.2083, 0.3333, 0.2083).$$

For $n = 20$, we have

$$\lambda P^{20} = (0.2727, 0.1818, 0.3030, 0.2424).$$

And, for $n = 1000$, we have

$$\lambda P^{1000} = (0.2727, 0.1818, 0.3030, 0.2424).$$

If we use the initial distribution $\lambda = (1, 0, 0, 0)$, then, for $n = 1$, we have

$$\lambda P = (0, 0.3333, 0.3333, 0.3333).$$

For $n = 2$, we have

$$\lambda P^2 = (0.3333, 0.1667, 0.3333, 0.1667).$$

For $n = 20$, we have

$$\lambda P^{20} = (0.2727, 0.1818, 0.3030, 0.2424).$$

And, for $n = 1000$, we have

$$\lambda P^{1000} = (0.2727, 0.1818, 0.3030, 0.2424).$$

Notice that the distributions appears to converge to the same distribution regardless of the choice of initial distribution.

When we define Markov chains on infinite (but countable) state spaces, we can still keep lots of the formalism from the finite state space setting. Because the state space is countable, it still makes sense to talk about $P_{i,j}$. However, now the matrix is infinite, so calculating values like P^n may be a bit more difficult (we cannot just enter the matrix into Matlab and take the n th power).

Example 3.0.7 (Transition Matrix for a Random Walk). For a random walk, we have the transition matrix $(P_{i,j})_{i \in \mathbb{Z}, j \in \mathbb{Z}}$, where

$$P_{i,j} = \begin{cases} p & \text{if } j = i + 1 \\ (1 - p) & \text{if } j = i - 1 \\ 0 & \text{otherwise} \end{cases}.$$

3.1 Simulating Markov Chains

We need some basic techniques for drawing from discrete distributions before we can start simulating Markov chains.

3.1.1 Drawing from a Discrete Uniform Distribution

We can simulate a random variables from the discrete uniform distribution on $\{1, \dots, L\}$ (i.e., $\boldsymbol{\mu} = (1/L, \dots, 1/L)$) by observing that if $U \sim \mathcal{U}(0, 1)$, then

$$\begin{aligned} \mathbb{P}(\lceil LU \rceil = 1) &= \mathbb{P}(LU \leq 1) = \mathbb{P}(U \leq 1/L) = 1/L, \\ \mathbb{P}(\lceil LU \rceil = 2) &= \mathbb{P}(1 < LU \leq 2) = \mathbb{P}(1/L < U \leq 2/L) = 1/L. \end{aligned}$$

and so on. This suggests that $\lceil LU \rceil$ is a random variable distributed uniformly on $\{1, \dots, L\}$.

Listing 3.2: Drawing uniformly from L values.

```
1 X = ceil(L * rand);
```

3.1.2 Drawing From A Discrete Distribution on a Small State Space

3.1.3 Simulating a Markov Chain

It is pretty straightforward to simulate finite state space Markov chains (provided that the state space is not too big).

Algorithm 3.1.1 (Simulating a Markov Chain).

- (i) Draw X_0 from λ . Set $i = 1$.
- (ii) Set $X_{i+1} = j$ with probability $P_{X_i,j}$.
- (iii) Set $i = i + 1$. If $i < n$ repeat from step 2.

Example 3.1.1. We can simulate the following Markov chain, where $\lambda = (1/3, 1/3, 1/3)$ and

$$P = \begin{bmatrix} 1/2 & 1/4 & 1/4 \\ 0 & 0 & 1 \\ 2/3 & 1/3 & 0 \end{bmatrix}.$$

Listing 3.3: Matlab Code

```

1 n = 10^3; X = zeros(n,1);
2 X(1) = ceil(rand*3); i =1;
3 P = [1/2 1/4 1/4; 0 0 1; 2/3 1/3 0];
4 while i<n
5     X(i+1) = min(find(rand<cumsum(P(X(i),:)))));
6     i = i+1;
7 end

```

Unfortunately, in the case of Markov chains with big (or infinite) state spaces, this approach does not work. However, there is often another way to simulate such Markov chains. We have already seen one such example for random walks.

3.2 Communication

It is very often the case that the distribution of X_n settles down to some fixed value as n grows large. We saw this in example 3.0.6. In order to talk about limiting behavior (and stationary distributions, which are closely related) we need the transition matrix to possess certain properties. For this reason, we introduce a number of definitions. We say that state i *leads to* state j , written $i \rightarrow j$, if

$$\mathbb{P}_i(X_n = j \text{ for some } n \geq 0) > 0.$$

That is, if it is possible to get to j from i (though not necessarily in a single step). Assuming that $\lambda_i > 0$, $\mathbb{P}_i(A) = \mathbb{P}(A | X_0 = i)$. If no λ is specified, then assume that $\lambda_i = 1$.

Definition 3.2.1 (Communication). We say two states $i, j \in \mathcal{X}$ *communicate*, denoted $i \leftrightarrow j$, if $i \rightarrow j$ and $j \rightarrow i$.

Obviously, communicating is reflexive. That is, $i \rightarrow i$ for all $i \in \mathcal{X}$. We can partition the state space \mathcal{X} into *communicating classes* as follows. We say i and j are in the communicating class C if they communicate.

Example 3.2.2 (Communicating Classes). Consider the Markov chain with the following graphical representation.

PICTURE HERE

There are two communicating classes $C_1 = \{1, 2\}$ and $C_2 = \{3\}$. In full, we have the following relationships. $1 \leftrightarrow 2$, $1 \rightarrow 3$ and $2 \rightarrow 3$.

Definition 3.2.3 (Irreducibility). We say a transition matrix P is *irreducible* if \mathcal{X} is a single communicating class (if it is always possible to get from one state to another, though not always in just 1 step).

3.3 The Strong Markov Property

An important and (slightly) stronger property than the normal Markov property is the *strong Markov property*. To introduce this, we give our first definition of a stopping time.

Definition 3.3.1 (Stopping Time: Discrete Time Version). A random variable $\tau \rightarrow \{0, 1, \dots\} \cup \{\infty\}$ is a stopping time if the event $\{\tau = n\}$ only depends on X_0, X_1, \dots, X_n , for $n \geq 0$.

Basically, in order for something to be a stopping time, we have to decide if it has happened or not without knowing the future. So, for example $\tau = \inf\{n \geq 0 : X_n = 1\}$ is a stopping time, because we can tell when it occurs without seeing into the future. On the other hand, $\tau = \sup\{n \geq 0 : X_n = 1\}$ is not a stopping time for a random walk (we cannot tell when it happens unless we know the future). Another example of something that is not a stopping time is $\tau = \inf\{n \geq 0 : X_{n+1} = i\}$.

A stopping time that we have already encountered was the first passage time for the random walk.

Theorem 3.3.2 (Strong Markov Property). Let $\{X_n\}_{n \geq 0}$ be a $(\mathbf{\lambda}, P)$ Markov chain and let τ be a stopping time of $\{X_n\}_{n \geq 0}$. Then, conditional on $\tau < \infty$ and $X_\tau = i$, $\{X_{\tau+n}\}_{n \geq 0}$ is a (δ_i, P) Markov chain (δ_i is a distribution with 1 at state i and 0 everywhere else) and is independent of X_0, X_1, \dots, X_τ .

Proof. See [1] for a proof. □

3.4 Recurrence and Transience

All states in a Markov chain have the property of being either *recurrent* or *transient*.

Definition 3.4.1 (Recurrent). Given a Markov chain $\{X_n\}_{n \geq 0}$ with transition matrix P , we say that a state i is recurrent if

$$\mathbb{P}_i(X_n = i \text{ for infinitely many } n) = 1.$$

Definition 3.4.2 (Transient). Given a Markov chain $\{X_n\}_{n \geq 0}$ with transition matrix P , we say that a state i is transient if

$$\mathbb{P}_i(X_n = i \text{ for infinitely many } n) = 0.$$

In order to establish some important facts about recurrence and transience, we need the following result about expectations of non-negative integer valued random variables (random variables taking values in the natural numbers).

Theorem 3.4.3. Given a random variable X taking values in \mathbb{N} ,

$$\mathbb{E}X = \sum_{i=0}^{\infty} \mathbb{P}(X > i).$$

Proof. We have

$$\begin{aligned} \mathbb{E}X &= \sum_{x=1}^{\infty} x \mathbb{P}(X = x) = \sum_{x=1}^{\infty} \sum_{i=0}^{x-1} \mathbb{P}(X = x) \\ &= \sum_{i=0}^{\infty} \sum_{x=i+1}^{\infty} \mathbb{P}(X = x) = \sum_{i=0}^{\infty} \mathbb{P}(X > i), \end{aligned}$$

where we can swap sums because of the non-negative summands. □

We also need to introduce a stopping time.

Definition 3.4.4 (First Passage Time Including Return). We define the first passage time into the state i (including the possibility of starting in i) as

$$\tilde{\tau}_i = \inf\{n \geq 1 : X_n = i\}.$$

We also define $m_i = \mathbb{E}_i \tilde{\tau}_i$, where \mathbb{E}_i denotes the expectation given the Markov chain starts in state i .

Theorem 3.4.5. The following holds

- (i) If $\mathbb{P}_i(\tilde{\tau}_{\{i\}} < \infty) = 1$ then i is recurrent and $\sum_{n=0}^{\infty} P_{i,i}^n = \infty$.
- (ii) If $\mathbb{P}_i(\tilde{\tau}_{\{i\}} < \infty) < 1$ then i is transient and $\sum_{n=0}^{\infty} P_{i,i}^n < \infty$.

Proof. See [1] for a proof. □

Theorem 3.4.6. Let C be a communicating class. Then either all states in C are transient or all are recurrent.

Proof. Take a pair of states i and j in C and assume i is transient. Because i and j are in the same communicating class, there must exist $n \geq 0$ and $m \geq 0$ so that $P_{i,j}^n > 0$ and $P_{j,i}^m > 0$. Now, it must be the case that

$$P_{i,i}^{n+r+m} \geq P_{i,j}^n P_{j,j}^r P_{j,i}^m$$

as this only describes the probability of one possible path from i back to i (such a path need not pass through j). Rearranging, we have

$$P_{j,j}^r \leq \frac{P_{i,i}^{n+r+m}}{P_{i,j}^n P_{j,i}^m}.$$

Summing over r we get

$$\sum_{r=0}^{\infty} P_{j,j}^r \leq \frac{\sum_{r=0}^{\infty} P_{i,i}^{n+r+m}}{P_{i,j}^n P_{j,i}^m}.$$

Now, because i is assumed to be transient, $\sum_{r=0}^{\infty} P_{i,i}^{n+r+m} < \infty$. As a result, $\sum_{r=0}^{\infty} P_{j,j}^r < \infty$, implying j is also transient. Thus, the only way a state can be recurrent is if all states are recurrent. □

Definition 3.4.7 (Closed Class). A communicating class C is closed if $i \in C$ and $i \rightarrow j$ implies $j \in C$.

Theorem 3.4.8. Every finite closed class is recurrent.

Proof. See [1]. □

As an irreducible Markov chain consists of one single closed class, this implies that all irreducible Markov chains on finite state spaces are recurrent.

3.4.1 Recurrence of Random Walks

We can use the criteria given in theorem 3.4.5 to establish facts about the recurrence properties of random walks. In order to establish recurrence, we would need $\sum_{n=0}^{\infty} P_{0,0}^n = \infty$. Thus, a first step is to find an expression for $P_{0,0}^n$. As we have already discussed, a random walk starting at 0 can only return to 0 in an even number of steps. Thus, it is sufficient for us to consider $P_{0,0}^{2n}$. Now in order for a random walk to end up at 0 after $2n$ steps, it needs to take exactly n up steps and n down steps. There are $\binom{2n}{n}$ ways of taking this many steps. This gives

$$P_{0,0}^{2n} = \binom{2n}{n} p^n (1-p)^n.$$

Lemma 3.4.9. The following bounds on $n!$ hold of $n \geq 1$.

$$\sqrt{2\pi} n^{n+1/2} e^{-n} \leq n! \leq e n^{n+1/2} e^{-n}$$

Given these bounds on $n!$ it is relatively straightforward to get bounds on $P_{0,0}^n$ that allow us to establish the transience or recurrence of the random walk.

Theorem 3.4.10. A symmetric random walk (i.e., a random walk with $p = (1-p) = 1/2$) is recurrent.

Proof. To show the random walk is recurrent, we need to show

$$\sum_{n=0}^{\infty} P_{0,0}^{2n} = \infty.$$

The first step is to get a bound on $\binom{2n}{n}$. By lemma 3.4.9, we have

$$\binom{2n}{n} = \frac{(2n)!}{n!n!} \geq \frac{\sqrt{2\pi} 2n^{2n+1/2} e^{-2n}}{e n^{n+1/2} e^{-n} e n^{n+1/2} e^{-n}} = \frac{2\sqrt{\pi} 4^n}{e^2 \sqrt{n}}.$$

So, we can bound the probabilities from below by

$$P_{0,0}^{2n} \geq \frac{2\sqrt{\pi}}{e^2} \frac{(4p(1-p))^n}{\sqrt{n}}.$$

This implies that

$$\sum_{n=0}^{\infty} P_{0,0}^{2n} \geq \frac{2\sqrt{\pi}}{e^2} \sum_{n=0}^{\infty} \frac{(4p(1-p))^n}{\sqrt{n}}.$$

If $p = (1-p)$ then $4p(1-p) = 1$, so

$$\sum_{n=0}^{\infty} P_{0,0}^{2n} \geq C \sum_{n=0}^{\infty} \frac{1}{\sqrt{n}} = \infty,$$

where $C > 0$ is a constant. □

3.5 Invariant Distributions

A topic of enormous interest, especially from a simulation perspective, is the behavior of the distribution of a Markov chain, $\{X_n\}_{n \geq 0}$ as n becomes large. It turns out there are actually a few different things we can mean by this. The best possible situation is that a Markov chain can have a *limiting distribution*. That is, there can exist a distribution, π , such that $\{X_n\}_{n \geq 0} \rightarrow \pi$ as $n \rightarrow \infty$ no matter what initial distribution, λ , we choose. We will discuss the conditions for a limiting distribution to exist later. It turns out, however, that even if a limiting distribution does not exist, the Markov chain may have a unique *stationary distribution*, sometimes also called an *invariant distribution* or an *equilibrium distribution*. This is a distribution, π , such that $\pi P = \pi$.

Definition 3.5.1 (Stationary Distribution). An invariant distribution of a Markov chain with transition matrix P is a distribution that satisfies

$$\pi P = \pi.$$

The importance of stationary distributions is made clear by the following lemma.

Lemma 3.5.2. Let $\{X_n\}_{n \geq 0}$ be Markov (π, P) and suppose π is a stationary distribution for P , then $X_n \sim \pi$ for all $n \geq 0$.

Proof. The distribution of X_n is given by πP^n . Now, for $n = 0$ (where we define $P^0 = I$) it is clearly true that $\pi P^0 = \pi$. We just need to show $\pi P^n = \pi$ for $n > 1$. We do this by induction. That is, assume $\pi P^n = \pi$. Then

$$\pi P^{n+1} = (\pi P) P^n = \pi P^n = \pi$$

by assumption. Thus, as this is true for $n = 0$, the result follows. \square

The obvious questions to ask when faced with a nice and interesting object like a stationary distribution are: does one exist? and, if so, is it unique?. The simple answer is that it is often the case that there is a unique stationary distribution (at least for many of the Markov chains we might be interested in), but that a number of conditions need to be fulfilled. One of these is *positive recurrence*.

Recurrent implies that if the Markov chain starts in state i it will return to state i with probability 1. However, this is no guarantee that the expected return time is finite. Recall that $m_i = \mathbb{E}_i \tilde{\tau}_i$, where $\tilde{\tau}_i = \inf\{n > 0 : X_n = i\}$.

Definition 3.5.3 (Positive Recurrence). We say a Markov chain is positive recurrent if $m_i < \infty$.

Definition 3.5.4 (Null Recurrence). We say a Markov chain is *null recurrent* if $m_i = \infty$.

It is straightforward to establish that if a Markov chain is irreducible and has a finite state space (i.e. $|\mathcal{X}| < \infty$), then every state is positive recurrent. It is not always easy to show that an infinite state space Markov chain has such a property.

Now that we have defined positive recurrence, we are able to give conditions that guarantee a Markov chain has a unique stationary distribution.

Theorem 3.5.5. Let P be irreducible. Then the following are equivalent:

- (i) Every state is positive recurrent.
- (ii) Some state i is positive recurrent.
- (iii) P has a (unique) invariant distribution, $\boldsymbol{\pi}$, and $m_i = 1/\pi_i$.

Proof. See [1]. □

Note that it is possible for a Markov chain to have a stationary distribution but not satisfy these conditions.

Example 3.5.6 (Calculating a stationary distribution). We can find a stationary distribution by solving $\boldsymbol{\pi}P = \boldsymbol{\pi}$. For example, consider the following Markov chain

$$P = \begin{bmatrix} 0 & 1/2 & 1/2 \\ 1/2 & 0 & 1/2 \\ 1/2 & 1/2 & 0 \end{bmatrix}.$$

So, we need to solve

$$(\pi_1, \pi_2, \pi_3) \begin{bmatrix} 0 & 1/2 & 1/2 \\ 1/2 & 0 & 1/2 \\ 1/2 & 1/2 & 0 \end{bmatrix} = (\pi_1, \pi_2, \pi_3). \quad (3.1)$$

This gives the following linear system

$$\begin{aligned} \frac{1}{2}\pi_2 + \frac{1}{2}\pi_3 &= \pi_1 \\ \frac{1}{2}\pi_1 + \frac{1}{2}\pi_3 &= \pi_2 \\ \frac{1}{2}\pi_1 + \frac{1}{2}\pi_2 &= \pi_3. \end{aligned}$$

Solving this, we have $\pi_1 = \pi_2 = \pi_3$. As we need $\pi_1 + \pi_2 + \pi_3 = 1$ if $\boldsymbol{\pi}$ is to be a distribution, we have $\pi_1 = \pi_2 = \pi_3 = 1/3$.

3.6 Limiting Distribution

The fact that a Markov chain has a unique stationary distribution does not guarantee that it has a limiting distribution. This can be shown using a simple example.

Example 3.6.1 (A Markov chain with a stationary distribution but no limiting distribution). Consider the Markov chain with graph

PICTURE HERE.

The transition matrix is given by

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

This chain is clearly irreducible and positive recurrent, so it has a unique stationary distribution, $\pi = (1/2, 1/2)$. However, $P^{2n} = I$ and $P^{2n+1} = P$. This means that if we start with certainty in a given state the distribution will not converge. For example, if $\lambda = \delta_i = (1, 0)$, then $\lim_{n \rightarrow \infty} \lambda P^{2n} = (1, 0)$ and $\lim_{n \rightarrow \infty} \lambda P^{2n+1} = (0, 1)$. Thus, no limit exists.

The problem with the Markov chain in example 3.6.1 is that it is periodic: it is only possible to get from state 1 back to state 1 in an even number of steps. In order for a Markov chain to have a limiting distribution, it should not have any periodicity. Unsurprisingly, this requirement is called *aperiodicity*.

Definition 3.6.2. A state is said to be aperiodic if $P_{i,i}^n > 0$ for all sufficiently large n . Equivalently, the set $\{n \geq 0 : P_{i,i}^n > 0\}$ has no common divisor other than 1.

In a irreducible Markov chain, all states are either aperiodic or periodic.

Theorem 3.6.3. Suppose P is irreducible and has an aperiodic state i . Then, all states are aperiodic.

A random walk is periodic with period 2.

Example 3.6.4.

PICTURE HERE.

PICTURE HERE.

If a Markov chain is irreducible, positive recurrent and aperiodic, then it has a limiting distribution.

Theorem 3.6.5. Let P be irreducible and aperiodic and suppose P has an invariant distribution π . Let λ be any distribution. and suppose $\{X_n\}_{n \geq 0}$ is Markov (λ, P) . Then, $\mathbb{P}(X_n = j) \rightarrow \pi_j$ as $n \rightarrow \infty$ for all $j \in \mathcal{X}$. In particular, $P_{i,j}^n \rightarrow \pi_j$ as $n \rightarrow \infty$ for all $i, j \in \mathcal{X}$.

Proof. See [1]. □

3.7 Reversibility

Stationary distributions play a number of very important roles in simulation. In particular, we often wish to create a Markov chain that has a specified stationary distribution. One important tool for finding such a Markov chain is to exploit a property called *reversibility*. Not all Markov chains are reversible but, as we shall see, those that are have some nice properties.

Definition 3.7.1. Let $\{X_n\}_{n \geq 0}$ be a Markov (π, P) with P irreducible. We say that $\{X_n\}_{n \geq 0}$ is reversible if, for all $N \geq 1$, $\{X_{N-n}\}_{n=0}^N$ is also Markov (π, P) .

The most important property of reversible Markov chains is that they satisfy the *detailed balance equations*.

Definition 3.7.2. A matrix P and a measure ν (a vector with all non-negative components) are in detailed balance if

$$\nu_i P_{i,j} = \nu_j P_{j,i} \quad \forall i, j \in \mathcal{X}.$$

The detailed balance equations are important because they allow us to find stationary distributions (and, as we will learn later, construct transition matrices with given stationary distributions).

Lemma 3.7.3. If P and the distribution π are in detailed balance, then π is invariant for P .

Proof. We have

$$(\pi P)_i = \sum_{j \in \mathcal{X}} \pi_j P_{j,i} = \sum_{j \in \mathcal{X}} \pi_i P_{i,j} = \pi_i.$$

□

The following theorem shows that reversibility is equivalent to satisfying the detailed balance equations.

Theorem 3.7.4. Let P be an irreducible stochastic matrix and let $\boldsymbol{\pi}$ be a distribution. Suppose $\{X_n\}_{n \geq 0}$ is Markov $(\boldsymbol{\pi}, P)$. Then, the following are equivalent.

- (i) $\{X_n\}_{n \geq 0}$ is reversible.
- (ii) P and $\boldsymbol{\pi}$ are in detailed balance (i.e., $\boldsymbol{\pi}$ is the stationary distribution of P).

Proof. See [1].

□

Chapter 4

The Poisson Process

The *Poisson process* is one of the fundamental stochastic processes in probability. We can use it to build many complex and interesting stochastic objects. It is a very simple model for processes such as the arrival of people in a queue, the number of cars arriving at a red traffic light, and the occurrence of insurance claims.

4.1 Point Processes on $[0, \infty)$

A Poisson process is a *point process* on $[0, \infty)$. We will not consider such point processes in their most general form, but rather a straightforward and easy to work with subset of point processes. We can use a number of equivalent definitions. The first of these is to see the point process in terms of a counting process. A counting process, $\{N_t\}_{t \geq 0}$ counts the number of events that have happened by time t .

Definition 4.1.1 (Point Process on $[0, \infty)$: Counting Process Version). A point process on $[0, \infty)$ is a process $\{N_t\}_{t \in [0, \infty)}$ taking values in \mathbb{N} that:

- (i) Vanishes at 0. That is, $N_{0-} = N_0 = 0$, where $N_{t-} = \lim_{u \uparrow t} N_u$.
- (ii) Is non-decreasing. That is, $N_0 \leq N_s \leq N_t$ for $0 \leq s \leq t$.
- (iii) Is right continuous. That is, $N_t = N_{t+}$, where $N_{t+} = \lim_{u \downarrow t} N_u$.
- (iv) Has unit jumps. That is, $N_t - N_{t-} \in \{0, 1\}$. Technically, a process where only one jump can occur at a given time t is called a *simple point process*.
- (v) Has an infinite limit. That is, $\lim_{t \rightarrow \infty} N_t = \infty$.

Note, again, that some of these requirements can be relaxed.

Definition 4.1.2 (Point Process on $[0, \infty)$: Jump Instances Version). A point process on $[0, \infty)$ is defined by a sequence $\{T_n\}_{n \geq 1}$ of random variables that are positive and increasing to infinity. That is,

$$0 < T_1 < T_2 < \cdots < \infty \quad \text{and} \quad \lim_{n \rightarrow \infty} T_n = \infty.$$

This defines the counting process

$$N_t = \sum_{n \geq 1} \mathbb{I}(T_n \leq t) = \sup\{n \geq 1 : T_n \leq t\}.$$

Definition 4.1.3 (Point Process on $[0, \infty)$: Inter-arrivals Version). A point process on $[0, \infty)$ is defined by a sequence $\{S_n\}_{n \geq 1}$ of positive random variables such that $\sum_{n \geq 1} S_n = \infty$. These define a sequence of jump instances by $S_1 = T_1$ and $S_n = T_n - T_{n-1}$ for $n \geq 2$.

PICTURE HERE

These three definitions of point processes suggest three possible ways to simulate them.

- (i) We can simulate the counting process $\{N_t\}_{t \geq 0}$ (or its increments).
- (ii) We can simulate the jump times $\{T_n\}_{n \geq 1}$.
- (iii) We can simulate the inter-arrival times $\{S_n\}_{n \geq 1}$.

The key properties of a Poisson process (aside from being a point process) are that it has *stationary increments* and *independent increments*.

Definition 4.1.4 (Stationary Increments). A stochastic process $\{X_t\}_{t \geq 0}$ has stationary increments if the distribution of $X_{t+h} - X_t$ depends only on h for $h \geq 0$.

Definition 4.1.5 (Independent Increments). A stochastic process $\{X_t\}_{t \geq 0}$ has independent increments if the random variables $\{X_{t_{i+1}} - X_{t_i}\}_{i=1}^n$ are independent whenever $0 \leq t_1 < t_2 < \cdots < t_n$ and $n \geq 1$.

A process that has stationary and independent increments is attractive from a simulation perspective because we can simulate it in 'parts' (the increments).

4.2 Poisson Process

Equipped with the necessary definitions, we can now define a Poisson process.

Definition 4.2.1 (Poisson Process on $[0, \infty)$). A (homogenous) Poisson process on $[0, \infty)$ with parameter $\lambda > 0$ is a point process on $[0, \infty)$ with stationary and independent increments and $N_t \sim \text{Poi}(\lambda t)$ for all $t \geq 0$. That is,

$$\mathbb{P}(N_t = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}.$$

Actually, we can define a Poisson process in a number of different ways.

Theorem 4.2.2 (Poisson process). The following definitions are equivalent definitions.

- (i) A Poisson process is a point process, $\{N_t\}_{t \geq 0}$, with stationary and independent increments with $N_t \sim \text{Poi}(\lambda t)$ for all $t \geq 0$.
- (ii) A Poisson process is a point process, $\{N_t\}_{t \geq 0}$, with independent increments and the property that, as $h \downarrow 0$, uniformly in t
 - (a) $\mathbb{P}(N_{t+h} - N_t = 0) = 1 - \lambda h + o(h)$.
 - (b) $\mathbb{P}(N_{t+h} - N_t = 1) = \lambda h + o(h)$.
 - (c) $\mathbb{P}(N_{t+h} - N_t > 1) = o(h)$.
- (iii) A Poisson process is a point process defined by its inter-arrival times, $\{S_n\}_{n \geq 1}$, which are i.i.d. $\text{Exp}(\lambda)$.

Before we prove anything, we need a few results.

Lemma 4.2.3 (Memoryless Property). We say that exponential random variables have the *memoryless property*. That is, for $t, s \geq 0$, if $X \sim \text{Exp}(\lambda)$, then

$$\mathbb{P}(X > t + s \mid X > s) = \mathbb{P}(X > t).$$

Proof. We can write $\mathbb{P}(X > t + s \mid X > s)$ as

$$\frac{\mathbb{P}(X > t + s, X > s)}{\mathbb{P}(X > s)}.$$

As $s < t$ then $\mathbb{P}(X > t + s, X > s) = \mathbb{P}(X > t)$, so

$$\frac{\mathbb{P}(X > t + s, X > s)}{\mathbb{P}(X > s)} = \frac{\mathbb{P}(X > t + s)}{\mathbb{P}(X > s)} = \frac{e^{-\lambda(t+s)}}{e^{-\lambda s}} = e^{-\lambda t} = \mathbb{P}(X > t).$$

□

Theorem 4.2.4 (Markov Property). If $\{N_t\}_{t \geq 0}$ is a Poisson process with rate $\lambda > 0$, then, for any $s > 0$, $\{N_{t+s} - N_s\}_{t \geq 0}$ is also a Poisson process with rate λ . Furthermore, this process is independent of $\{N_r\}_{r \leq s}$.

Proof. First, note that the event $\{N_s = i\}$ can be written as

$$\{N_s = i\} = \{T_i \leq s\} \cap \{S_{i+1} > s - T_i\}$$

and, given this, for $r \leq s$,

$$N_r = \sum_{j=1}^i \mathbb{I}(S_j \leq r).$$

Define the process starting at time s by $\tilde{N}_t = N_{t+s} - N_s$. Then, given $\{N_s = i\}$, $\tilde{S}_1 = S_{i+1} - (s - T_i)$ and $\tilde{S}_n = S_{i+n}$ for $n \geq 2$. Now, by the memoryless property, \tilde{S}_1 is an $\text{Exp}(\lambda)$ random variable. Thus the $\{\tilde{S}_n\}_{n \geq 1}$ are i.i.d. $\text{Exp}(\lambda)$ random variables independent of S_1, \dots, S_n . Thus, conditional of $\{N_s = i\}$, $\{\tilde{N}_t\}_{t \geq 0}$ is a Poisson process independent of $\{X_r\}_{r \leq s}$. \square

Now, we can prove Theorem 4.2.2.

Proof. I give a number of proofs of equivalences here. The rest will be left for exercises or self-study.

Part 1. First, we will show that (i) implies (ii). Now, as the increments are stationary, we know $N_{t+h} - N_t$ has the same distribution as $N_h - N_0$. So,

$$\mathbb{P}(N_{t+h} - N_t = 0) = \mathbb{P}(N_h - N_0 = 0) = \mathbb{P}(N_h = 0) = e^{-\lambda h}.$$

Taking a Taylor expansion of the exponential function, we get $e^{-\lambda h} = 1 - \lambda h + o(h)$. Likewise,

$$\mathbb{P}(N_{t+h} - N_t = 1) = \mathbb{P}(N_h = 1) = \lambda h e^{-\lambda h} = \lambda h + o(h).$$

and $\mathbb{P}(N_{t+h} - N_t > 1) = o(h)$.

Part 2. We will show that (ii) implies (i). We do this by solving some differential equations. First, let us define $p_j(t) = \mathbb{P}(N_t = j)$. Then,

$$p_0(t+h) = \mathbb{P}(N_{t+h} = 0) = \mathbb{P}(N_{t+h} - N_t = 0) \mathbb{P}(N_t = 0) = (1 - \lambda h + o(h)) p_0(t).$$

Rearranging, we have

$$\frac{p_0(t+h) - p_0(t)}{h} = -\lambda p_0(t) + \frac{o(h)}{h}.$$

Because this holds for all t , we also get

$$\frac{p_0(t) - p_0(t-h)}{h} = -\lambda p_0(t-h) + \frac{o(h)}{h}.$$

Letting $h \rightarrow 0$ shows that $p_0(t)$ has a derivative (as the limit exists). This gives us

$$p_0(t)' = -\lambda p_0(t) \Rightarrow p_0(t) = Ae^{-\lambda t}.$$

Now, because $N_0 = 0$, we know that $p_0(0) = 1$ so $A = 1$ (i.e., $p_0(t) = e^{-\lambda t}$). Doing the same for $p_j(t)$ we have

$$\begin{aligned} p_j(t+h) &= \sum_{i=0}^j \mathbb{P}(N_{t+h} - N_t = i) \mathbb{P}(N_t = j-i) \\ &= \mathbb{P}(N_{t+h} - N_t = 0) \mathbb{P}(N_t = j) + \mathbb{P}(N_{t+h} - N_t = 1) \mathbb{P}(N_t = j-1) \\ &\quad + \sum_{i=2}^j \mathbb{P}(N_{t+h} - N_t = i) \mathbb{P}(N_t = j-i) \\ &= (1 - \lambda h + o(h)) p_j(t) + (\lambda h + o(h)) p_{j-1}(t) + o(h). \end{aligned}$$

Rearranging, we have

$$\frac{p_j(t+h) - p_j(t)}{h} = -\lambda p_j(t) + \lambda p_{j-1}(t) + \frac{o(h)}{h}.$$

By a similar argument to the one above, we get

$$p_j(t)' = -\lambda p_j(t) + \lambda p_{j-1}(t).$$

Now, remember that the product rule tells us that $(fg)' = f'g + g'f$. If we use the integrating factor $e^{\lambda t}$, we get

$$\begin{aligned} p_j(t)' &= -\lambda p_j(t) + \lambda p_{j-1}(t) \Rightarrow e^{\lambda t} p_j(t) = -\lambda e^{\lambda t} p_j(t) + \lambda e^{\lambda t} p_{j-1}(t) \\ &\Rightarrow e^{\lambda t} p_j(t)' + \lambda e^{\lambda t} p_j(t) = \lambda e^{\lambda t} p_{j-1}(t) \Rightarrow (e^{\lambda t} p_j(t))' = \lambda e^{\lambda t} p_{j-1}(t). \end{aligned}$$

We can solve this by induction. We start with $j = 1$. This gives

$$(e^{\lambda t} p_1(t))' = \lambda e^{\lambda t} p_0(t) = \lambda e^{\lambda t} e^{-\lambda t} = \lambda.$$

Integrating, we get

$$e^{\lambda t} p_1(t) = \lambda t + A \Rightarrow p_1(t) = t\lambda e^{\lambda t} + Ae^{\lambda t}.$$

Now, $p_j(0) = 0$ for $j > 0$, so $A = 0$ and $p_1(t) = t\lambda e^{\lambda t}$. Repeating in this way, we get

$$p_j(t) = \mathbb{P}(N_t = j) = e^{-\lambda t} \frac{(\lambda t)^j}{j!}.$$

Part 3. We show that (iii) implies (ii). This is just like the proof that (i) implies (ii). Using the Markov property (which was based on the interarrivals definition) we observe that $N_{t+h} - N_t$ has the same distribution as N_h , so

$$\mathbb{P}(N_{t+h} - N_t = 0) = \mathbb{P}(N_h = 0) = \mathbb{P}(S_1 > h) = e^{-\lambda h} = 1 - \lambda h + o(h),$$

and

$$\begin{aligned} \mathbb{P}(N_{t+h} - N_t = 1) &= \mathbb{P}(N_h = 1) = \mathbb{P}(S_1 < h, S_1 + S_2 > h) \\ &= \int_0^h e^{-\lambda(h-u)} \lambda e^{-\lambda u} du = \int_0^h \lambda e^{-\lambda h} du = \lambda h e^{-\lambda h} = \lambda h + o(h). \end{aligned}$$

It is also straightforward to see that

$$\mathbb{P}(N_{t+h} - N_t > 1) = \mathbb{P}(S_1 < h, S_1 + S_2 < h) \leq \mathbb{P}(S_1 < h) \mathbb{P}(S_2 < h) = o(h).$$

□

4.2.1 Order Statistics and the Distribution of Arrival Times

Order Statistics

Consider identically distributed random variables X_1, \dots, X_n with distribution $F(x)$ (that is $\mathbb{P}(X < x) = F(x)$). The *order statistics* of these random variables are simply the random variables ordered from smallest to largest. We write these as $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$. Two of the most important order statistics are the minimum, $X_{(1)}$, and the maximum, $X_{(n)}$.

Lemma 4.2.5. The distribution of the minimum is given by

$$\mathbb{P}(X_{(1)} \leq x) = 1 - (1 - F(x))^n.$$

Proof. We have

$$\begin{aligned} \mathbb{P}(X_{(1)} \leq x) &= 1 - \mathbb{P}(X_{(1)} > x) = 1 - \mathbb{P}(X_1 > x, \dots, X_n > x) \\ &= 1 - \mathbb{P}(X_1 > x) \cdots \mathbb{P}(X_n > x) = 1 - (\mathbb{P}(X > x))^n = 1 - (1 - F(x))^n. \end{aligned}$$

□

Lemma 4.2.6. The distribution of the maximum is given by

$$\mathbb{P}(X_{(n)} \leq x) = F(x)^n$$

Proof. We have

$$\mathbb{P}(X_{(n)} \leq x) = \mathbb{P}(X_1 \leq x, \dots, X_n \leq x) = \mathbb{P}(X \leq x)^n = F(x)^n.$$

□

Lemma 4.2.7. The density of the order statistics of n random variables U_1, \dots, U_n , with distribution $\mathcal{U}(a, b)$, is given by

$$f(u_1, \dots, u_n) = \frac{n!}{(b-a)^n} \mathbb{I}(a \leq u_1 \leq \dots \leq u_n \leq b)$$

Proof. It is easiest to begin with the distribution and then take derivatives to get the density. We wish to calculate $\mathbb{P}(U_{(1)} \leq u_1, \dots, U_{(n)} \leq u_n)$. Note that there are $n!$ orderings of the uniform random variables, each of which is equally likely. So,

$$\begin{aligned} \mathbb{P}(U_{(1)} \leq u_1, \dots, U_{(n)} \leq u_n) &= n! \mathbb{P}(U < u_1) \cdots \mathbb{P}(U < u_n) \\ &= n! \frac{u_1 - a}{b - a} \cdots \frac{u_n - a}{b - a}. \end{aligned}$$

Taking derivatives, we get

$$\begin{aligned} f(u_1, \dots, u_n) &= \frac{\partial}{\partial u_1} \cdots \frac{\partial}{\partial u_n} \mathbb{P}(U_{(1)} \leq u_1, \dots, U_{(n)} \leq u_n) \\ &= \frac{n!}{(b-a)^n} \mathbb{I}(a \leq u_1 \leq \dots \leq u_n \leq b). \end{aligned}$$

□

4.2.2 Distribution of Arrival Times

As it turns out, given that we know how many arrivals a Poisson process has had in an interval $[0, t]$ (that is, we know N_t), the arrival times will be uniformly distributed in the interval. This implies that the points of a Poisson process have very little structure to them (in some sense, it is a process that puts points as arbitrarily as possible on a line).

Theorem 4.2.8. Let $\{N_t\}_{t \geq 0}$ be a Poisson process. Then, conditional on $\{N_t = n\}$, T_1, \dots, T_n have the joint density function

$$f(t_1, \dots, t_n) = \frac{n!}{t^n} \mathbb{I}(0 \leq t_1 \leq \dots \leq t_n \leq t).$$

This is the density of the order statistics of i.i.d. uniform random variables on the interval $[0, t]$. This means the arrival times are distributed uniformly on this interval.

Proof. Consider the event $\{T_1 = t_1, \dots, T_n = t_n, N_t = n\}$. Because there is a bijection between arrival times and inter-arrival times, this should have the same probability density as the event

$$\{S_1 = t_1, S_2 = t_2 - t_1, \dots, S_n = t_n - t_{n-1}, S_{n+1} > t - t_n\}.$$

Because this is the joint density of i.i.d. exponential random variables, we can write this explicitly as

$$\lambda e^{-\lambda u_1} \lambda e^{-\lambda(u_2 - u_1)} \dots \lambda e^{-\lambda(u_n - u_{n-1})} e^{-\lambda(t - u_n)} = \lambda^n e^{-\lambda t}.$$

We then get the conditional density we wish by dividing by the probability of the event $\{N_t = n\}$,

$$\begin{aligned} f(t_1, \dots, t_n) &= \frac{\lambda^n e^{-\lambda t}}{(\lambda t)^n e^{-\lambda t} / n!} \mathbb{I}(0 \leq t_1 \leq \dots \leq t_n \leq t) \\ &= \frac{n!}{t^n} \mathbb{I}(0 \leq t_1 \leq \dots \leq t_n \leq t). \end{aligned}$$

□

4.3 Simulating Poisson Processes

As already mentioned, there are at least three ways of simulating a Poisson process. These follow directly from the different definitions we have used.

4.3.1 Using the Infinitesimal Definition to Simulate Approximately

The infinitesimal definition gives us a way to simulate the counting process $\{N_t\}_{t \geq 0}$ directly. This simulation is approximate but becomes increasingly good as $h \downarrow 0$. The idea is to slice the interval $[0, t]$ up into little pieces of length (sometimes called mesh size) h . In each one of these slices, we increase $\{N_t\}_{t \geq 0}$ with probability λh .

Listing 4.1: Matlab code

```

1 lambda = 4; t = 1; h = 0.0001;
2 mesh = 0:h:t; N = zeros(1, length(mesh));
3 S = []; jump_indices = [];
4
5 N(1) = 0;
6
```

```

7  for i = 2:length(mesh)
8      if rand < lambda * h
9          jump_indices = [jump_indices i];
10         N(i) = N(i-1) + 1;
11     else
12         N(i) = N(i-1);
13     end
14 end
15
16 if isempty(jump_indices)==0
17     Ts = (jump_indices - 1)*h;
18     S(1) = Ts(1);
19     if length(jump_indices) > 1
20         for i = 2:length(jump_indices)
21             S(i) = Ts(i) - Ts(i-1);
22         end
23     end
24 end

```

4.3.2 Simulating the Arrival Times

The idea of this approach is to simulate the arrival times directly. Given an interval $[0, t]$, we know that we have a $\text{Poi}(\lambda t)$ random number of arrivals. These are then distributed uniformly in $[0, t]$.

Listing 4.2: Matlab code

```

1  t = 5; lambda = 2;
2
3  T = [];
4  n = poissrnd(lambda * t);
5
6  if n~=0
7      T = sort(t * rand(n,1));
8      S = zeros(n,1);
9      S(1) = T(1);
10     if n > 1
11         for i = 2:n
12             S(i) = T(i) - T(i-1);
13         end
14     end
15 end

```

4.3.3 Simulating the Inter-Arrival Times

The idea here is to simulate the inter-arrival times, which are i.i.d. $\text{Exp}(\lambda)$ random variables. Note that, if $U \sim \mathcal{U}(0, 1)$, then $-\log(U)/\lambda$ is $\text{Exp}(\lambda)$.

Listing 4.3: Matlab code

```

1 lambda = 4; h = 0.0001; t = 1;
2
3 s = 0;
4
5 S = []; Ts = [];
6
7 while s <= t
8     inter_time = - log(rand) / lambda;;
9     s = s + inter_time;
10    if s > t
11        break;
12    else
13        Ts = [Ts s];
14        S = [S inter_time];
15    end
16 end

```

4.4 Inhomogenous Poisson Processes

For many of the processes that Poisson processes are used to model, such as queues and traffic, the assumption that events occur at a constant rate (i.e., λ is constant) is very unrealistic. If you think about traffic (either on the internet or on a road) it tends to be heavier in some time periods and lighter in others. Inhomogenous Poisson processes modify the definition of a Poisson process so that it can incorporate time-dependent arrivals. Inhomogenous Poisson processes can be defined in a number of ways. Note, however, that it is no longer straightforward to use a definition based on inter-arrival times.

Definition 4.4.1 (Inhomogenous Poisson Process). A point process $\{N_t\}_{t \geq 0}$ is said to be an inhomogenous Poisson process with intensity function $\lambda(t) \geq 0 \forall t \geq 0$.

- (i) $N_0 = 0$.
- (ii) For each $t \geq 0$, N_t has a Poisson distribution with parameter $\Lambda = \int_0^t \lambda(s) ds$.

- (iii) For each $0 \leq t_1 < t_2 < \dots < t_m$, the random variables $N_{t_1}, \dots, N_{t_m} - N_{t_{m-1}}$ are independent (that is, $\{N_t\}_{t \geq 0}$ has independent increments).

Definition 4.4.2 (Inhomogenous Poisson Process (Infinitesimal Definition)). A point process $\{N_t\}_{t \geq 0}$ is said to be an inhomogenous Poisson process with intensity function $\lambda(t) \geq 0 \forall t \geq 0$ if, as $h \downarrow 0$,

- (i) $\{N_t\}_{t \geq 0}$ has independent increments.
- (ii) $\mathbb{P}(N_{t+h} - N_t = 0) = 1 - \lambda(t)h + o(h)$.
- (iii) $\mathbb{P}(N_{t+h} - N_t = 1) = \lambda(t)h + o(h)$.
- (iv) $\mathbb{P}(N_{t+h} - N_t > 1) = o(h)$.

4.5 Simulating an Inhomogenous Poisson Process

There are at least two ways to simulate an inhomogenous Poisson process.

4.5.1 Acceptance-Rejection

One way to simulate an inhomogenous Poisson process on an interval $[0, t]$ is to simulate a homogenous Poisson process with parameter

$$\lambda^* = \max\{\lambda(s) : 0 \leq s \leq t\}$$

then ‘thin’ this process by only accepting arrivals with a certain probability. If an arrival / jump occurs at time T_i it should only be accepted with probability $\lambda(T_i)/\lambda^*$. It is not hard to check that this method works using the infinitesimal definition.

Listing 4.4: Matlab code

```

1 t = 10; lambda_star = 1;
2
3 T = [];
4 n = poissrnd(lambda_star * t);
5
6 if n~=0
7     point_count = 0;
8     for i = 1:n
9         T_temp = t * rand;
10        if rand < sin(T_temp)^2 / lambda_star

```

```

11         point_count = point_count + 1;
12         T(point_count) = T_temp;
13     end
14 end
15 if point_count ~= 0
16     T = sort(T);
17     S = zeros(point_count,1);
18     S(1) = T(1);
19     if point_count > 1
20         for i = 2:point_count
21             S(i) = T(i) - T(i-1);
22         end
23     end
24 end
25 end

```

4.5.2 Infinitesimal Approach (Approximate)

As in the homogenous Poisson process case, we can simulate an inhomogenous Poisson process approximately using its infinitesimal definition. This approximation becomes better as $h \downarrow 0$.

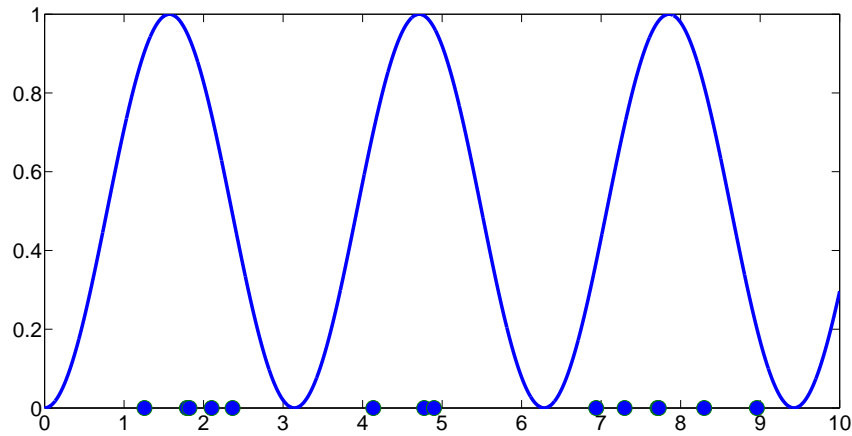


Figure 4.5.1: A realization of a inhomogenous Poisson process with the intensity function $\lambda(t) = 3 \sin^2(t)$ plotted.

Listing 4.5: Matlab code

```

1 lambda = 1; t = 10; h = 0.0001;

```

```

2 mesh = 0:h:t; N = zeros(1, length(mesh));
3 S = []; jump_indices = [];
4
5 N(1) = 0;
6
7 for i = 2:length(mesh)
8     if rand < 3*sin(h*(i-1))^2 * h
9         jump_indices = [jump_indices i];
10        N(i) = N(i-1) + 1;
11    else
12        N(i) = N(i-1);
13    end
14 end
15
16 if isempty(jump_indices)==0
17     Ts = (jump_indices - 1)*h;
18     S(1) = Ts(1);
19     if length(jump_indices) > 1
20         for i = 2:length(jump_indices)
21             S(i) = Ts(i) - Ts(i-1);
22         end
23     end
24 end

```

4.6 Compound Poisson Processes

A very useful extension of a Poisson process is what is called a *compound Poisson process*. A compound Poisson process replaces the unit jumps of a homogenous Poisson process with random jump sizes.

Definition 4.6.1 (Compound Poisson Process). Given a homogenous Poisson process, $\{N_t\}_{t \geq 0}$ and a jump distribution G , we say

$$X_t = \sum_{i=1}^{N_t} J_i,$$

is a compound Poisson process, where the jumps $\{J_n\}_{n \geq 0}$ are i.i.d. draws from the distribution G .

Example 4.6.2. A street musician plays the accordion in the main street of Ulm for three hours. He hopes to earn enough for a beer, which costs €3.50. Throughout the three hours, people give him coins at random. There does

not seem to be any pattern to when people give him money, so a Poisson process is a good model. The amount of money each person gives is random, with distribution

$$\mathbb{P}(\text{€}0.05) = 2/5$$

$$\mathbb{P}(\text{€}0.10) = 2/5$$

$$\mathbb{P}(\text{€}0.20) = 1/5.$$

On average, 5 people per hour give the street musician money. This implies that the Poisson process has intensity $\lambda = 5$. What is the probability the musician gets his beer? That is, what is $\ell = \mathbb{P}(X_3 \geq 3.50)$. We can estimate this easily using Monte Carlo.

Listing 4.6: Matlab code

```

1 t = 3; lambda = 5; N = 10^6;
2 beer = zeros(N,1); beer_price = 350;
3
4 for i = 1:N
5
6     n = poissrnd(lambda * t);
7
8     if n~=0
9         coins = zeros(n,1);
10        for j = 1:n
11            U = rand;
12            coins(j) = (U <= 2/5)*5 + ...
13                    (U > 2/5 && U <= 4/5)*10 + (U > 4/5)*20;
14        end
15    end
16
17    beer(i) = (sum(coins) >= beer_price);
18 end
19
20 ell_hat = mean(beer)
21 re_hat = std(beer) / (ell_hat * sqrt(N))

```

Chapter 5

Continuous Time Markov Chains

Continuous time Markov chains (CTMCS) — or, at least, the ones we consider — behave much like discrete time Markov chains, with the key difference that the jumps between states take place at random times rather than at fixed steps. Continuous time Markov chains are trickier to work with than discrete time Markov chains, because there are a number of technical issues and strange / pathological things that can happen. This behavior is mainly related to situations where an infinite number of jumps can happen in a finite time. In this course, we will not consider chains where such things happen. This is because these chains can not really be simulated and considering a smaller subset of continuous time Markov chains will be sufficient for modeling most real world phenomena we might be interested in.

The chains we consider can be described in two different but equivalent ways: *transition functions* and *infinitesimal generators*. Warning: these descriptions are not always equivalent or even valid when working with more general classes of CTMCs.

5.1 Transition Function

Transition functions are the continuous time equivalent of the transition matrix, P , that we have already encountered when discussing discrete time Markov chains. That is, they are functions that allow us to calculate probabilities of the form $\mathbb{P}_i(X_t = j)$. More formally, a transition function is defined as follows.

Definition 5.1.1 (Transition function). A transition function $p_t(i, j)$ with

$i, j \in \mathcal{X}$ and $t \geq 0$ is a real-valued function satisfying the following properties

- (i) $p_t(i, j) \geq 0$ for all $i, j \in \mathcal{X}$ and $t \geq 0$.
- (ii) $\sum_{j \in \mathcal{X}} p_t(i, j) = 1$ for all $i \in \mathcal{X}$ and $t \geq 0$.
- (iii) $\lim_{t \downarrow 0} p_t(i, i) = p_0(i, i) = 1$ for all $i \in \mathcal{X}$.
- (iv) The Chapman-Kolmogorov equations:

$$p_{s+t}(i, j) = \sum_{k \in \mathcal{X}} p_s(i, k) p_t(k, j)$$

for all $i, j \in \mathcal{X}$ and $t \geq 0$.

We can think of transitions functions as a family of matrices, $\{P(t)\}_{t \geq 0}$, indexed by t . For a fixed value of t , $(P(t))_{i,j} = p_t(i, j)$. We can write the properties out again in terms of matrices:

Definition 5.1.2 (Transition function: Matrix Version). A transition function is a family of matrices, $\{P(t)\}_{t \geq 0}$ with the following properties.

- (i) They are non-negative, real-valued and $\sum_{j \in \mathcal{X}} P_{i,j} = 1$ for all $i \in \mathcal{X}$ and $t \geq 0$. Matrices that are non-negative with unit row sums are called *stochastic matrices*.
- (ii) They satisfy the Chapman-Kolmogorov equations: $P(s+t) = P(s)P(t)$.

5.2 Infinitesimal Generator

Definition 5.2.1. The infinitesimal generator (or Q -matrix) is a real valued matrix satisfying the following properties

- (i) $q_{i,j} \geq 0$ for all $i \neq j$.
- (ii) $\sum_{j \in \mathcal{X}} q_{i,j} = 0$.

Thus, we require that $q_{i,i} = -\sum_{j \neq i} q_{i,j}$. Because it plays an important role in our treatment of CTMCs, we denote $q_i = -q_{i,i}$.

INF DEF. COMPETING POISSON PROCESS VERSION.
EXAMPLE.

5.3 Continuous Time Markov Chains

FORMALLY. OUR VERSION.

5.4 The Jump Chain and Holding Times

If we only consider the sequence of states that a CTMC visits (by looking at the value the CTMC takes immediately after each jump), we get a discrete time Markov chain $\{Y_n\}_{n \geq 0}$ called the *jump chain* of the CTMC. We call the time spent in each state a *holding time*. Associate with a path of the jump chain Y_0, \dots, Y_n is a sequence of holding times S_1, \dots, S_{n+1} , where S_1 is the time spent in Y_0 before jumping to Y_1 , S_2 is the time spent in Y_1 before jumping to Y_2 , and so on. The definitions of the jump chain and holding times are made clear in figure (REF).

PICTURE HERE

The jump chain $\{Y_n\}_{n \geq 0}$ is a discrete time Markov chain, so it can be described by a transition matrix J (we do not use P in order to avoid confusion with the transition function P). As it turns out, J can be written in terms of the Q -matrix.

Lemma 5.4.1. Given a CTMC with infinitesimal generator Q , the jump matrix J is defined by

$$J_{i,i} = 0$$

for all $i \in \mathcal{X}$ and

$$J_{i,j} = \frac{q_{i,j}}{q_i}$$

for all $i, j \in \mathcal{X}$ such that $j \neq i$.

The holding times are also defined by the Q -matrix.

Lemma 5.4.2. Given Y_0, Y_1, \dots, Y_n , the holding times S_1, \dots, S_{n+1} are exponential random variables with parameters $q_{Y_0}, q_{Y_1}, \dots, q_{Y_n}$.

Example 5.4.3. Consider the following CTMC.

PICTURE HERE

The Q matrix of this chain is given by

$$Q = \begin{bmatrix} -\lambda & \lambda & 0 \\ \mu & -(\lambda + \mu) & \mu \\ 0 & \mu & -\mu \end{bmatrix}.$$

Thus, the transition matrix of the jump chain is

$$J = \begin{bmatrix} 0 & 1 & 0 \\ \frac{\mu}{\lambda + \mu} & 0 & \frac{\lambda}{\lambda + \mu} \\ 0 & 1 & 0 \end{bmatrix},$$

and the amount of time spent in state 1 is $\text{Exp}(\lambda)$, the amount of time spent in state 2 is $\text{Exp}(\lambda + \mu)$ and the amount of time spent in state 3 is $\text{Exp}(\mu)$.

5.5 Examples of Continuous Time Markov Chains

5.5.1 Poisson Process

We are already quite familiar with one CTMC: the Poisson process. This process has exponential holding times (always with parameter λ) and a jump chain that always increases by one. That is, the jump chain is of the form

$$J = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 1 & 0 & 0 & \cdots \\ 0 & 0 & 0 & 1 & 0 & \cdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \end{bmatrix}.$$

This gives a Q -matrix of the form

$$Q = \begin{bmatrix} -\lambda & \lambda & 0 & 0 & 0 & \cdots \\ 0 & -\lambda & \lambda & 0 & 0 & \cdots \\ 0 & 0 & -\lambda & \lambda & 0 & \cdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \end{bmatrix}.$$

5.5.2 Birth-Death Process

A Poisson process is sometimes called a pure birth process, as it models a population that is constantly growing. More generally, birth-death processes are simple models of populations where births and deaths happen at random. A simple example is the following

PICTURE HERE

The Q -matrix for this chain is of the form

$$Q = \begin{bmatrix} -\lambda & \lambda & 0 & 0 & 0 & \cdots \\ \mu & -(\mu + \lambda) & \lambda & 0 & 0 & \cdots \\ 0 & \mu & -(\mu + \lambda) & \lambda & 0 & \cdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots \end{bmatrix}.$$

5.6 Simulating Continuous Time Markov Chains

If we think about CTMCs in terms of jump chains and holding times, then the generic method for simulating them becomes obvious. We simply wait an exponential rate of time in each state, determined by the Q -matrix,

the jump to the next state according to the J matrix. The jump chain itself is simulated just as in the section on discrete time Markov chains. As with discrete time Markov chains, it is important to consider whether your chain is finite or infinite, as this usually changes the choice of how to encode J . If the chain is infinite, then it obviously is not possible to write out the J matrix.

Example 5.6.1 (A finite state space CTMC). Consider the following CTMC

PICTURE HERE,

with initial distribution $\lambda = (1/2, 1/2, 0, 0)$. This chain has Q -matrix

$$Q = \begin{bmatrix} -9 & 2 & 3 & 4 \\ 1 & -4 & 3 & 0 \\ 1 & 2 & -7 & 4 \\ 1 & 0 & 3 & -4 \end{bmatrix}.$$

and jump matrix

$$J = \begin{bmatrix} 0 & 2/9 & 3/9 & 4/9 \\ 1/4 & 0 & 3/4 & 0 \\ 1/7 & 2/7 & 0 & 4/7 \\ 1/4 & 0 & 3/4 & 0 \end{bmatrix}.$$

The Matlab code for simulating this is as follows.

Listing 5.1: Matlab code

```

1 T = 10; t = 0;
2
3 Q = [-9 2 3 4; 1 -4 3 0; 1 2 -7 4; 1 0 3 -4];
4 J = [0 2/9 3/9 4/9; 1/4 0 3/4 0; 1/7 2/7 0 4/7; 1/4 0 3/4 0];
5
6 U = rand; X = (U <= 1/2) + 2*(U > 1/2);
7 jump_times = [0]; jump_chain = [X];
8
9 while t <= T
10     t = t + log(rand) / Q(X,X);
11     if t > T
12         break;
13     else
14         jump_times = [jump_times t];
15         X = min(find(rand < cumsum(J(X,:)))));
16         jump_chain = [jump_chain X];
17     end
18 end

```

Example 5.6.2 (A birth-death process). Consider the birth-death process pictured in figure (REF) such that $\mathbb{P}(X_0 = 0) = 1$.

PICTURE HERE

We can implement this in Matlab by encoding the transition rule itself.

Listing 5.2: Matlab code

```

1 T = 100000; t = 0;
2 lambda = 2; mu = 0.25;
3
4 X = 0; jump_times = [0]; jump_chain = [X];
5
6 while t <= T
7     t = t - log(rand) / (lambda + X*mu);
8     if t > T
9         break;
10    else
11        jump_times = [jump_times t];
12
13        if rand < lambda / (lambda + X*mu)
14            X = X + 1;
15        else
16            X = X - 1;
17        end
18        jump_chain = [jump_chain X];
19    end
20 end

```

5.7 The Relationship Between P and Q in the Finite Case

In the finite case, we can write P directly in terms of Q . Remember that P must satisfy $P(s + t) = P(t)P(s)$. This implies that P should be some kind of exponential like function. In fact, we can write

$$P(t) = e^{tQ}$$

where e^{tQ} is the *matrix exponential*. It might at first be unclear how to define a matrix exponential. However, if we remember the Taylor series definition

of the exponential, then we realise we can write something equivalent using matrices. That is

$$P(t) = \sum_{k=0}^{\infty} \frac{(tQ)^k}{k!}.$$

Warning! It is often very computationally challenging to compute a matrix exponential. The answer the computer gives you might not always be correct. In Matlab you need to use ‘expm’ rather than ‘exp’ if you wish to use the matrix exponential.

Example 5.7.1. Consider the CTMC given earlier with Q -matrix

$$Q = \begin{bmatrix} -9 & 2 & 3 & 4 \\ 1 & -4 & 3 & 0 \\ 1 & 2 & -7 & 4 \\ 1 & 0 & 3 & -4 \end{bmatrix}$$

and initial distribution $\lambda = (1/2, 1/2, 0, 0)$. We can calculate

We can calculate e^{tQ} easily using Matlab. For $t = 0.01, 0.1, 1, 10, 100$ we have

Listing 5.3: Matlab code

```

1 [1/2 1/2 0 0] * expm(.01 * Q) = 0.4619 0.4901 0.0285 0.0194
2 [1/2 1/2 0 0] * expm(.1 * Q) = 0.2472 0.4112 0.1896 0.1520
3 [1/2 1/2 0 0] * expm(1 * Q) = 0.1000 0.2061 0.3000 0.3939
4 [1/2 1/2 0 0] * expm(10 * Q) = 0.1000 0.2000 0.3000 0.4000
5 [1/2 1/2 0 0] * expm(100 * Q) = 0.1000 0.2000 0.3000 0.4000.
```

Looking at the example, we see that the distribution of the chain seems to converge to a stationary distribution.

5.8 Irreducibility, Recurrence and Positive Recurrence

As in the case of discrete time Markov chains, we want to establish conditions under which a CTMC has a unique stationary distribution. As it transpires, these conditions are basically the same as those for discrete time Markov chains. This is because it is the structure of the jump chain that determines whether a stationary distribution exists or not.

We say a state i leads to a state j (written $i \rightarrow j$) if

$$\mathbb{P}_i(X_t = j \text{ for some } t \geq 0) > 0.$$

We say i and j communicate (or $i \leftrightarrow j$) if $i \rightarrow j$ and $j \rightarrow i$. The following theorem shows that it is sufficient to look at the jump chain in order to determine communicating classes.

Theorem 5.8.1. Consider a CTMC $\{X_t\}_{t \geq 0}$ with Q -matrix Q and jump chain $\{Y_n\}_{n \geq 0}$. For distinct states i and j the following are equivalent

- (i) $i \rightarrow j$.
- (ii) $i \rightarrow j$ for the jump chain.
- (iii) $q_{i,k_1} q_{k_1,k_2} \cdots q_{k_n,j} > 0$ for some $n > 0$ and states k_1, \dots, k_n .
- (iv) $p_t(i, j) > 0$ for all $t > 0$.
- (v) $p_t(i, j) > 0$ for some $t > 0$.

Proof. See [1]. □

As in the discrete time case, a CTMC is irreducible if $i \leftrightarrow j$ for all $i, j \in \mathcal{X}$.

We say a state i is *recurrent* if

$$\mathbb{P}_i(\{t \geq 0 : X_t = i\} \text{ is unbounded}) = 1.$$

We say a state j is *transient* if

$$\mathbb{P}_i(\{t \geq 0 : X_t = i\} \text{ is unbounded}) = 0.$$

Again, it is sufficient to look at the jump chain to establish that these properties hold.

Theorem 5.8.2. Consider a CTMC $\{X_t\}_{t \geq 0}$ with jump chain $\{Y_n\}_{n \geq 0}$.

- (i) If i is recurrent for $\{Y_n\}_{n \geq 0}$ it is recurrent for $\{X_t\}_{t \geq 0}$.
- (ii) If i is transient for $\{Y_n\}_{n \geq 0}$ it is recurrent for $\{X_t\}_{t \geq 0}$.

Proof. See [1]. □

A state i is *positive recurrent* if it is positive recurrent for the jump chain.

5.9 Invariant Measures and Stationary Distribution

Definition 5.9.1 (Invariant Measure). We say a measure (remember this is a vector with non-negative elements) μ is invariant for a Q -matrix Q if

$$\mu Q = \mathbf{0}.$$

Invariant measures of CTMCs are closely related to invariant measures for the associated jump chain.

Theorem 5.9.2. Given a CTMC $\{X_t\}_{t \geq 0}$ with Q -matrix Q and jump matrix J , the following are equivalent

- (i) μ is invariant for Q .
- (ii) $\nu J = \nu$, where $\nu_i = \mu_i q_i$ for all $i \in \mathcal{X}$.

Note that $1/q_i$ is the expected holding time in state i , so μ has elements $\mu_i = \nu_i/q_i$ and is thus a invariant measure for the jump chain reweighted by the expected time spent in each state. Note that neither μ or ν is necessarily a probability distribution. Furthermore, even if μ is a probability distribution, the ν defined in respect to it will probably not be a probability distribution (until it is normalized). This is also the case if ν is a probability distribution.

We are now able to give conditions for the existence of a unique stationary measure for a CTMC $\{X_t\}_{t \geq 0}$. Unsurprisingly, these are the same conditions as for discrete time Markov chains.

Theorem 5.9.3. Let Q be an irreducible Q matrix. Then, the following are equivalent.

- (i) Every state is positive recurrent.
- (ii) Some state, i , is positive recurrent.
- (iii) Q has an invariant distribution ρ .

Example 5.9.4. Given that a stationary distribution exists, we can find it by solving the linear system $\rho Q = \mathbf{0}$. Consider the example from earlier with Q matrix

$$Q = \begin{bmatrix} -9 & 2 & 3 & 4 \\ 1 & -4 & 3 & 0 \\ 1 & 2 & -7 & 4 \\ 1 & 0 & 3 & -4 \end{bmatrix}.$$

Solving this system, we find $\rho = (1/10, 2/10, 3/10, 4/10)$.

MENTION LIMITING DISTRIBUTION

5.9.1 Reversibility and Detailed Balance

As is the case with discrete time Markov chains, reversible chains are particularly nice to work with. This largely due to detailed balance equations holding.

Theorem 5.9.5. Given a CTMC $\{X_t\}_{t \geq 0}$ with irreducible Q -matrix Q and invariant distribution $\boldsymbol{\rho}$ (which also serves as the initial distribution of $\{X_t\}_{t \geq 0}$). The process $\{\hat{X}_t\}_{t \geq 0}$ defined by $\hat{X}_t = X_{T-t}$ is Markov $(\boldsymbol{\rho}, \hat{Q})$. Additionally, \hat{Q} is irreducible and has invariant distribution $\boldsymbol{\rho}$.

Proof. See [1]. □

We say a CTMC is *reversible* if $Q = \hat{Q}$.

Definition 5.9.6 (Detailed Balance Equations). A Q -matrix, Q , and a measure $\boldsymbol{\mu}$ are in *detailed balance* if

$$\mu_i q_{i,j} = \mu_j q_{j,i} \text{ for all } i, j \in \mathcal{X}.$$

The following theorem is more or less identical to the discrete version.

Theorem 5.9.7. If a Q -matrix, Q , and a measure $\boldsymbol{\mu}$ are in detailed balance then $\boldsymbol{\mu}$ is invariant for Q .

Proof. For a given $i \in \mathcal{X}$, we have

$$(\boldsymbol{\mu}Q)_i = \sum_{j \in \mathcal{X}} \mu_j q_{j,i} = \sum_{j \in \mathcal{X}} \mu_i q_{i,j} = 0.$$

□

The only CTMCs that satisfy detailed balance are reversible ones.

Theorem 5.9.8. Let Q be an irreducible Q -matrix and $\boldsymbol{\rho}$ a distribution. Suppose $\{X_t\}_{t \geq 0}$ is Markov $(\boldsymbol{\rho}, Q)$. Then, the following are equivalent

- (i) $\{X_t\}_{t \geq 0}$ is reversible.
- (ii) Q and $\boldsymbol{\rho}$ are in detailed balance.

Proof. See [1]. □

Chapter 6

Gaussian Processes

Gaussian processes are a reasonably large class of processes that have many applications, including in finance, time-series analysis and machine learning. Certain classes of Gaussian processes can also be thought of as spatial processes. We will use these as a vehicle to start considering more general spatial objects.

Definition 6.0.9 (Gaussian Process). A stochastic process $\{X_t\}_{t \geq 0}$ is Gaussian if, for any choice of times t_1, \dots, t_n , the random vector $(X_{t_1}, \dots, X_{t_n})$ has a *multivariate normal distribution*.

6.1 The Multivariate Normal Distribution

Because Gaussian processes are defined in terms of the multivariate normal distribution, we will need to have a pretty good understanding of this distribution and its properties.

Definition 6.1.1 (Multivariate Normal Distribution). A vector $\mathbf{X} = (X_1, \dots, X_n)$ is said to be multivariate normal (multivariate Gaussian) if all linear combinations of \mathbf{X} , i.e. all random variables of the form

$$\sum_{k=1}^n \alpha_k X_k$$

have univariate normal distributions.

This is quite a strong definition. Importantly, it implies that, even if all of its components are normally distributed, a random vector is not necessarily multivariate normal.

Example 6.1.2 (A random vector with normal marginals that is not multivariate normal). Let $X_1 \sim \mathbf{N}(0, 1)$ and

$$X_2 = \begin{cases} X_1 & \text{if } |X_1| \leq 1 \\ -X_1 & \text{if } |X_1| > 1 \end{cases}.$$

Note that $X_2 \sim \mathbf{N}(0, 1)$. However, $X_1 + X_2$ is not normally distributed, because $|X_1 + X_2| \leq 2$, which implies $X_1 + X_2$ is bounded and, hence, cannot be normally distributed.

Linear transformations of multivariate normal random vectors are, again, multivariate normal.

Lemma 6.1.3. Suppose $\mathbf{X} = (X_1, \dots, X_n)$ is multivariate normal and A is an $m \times n$ real-valued matrix. Then, $\mathbf{Y} = A\mathbf{X}$ is also multivariate normal.

Proof. Any linear combination of Y_1, \dots, Y_m is a linear combination of linear combinations of X_1, \dots, X_n and, thus, univariate normal. \square

Theorem 6.1.4. A multivariate normal random vector $\mathbf{X} = (X_1, \dots, X_n)$ is completely described by a mean vector $\boldsymbol{\mu} = \mathbb{E}\mathbf{X}$ and a covariance matrix $\Sigma = \text{Var}(\mathbf{X})$.

Proof. The distribution of \mathbf{X} is described by its characteristic function which is

$$\mathbb{E} \exp \{i\boldsymbol{\theta}^\top \mathbf{X}\} = \mathbb{E} \exp \left\{ i \sum_{i=1}^n \theta_i X_i \right\}.$$

Now, we know $\sum_{i=1}^n \theta_i X_i$ is a univariate normal random variable (because \mathbf{X} is multivariate normal). Let $m = \mathbb{E} \sum_{i=1}^n \theta_i X_i$ and $\sigma^2 = \text{Var}(\sum_{i=1}^n \theta_i X_i)$. Then,

$$\mathbb{E} \left\{ i \sum_{i=1}^n \theta_i X_i \right\} = \exp \left\{ im - \frac{1}{2} \sigma^2 \right\}.$$

Now

$$m = \mathbb{E} \sum_{i=1}^n \theta_i X_i = \sum_{i=1}^n \theta_i \mu_i$$

and

$$\sigma^2 = \text{Var} \left(\sum_{i=1}^n \theta_i X_i \right) = \sum_{i=1}^n \sum_{j=1}^n \theta_i \theta_j \text{Cov}(X_i, X_j) = \sum_{j=1}^n \sum_{i=1}^n \theta_i \theta_j \Sigma_{i,j}.$$

So, everything is specified by $\boldsymbol{\mu}$ and Σ . \square

INDEPENDENCE

There is nothing too difficult about dealing with the mean vector $\boldsymbol{\mu}$. However, the covariance matrix makes things pretty difficult, especially when we wish to simulate a high dimensional random vector. In order to simulate Gaussian processes effectively, we need to exploit as many properties of covariance matrices as possible.

6.1.1 Symmetric Positive Definite and Semi-Positive Definite Matrices

Covariance matrices are members of a family of matrices called *symmetric positive definite matrices*.

Definition 6.1.5 (Positive Definite Matrices (Real-Valued)). An $n \times n$ real-valued matrix, A , is positive definite if and only if

$$\mathbf{x}^\top A \mathbf{x} > 0 \text{ for all } \mathbf{x} \neq \mathbf{0}.$$

If A is also symmetric, then A is called symmetric positive definite (SPD). Some important properties of SPD matrices are

- (i) $\text{rank}(A) = n$.
- (ii) $|A| > 0$.
- (iii) $A_{i,i} > 0$.
- (iv) A^{-1} is SPD.

Lemma 6.1.6 (Necessary and sufficient conditions for an SPD). The following are necessary and sufficient conditions for an $n \times n$ matrix A to be SPD

- (i) All the eigenvalues $\lambda_1, \dots, \lambda_n$ of A are strictly positive.
- (ii) There exists a unique matrix C such that $A = CC^\top$, where C is a real-valued lower-triangular matrix with positive diagonal entries. This is called the *Cholesky decomposition* of A .

Definition 6.1.7 (Positive Semi-definite Matrices (Real-Valued)). An $n \times n$ real-valued matrix, A , is positive semi-definite if and only if

$$\mathbf{x}^\top A \mathbf{x} \geq 0 \text{ for all } \mathbf{x} \neq \mathbf{0}.$$

If A is also symmetric, then A is called symmetric positive semi-definite (SPSD). Note that if A is an SPD then there is a real-valued decomposition

$$A = LL^\top,$$

though it is not necessarily unique and L may have zeroes on the diagonals.

Lemma 6.1.8. Covariance matrices are SPD.

Proof. Given an $n \times 1$ real-valued vector \mathbf{x} and a random vector \mathbf{Y} with covariance matrix Σ , we have

$$\text{Var}(\mathbf{x}^\top \mathbf{Y}) = \mathbf{x}^\top \text{Var}(\mathbf{Y}) \mathbf{x}.$$

Now this must be non-negative (as it is a variance). That is, it must be the case that

$$\mathbf{x}^\top \text{Var}(\mathbf{Y}) \mathbf{x} \geq 0,$$

so Σ is positive semi-definite. Symmetry comes from the fact that $\text{Cov}(X, Y) = \text{Cov}(Y, X)$. \square

Lemma 6.1.9. SPD matrices are covariance matrices.

Proof. Let A be an SPD matrix and \mathbf{Z} be a vector of random variables with $\text{Var}(\mathbf{Z}) = I$. Now, as A is SPD, $A = LL^\top$. So,

$$\text{Var}(L\mathbf{Z}) = L\text{Var}(\mathbf{Z})L^\top = LIL^\top = A,$$

so A is the covariance matrix of $L\mathbf{Z}$. \square

COVARIANCE POS DEF ...

6.1.2 Densities of Multivariate Normals

If $\mathbf{X} = (X_1, \dots, X_n)$ is $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ and Σ is positive definite, then \mathbf{X} has the density

$$f(\mathbf{x}) = f(x_1, \dots, x_n) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

6.1.3 Simulating Multivariate Normals

One possible way to simulate a multivariate normal random vector is using the Cholesky decomposition.

Lemma 6.1.10. If $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, I)$, $\Sigma = AA^\top$ is a covariance matrix, and $\mathbf{X} = \boldsymbol{\mu} + A\mathbf{Z}$, then $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$.

Proof. We know from lemma 6.1.3 that $A\mathbf{Z}$ is multivariate normal and so is $\boldsymbol{\mu} + A\mathbf{Z}$. Because a multivariate normal random vector is completely described by its mean vector and covariance matrix, we simply need to find $\mathbb{E}\mathbf{X}$ and $\text{Var}(\mathbf{X})$. Now,

$$\mathbb{E}\mathbf{X} = \mathbb{E}\boldsymbol{\mu} + \mathbb{E}A\mathbf{Z} = \boldsymbol{\mu} + A\mathbf{0} = \boldsymbol{\mu}$$

and

$$\text{Var}(\mathbf{X}) = \text{Var}(\boldsymbol{\mu} + A\mathbf{Z}) = \text{Var}(A\mathbf{Z}) = A\text{Var}(\mathbf{Z})A^\top = AIA^\top = AA^\top = \Sigma.$$

□

Thus, we can simulate a random vector $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ by

- (i) Finding A such that $\Sigma = AA^\top$.
- (ii) Generating $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, I)$.
- (iii) Returning $\mathbf{X} = \boldsymbol{\mu} + A\mathbf{Z}$.

Matlab makes things a bit confusing because its function ‘chol’ produces a decomposition of the form $B^\top B$. That is, $A = B^\top$. It is important to be careful and think about whether you want to generate a row or column vector when generating multivariate normals.

The following code produces column vectors.

Listing 6.1: Matlab code

```

1 mu = [1 2 3]';
2 Sigma = [3 1 2; 1 2 1; 2 1 5];
3 A = chol(Sigma);
4 X = mu + A' * randn(3,1);

```

The following code produces row vectors.

Listing 6.2: Matlab code

```

1 mu = [1 2 3];
2 Sigma = [3 1 2; 1 2 1; 2 1 5];
3 A = chol(Sigma);
4 X = mu + randn(1,3) * A;

```

The complexity of the Cholesky decomposition of an arbitrary real-valued matrix is $O(n^3)$ floating point operations.

6.2 Simulating a Gaussian Processes Version 1

Because Gaussian processes are defined to be processes where, for any choice of times t_1, \dots, t_n , the random vector $(X_{t_1}, \dots, X_{t_n})$ has a multivariate normal density and a multivariate normal density is completely describe by a mean vector $\boldsymbol{\mu}$ and a covariance matrix Σ , the probability distribution of Gaussian process is completely described if we have a way to construct the mean vector and covariance matrix for an arbitrary choice of t_1, \dots, t_n .

We can do this using an *expectation function*

$$\mu(t) = \mathbb{E}X_t$$

and a covariance function

$$r(s, t) = \text{Cov}(X_s, X_t).$$

Using these, we can simulate the values of a Gaussian process at fixed times t_1, \dots, t_n by calculating $\boldsymbol{\mu}$ where $\mu_i = \mu(t_i)$ and Σ , where $\Sigma_{i,j} = r(t_i, t_j)$ and then simulating a multivariate normal vector.

In the following examples, we simulate the Gaussian processes at evenly spaced times $t_1 = 1/h, t_2 = 2/h, \dots, t_n = 1$.

Example 6.2.1 (Brownian Motion). *Brownian motion* is a very important stochastic process which is, in some sense, the continuous time continuous state space analogue of a simple random walk. It has an expectation function $\mu(t) = 0$ and a covariance function $r(s, t) = \min(s, t)$.

Listing 6.3: Matlab code

```

1 %use mesh size h
2 h = 1000; t = 1/h : 1/h : 1;
3 n = length(t);
4
5 %Make the mean vector
6 mu = zeros(1,n);
7 for i = 1:n
8     mu(i) = 0;
9 end
10
```

```

11 %Make the covariance matrix
12 Sigma = zeros(n,n);
13 for i = 1:n
14     for j = 1:n
15         Sigma(i,j) = min(t(i),t(j));
16     end
17 end
18
19 %Generate the multivariate normal vector
20 A = chol(Sigma);
21 X = mu + randn(1,n) * A;
22
23 %Plot
24 plot(t,X);

```

Example 6.2.2 (Ornstein-Uhlenbeck Process). Another very important Gaussian process is the Ornstein-Uhlenbeck process. This has expectation function $\mu(t) = 0$ and covariance function $r(s, t) = e^{-\alpha|s-t|/2}$.

Listing 6.4: Matlab code

```

1 %use mesh size h
2 h = 1000; t = 1/h : 1/h : 1;
3 n = length(t);
4
5 %paramter of OU process
6 alpha = 10;
7
8 %Make the mean vector
9 mu = zeros(1,n);
10 for i = 1:n
11     mu(i) = 0;
12 end
13
14 %Make the covariance matrix
15 Sigma = zeros(n,n);
16 for i = 1:n
17     for j = 1:n
18         Sigma(i,j) = exp(-alpha * abs(t(i) - t(j)) / 2);
19     end
20 end
21
22 %Generate the multivariate normal vector
23 A = chol(Sigma);

```

```

24 X = mu + randn(1,n) * A;
25
26 %Plot
27 plot(t,X);

```

Example 6.2.3 (Fractional Brownian Motion). Fractional Brownian motion (fBm) is a generalisation of Brownian Motion. It has expectation function $\mu(t) = 0$ and covariance function $\text{Cov}(s, t) = 1/2(t^{2H} + s^{2H} - |t - s|^{2H})$, where $H \in (0, 1)$ is called the Hurst parameter. When $H = 1/2$, fBm reduces to standard Brownian motion. Brownian motion has independent increments. In contrast, for $H > 1/2$ fBm has positively correlated increments and for $H < 1/2$ fBm has negatively correlated increments.

Listing 6.5: Matlab code

```

1 %Hurst parameter
2 H = .9;
3
4 %use mesh size h
5 h = 1000; t = 1/h : 1/h : 1;
6 n = length(t);
7
8 %Make the mean vector
9 mu = zeros(1,n);
10 for i = 1:n
11     mu(i) = 0;
12 end
13
14 %Make the covariance matrix
15 Sigma = zeros(n,n);
16 for i = 1:n
17     for j = 1:n
18         Sigma(i,j) = 1/2 * (t(i)^(2*H) + t(j)^(2*H)...
19             - (abs(t(i) - t(j)))^(2 * H));
20     end
21 end
22
23 %Generate the multivariate normal vector
24 A = chol(Sigma);
25 X = mu + randn(1,n) * A;
26
27 %Plot
28 plot(t,X);

```

6.3 Stationary and Weak Stationary Gaussian Processes

Gaussian processes (and stochastic processes in general) are easier to work with if they are *stationary stochastic processes*.

Definition 6.3.1 (Stationary Stochastic Process). A stochastic process $\{X_t\}_{t \geq 0}$ is said to be stationary if the random vectors $(X_{t_1}, X_{t_2}, \dots, X_{t_n})$ and $(X_{t_1+s}, X_{t_2+s}, \dots, X_{t_n+s})$ have the same distribution for all choices of s, n and t_1, t_2, \dots, t_n .

An example of such a process would be an irreducible positive recurrent continuous time Markov chain started from its stationary distribution.

The conditions for a stochastic process to be stationary are quite strict. Often, a slightly weaker version of stationarity, called *weak stationarity* is required instead.

Definition 6.3.2 (Weak Stationary Stochastic Process). A stochastic process $\{X_t\}_{t \geq 0}$ is said to be weak stationary (sometimes *wide sense stationary* or *second order stationary*) if $\mathbb{E}X_t = c$ for all $t \geq 0$ and $\text{Cov}(X_t, X_{t+s})$ does not depend on t .

An example of a weak stationary process is the Ornstein-Uhlenbeck process, as $\mathbb{E}X_t = \mu(t) = 0$ and $\text{Cov}(X_t, X_{t+s}) = r(t, t+s) = e^{-\alpha|t-(t+s)|/2} = e^{-\alpha s/2}$.

Lemma 6.3.3. Gaussian processes that are weak stationary are stationary.

Proof. We know from theorem 6.1.4 that Gaussian distributions are entirely determined by their mean vector and covariance matrix. Since the mean and covariance of a weakly stationary process do not change when the times are all shifted by s , a weakly stationary Gaussian process is stationary. \square

An Ornstein-Uhlenbeck process is a weak stationary Gaussian process, so it is a stationary stochastic process.

6.4 Finite Dimensional Distributions

A very nice property of Gaussian processes is that we know their *finite dimensional distributions*.

Definition 6.4.1 (Finite Dimensional Distributions). The finite dimensional distributions of a stochastic process $\{X_t\}_{t \geq 0}$ are the distributions of all vectors of the form $(X_{t_1}, \dots, X_{t_n})$ with $n > 0$ and $0 \leq t_1 \leq \dots \leq t_n$.

The finite dimensional distributions tell us a lot about the behaviour of a stochastic process. It is worth noting, however, that they do not fully specify stochastic processes. For example, Brownian motion is almost surely continuous but this property does not follow simply from specifying the finite dimensional distributions.

We can simulate the finite dimensional skeletons of a Gaussian process exactly. That is, we can generate X_{t_1}, \dots, X_{t_n} for any choice of n and t_1, \dots, t_n . This is not always true for other stochastic processes. However, we do encounter a new form of error, *discretization error*.

Definition 6.4.2 (Discretization Error). Discretization error is the error that arises from replacing a continuous object with a discrete object.

For example, if we wish to calculate the variance of the proportion of time that a Gaussian process spends above 0, or the expectation of the first time a process hits a set, A , then we will encounter an error in considering the process only at a fixed number of points.

Discretization error needs to be considered when we decide on a sampling budget. Consider, for example, an evenly spaced mesh of points $t_1 = 1/m, t_2 = 2/m, \dots, t_n = 1$. As m gets bigger, the mesh gets finer and the discretization error gets smaller. We still have statistical error, however, so we need to make sure we generate a large enough sample of realizations of the stochastic process (determined by the sample size N). The total work done by our simulation is then given by

$$\text{work} = \text{number of samples} \times \text{work to make one sample} = Nf(m).$$

Usually, f grows at least linearly in m . In the case of Cholesky decomposition, for example, it is $O(m^3)$. For a fixed level of work, we need to decide on how much effort to allocate to reducing discretization error (how large m should be) and how much effort to allocate to reducing statistical error (how large N should be). Finding the optimal tradeoff can be difficult. We will consider such tradeoffs in a number of situations.

6.5 Marginal and Conditional Multivariate Normal Distributions

The multivariate normal distribution has many attractive properties. In particular, its marginal distributions are also multivariate normal. In addition, if we condition on part of the a multivariate normal vector, the remaining values are also multivariate normal.

To see this, we need to write normal random vectors in the appropriate form. Let $\mathbf{X} \sim \mathbf{N}(\boldsymbol{\mu}, \Sigma)$. We can decompose \mathbf{X} into two parts, $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)^\top$. We can then write $\boldsymbol{\mu}$ as $(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2)^\top$ and Σ as

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}.$$

Theorem 6.5.1 (Multivariate Normal Marginal Distributions). Given $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2) \sim \mathbf{N}(\boldsymbol{\mu}, \Sigma)$,

$$\mathbf{X}_1 \sim \mathbf{N}(\boldsymbol{\mu}_1, \Sigma_{11}),$$

and

$$\mathbf{X}_2 \sim \mathbf{N}(\boldsymbol{\mu}_2, \Sigma_{22}).$$

Theorem 6.5.2 (Multivariate Normal Conditional Distributions). Given $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2) \sim \mathbf{N}(\boldsymbol{\mu}, \Sigma)$, \mathbf{X}_2 conditional on \mathbf{X}_1 is multivariate normal with

$$\mathbb{E}[\mathbf{X}_2 | \mathbf{X}_1] = \boldsymbol{\mu}_2 + \Sigma_{21}\Sigma_{11}^{-1}(\mathbf{X}_1 - \boldsymbol{\mu}_1)$$

and

$$\text{Var}(\mathbf{X}_2 | \mathbf{X}_1) = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}.$$

6.6 Interpolating Gaussian Processes

Because we know the finite dimensional distributions of Gaussian processes and know a lot about working with the multivariate normal distribution, we are able to interpolate between already simulated points of a Gaussian process. We can do this by simply drawing the new points conditional on the values that we have already generate. For example, if we have generate values for the process at 0.5 and 1, then we can generate values at the points 0.25 and 0.75 conditional on $X_{0.5}$ and X_1 .

There are a number of reasons why interpolation might be useful. One is that it might not make sense to simulate the whole stochastic process on a fine mesh (which is expensive), but rather to simulate a coarse path of the stochastic process first and then focus our efforts on a particular region of this path. For example, if we are trying to estimate the first time a stochastic process hits a particular value, then we might want to focus our simulation efforts on the part of the stochastic process that is closest to this value. We should be careful, however, as simulating in this way could introduce a bias.

Example 6.6.1 (Iteratively Updating Brownian Motion). Consider an example where we update Brownian motion in an iterative fashion. Remember, Brownian motion has mean function $\mu(t) = 0$ and covariance function

$r(s, t) = \min(s, t)$. We interpolate between points to simulate a process on an increasingly fine mesh.

Listing 6.6: Matlab code

```

1 num_levels = 10;
2
3 %Make the first two points (at 0.5 and 1)
4 t = [.5 1]; n = length(t);
5 Sigma = zeros(n,n);
6 for i = 1:n
7     for j = 1:n
8         Sigma(i,j) = min(t(i),t(j));
9     end
10 end
11 X = chol(Sigma)' * randn(2,1);
12
13 plot([0; t'],[0; X]);
14 axis([0 1 -2.5 2.5]);
15
16 %Interpolate
17 for level = 2:num_levels
18     %Make the additional mesh points
19     t_new = 1/2^level : 2/2^level : (2^level-1)/(2^level);
20     n_new = length(t_new);
21
22     %Record the time points for the whole process
23     t_temp = [t t_new];
24     n_temp = length(t_temp);
25
26     %Make a covariance matrix for the whole thing
27     Sigma_temp = zeros(n_temp,n_temp);
28     for i = 1:n_temp
29         for j = 1:n_temp
30             Sigma_temp(i,j) = min(t_temp(i),t_temp(j));
31         end
32     end
33
34     %Make the separate Sigma components
35     Sigma_11 = Sigma;
36     Sigma_21 = Sigma_temp(n+1:n_temp, 1:n);
37     Sigma_12 = Sigma_temp(1:n, n+1:n_temp);
38     Sigma_22 = Sigma_temp(n+1:n_temp, n+1:n_temp);
39

```

```

40 temp_mean = Sigma_21 * inv(Sigma_11) * X;
41 Sigma_new = Sigma_22 - Sigma_21 * inv(Sigma_11) * Sigma_12;
42 X_new = temp_mean + chol(Sigma_new)' * randn(n_new,1);
43 X = [X; X_new];
44 t = t_temp;
45 n = n_temp;
46 Sigma = Sigma_temp;
47 [dummy index] = sort(t);
48 another_dummy = waitforbuttonpress;
49
50 plot([0; t(index)'], [0; X(index)]);
51 axis([0 1 -2.5 2.5]);
52 end

```

6.7 Markovian Gaussian Processes

If a Gaussian process, $\{X_t\}_{t \geq 0}$, is Markovian, we can exploit this structure to simulate the process much more efficiently. Because $\{X_t\}_{t \geq 0}$ is Markovian, we can use we only need to know the value of X_{t_i} in order to generate $X_{t_{i+1}}$. Define

$$\sigma_{i,i+1} = \text{Cov}(X_{t_i}, X_{t_{i+1}})$$

and

$$\mu_i = \mathbb{E}X_{t_i}.$$

By theorem 6.5.1, we know that $(X_{t_i}, X_{t_{i+1}})^\top$ has a multivariate normal distribution. In particular,

$$\begin{pmatrix} X_{t_i} \\ X_{t_{i+1}} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mu_i \\ \mu_{i+1} \end{pmatrix}, \begin{pmatrix} \sigma_{i,i} & \sigma_{i,i+1} \\ \sigma_{i,i+1} & \sigma_{i+1,i+1} \end{pmatrix} \right).$$

Using theorem 6.5.2, we have

$$X_{t_{i+1}} | X_{t_i} = x_i \sim \mathcal{N} \left(\mu_i + \frac{\sigma_{i,i+1}}{\sigma_{i,i}}(x_i - \mu_i), \sigma_{i+1,i+1} - \frac{\sigma_{i,i+1}^2}{\sigma_{i,i}} \right).$$

Algorithm 6.7.1 (Generating a Markovian Gaussian Process).

- (i) Draw $Z \sim \mathcal{N}(0, 1)$. Set $X_{t_1} = \mu_1 + \sqrt{\sigma_{i,i}}Z$.
- (ii) For $i = 1, \dots, m - 1$ draw $Z \sim \mathcal{N}(0, 1)$ and set

$$X_{t_{i+1}} = \mu_i + \frac{\sigma_{i,i+1}}{\sigma_{i,i}}(x_i - \mu_i) + \left(\sqrt{\sigma_{i+1,i+1} - \frac{\sigma_{i,i+1}^2}{\sigma_{i,i}}} \right) Z.$$

Example 6.7.1 (Brownian Motion). Consider Brownian motion, which is a Markov process. Now,

$$\mu_i = \mathbb{E}X_{t_i} = 0$$

and

$$\sigma_{i,i+1} = \min(t_i, t_{i+1}) = t_i.$$

So, our updating formula is

$$X_{t_{i+1}} = X_{t_i} + \left(\sqrt{t_{i+1} - t_i} \right) Z.$$

Listing 6.7: Matlab code

```

1 m = 10^4; X = zeros(m,1);
2 h = 1/m;
3 X(1) = sqrt(h)*randn;
4
5 for i = 1:m-1
6     X(i+1) = X(i) + sqrt(h) * randn;
7 end
8
9 plot(h:h:1,X);

```

Example 6.7.2 (Ornstein-Uhlenbeck Process). The Ornstein-Uhlenbeck process has

$$\mu_i = \mathbb{E}X_{t_i} = 0$$

with

$$\sigma_{i,i+1} = \exp\{\alpha|t_i - t_{i+1}|/2\}$$

and $\sigma_{i,i} = 1$. So, our updating formula is

$$X_{t_{i+1}} = \exp\{-\alpha|t_i - t_{i+1}|/2\}X_{t_i} + \left(\sqrt{1 - \exp\{\alpha|t_i - t_{i+1}|\}} \right) Z.$$

Listing 6.8: Matlab code

```

1 alpha = 50; m = 10^4;
2
3 X = zeros(m,1); h = 1/m;
4 X(1) = randn;
5
6 for i = 1:m-1
7     X(i+1) = exp(-alpha * h / 2)*X(i)...
8         + sqrt(1 - exp(-alpha * h ))*randn;
9 end
10
11 plot(h:h:1,X);

```

6.8 Brownian Motion

One of the most fundamental stochastic processes. It is a Gaussian process, a Markov process, a Lévy process, a Martingale and a process that is closely linked to the study of harmonic functions (do not worry if you do not know all these terms). It can be used as a building block when considering many more complicated processes. For this reason, we will consider it in much more depth than any other Gaussian process.

Definition 6.8.1 (Brownian Motion). A stochastic process $\{W_t\}_{t \geq 0}$ is called Brownian motion (a Wiener process) if:

- (i) It has independent increments.
- (ii) It has stationary increments.
- (iii) $W_t \sim \mathbf{N}(0, t)$ for all $t \geq 0$.
- (iv) It has almost surely continuous sample paths. That is,

$$\mathbb{P}(\{\omega : X(t, \omega) \text{ is continuous in } t\}) = 1.$$

This definition implies that the increments of Brownian motion are normally distributed. Specifically, $W_{t+s} - W_t \sim \mathbf{N}(0, s)$. This implies the following simulation scheme.

Listing 6.9: Matlab code

```

1 m = 10^3; h = 1/m;
2 X = cumsum(sqrt(h)*randn(m,1));
3 plot(h:h:1,X);

```

The first three parts of the definition of Brownian motion are equivalent to saying Brownian motion is a Gaussian process with $\text{Cov}(X_t, X_s) = \min(t, s)$ and $\mathbb{E}X_t = 0$. However, the almost sure continuity of the paths of Brownian motion does not follow from this.

Theorem 6.8.2. The following two statements are equivalent for a stochastic process $\{X_t\}_{t \geq 0}$.

- (i) $\{X_t\}_{t \geq 0}$ has stationary independent increments and $X_t \sim \mathbf{N}(0, t)$.
- (ii) $\{X_t\}_{t \geq 0}$ is a Gaussian process with $\mu(t) = \mathbb{E}X_t = 0$ and $r(s, t) = \text{Cov}(X_t, X_s) = \min(t, s)$.

Proof.

Part 1. First, we show (i) implies (ii). In order to show this, we need to show that $(X_{t_1}, \dots, X_{t_n})$ is multivariate normal for all choices of $0 \leq t_1 \leq \dots \leq t_n$ and $n \geq 1$. In order to X_{t_1}, \dots, X_{t_n} to be multivariate normal, we need $\sum_{k=1}^n \alpha_k X_{t_k}$ to be univariate normal (for all choices of n etc.). Now,

$$\begin{aligned}
\sum_{k=1}^n \alpha_k X_{t_k} &= \alpha_1 X_{t_1} + \sum_{k=2}^n \alpha_k X_{t_k} \\
&= \alpha_1 X_{t_1} + \sum_{k=2}^n \alpha_k \left(\sum_{j=1}^k X_{t_j} - \sum_{j=1}^{k-1} X_{t_j} \right) \\
&= \alpha_1 X_{t_1} + \sum_{k=2}^n \alpha_k \left(\sum_{j=1}^k X_{t_j} - \sum_{j=2}^k X_{t_{j-1}} \right) \\
&= \sum_{k=1}^n \alpha_k X_{t_1} + \sum_{k=2}^n \sum_{j=2}^k \alpha_k (X_{t_j} - X_{t_{j-1}}) \\
&= \left(\sum_{k=1}^n \alpha_k \right) X_{t_1} + \sum_{j=2}^n \left(\sum_{k=j}^n \alpha_k \right) (X_{t_j} - X_{t_{j-1}}) \\
&= \beta_1 X_{t_1} + \sum_{j=2}^n \beta_j (X_{t_j} - X_{t_{j-1}}).
\end{aligned}$$

Because X_{t_1} and $X_{t_2} - X_{t_1}, X_{t_3} - X_{t_2}, \dots$ are independent random variables, it follows that the final equation results in a univariate normal random variable (sums of independent normals are normal). Now, we just need to check the expectation and covariance. It is easy to see that $\mathbb{E}X_t = 0$. In order to calculate the covariance, assume that $s < t$. We have that,

$$\text{Cov}(X_s, X_t) = \mathbb{E}X_s X_t - \mathbb{E}X_s \mathbb{E}X_t = \mathbb{E}X_s X_t,$$

and

$$\mathbb{E}X_s X_t = \mathbb{E}X_s (X_t - X_s + X_s) = \mathbb{E}X_s^2 + \mathbb{E}X_s (X_t - X_s)$$

and, as the independent increments property implies $\mathbb{E}X_s (X_t - X_s) = 0$,

$$\text{Cov}(X_s, X_t) = \mathbb{E}X_s^2 = s.$$

Repeating this with $t \geq s$, we get $\text{Cov}(X_s, X_t) = \min(s, t)$.

Part 2. We show that (ii) implies (i). It follows immediately from the definition in (ii) that $X_t \sim \mathbf{N}(0, t)$. It is also immediate that $X_{t+s} - X_t$ is univariate normal (as this is a linear combination of multivariate normal

random variables). We can thus show this increment is stationary by showing that mean is constant and the variance only depends on s . Now, $\mathbb{E}X_{t+s} - X_t = 0$ and

$$\text{Var}(X_{t+s} - X_s) = \text{Var}(X_t) + \text{Var}(X_s) - 2\text{Cov}(X_s, X_t) = t + s - 2t = s.$$

Thus, the increments are stationary. Because the increments are multivariate normal, we can show the increments of the process are independent if we can show the covariance between increments is 0. So, take $u < v \leq s < t$. Then,

$$\begin{aligned} & \text{Cov}(X_v - X_u, X_t - X_s) \\ &= \text{Cov}(X_v, X_t) - \text{Cov}(X_v, X_s) - \text{Cov}(X_u, X_t) + \text{Cov}(X_u, X_s) \\ &= \min(v, t) - \min(v, s) - \min(u, t) + \min(u, s) \\ &= v - v - u + u = 0. \end{aligned}$$

So, non-overlapping increments are independent.

□

Chapter 7

Random Fields

At the beginning of these notes we defined a stochastic process as a set of random variables $\{X_i\}_{i \in I}$ taking values in a state space \mathcal{X} with index set $I \subset \mathbb{R}$. It seems natural to ask what would happen if we took I to be something “larger” than a subset of \mathbb{R} — for example, a subset of \mathbb{R}^N . Such generalizations are called *random fields*. Intuitively, we can think of them as stochastic processes evolving over space instead of time. Sometimes, they can be thought of as spatial processes that are also evolving over time (for example, if $I = [a, b]^3 \times [0, T]$).

Definition 7.0.3 (Random Field). Let $I \subset \mathbb{R}^N$ and $\{\mathbf{X}_i\}_{i \in I}$ be a set of d -dimensional random vectors indexed by I . We say $\{X_i\}_{i \in I}$ is a (N, d) random field.

Note that stochastic processes are random fields.

7.1 Gaussian Random Fields

Some of the simplest random fields to work with are *Gaussian random fields*. These are a natural extension of our definition of Gaussian processes.

Definition 7.1.1. A random field $\{X_i\}_{i \in I}$ is a Gaussian random field if $(X_{i_1}, \dots, X_{i_n})$ has a multivariate normal distribution for all $1 \leq n < \infty$ and $(i_1, \dots, i_n) \in I^n$.

As with Gaussian processes, we can describe the finite dimensional distributions of Gaussian random fields using a mean function $\mu(i) = \mathbb{E}X_i$ and a covariance function $r(i, j) = \text{Cov}(X_i, X_j)$. If we wish to consider the Gaussian random field at a finite number of points, we can treat its values as a

random vector. We can use the mean function and covariance function to construct the mean vector, $\boldsymbol{\mu}$, and covariance matrix, Σ , of these points.

Using these functions, and treating a finite number of points of the random field as a vector, we can generate a realization of the field

Example 7.1.2 (Gaussian White Noise). A trivial example of a Gaussian random field is *Gaussian white noise*. On $I = \{1, \dots, m\} \times \{1, \dots, m\}$, this is the process, $\{X_i\}_{i \in I}$ with mean function

$$\mu(i) = \mathbb{E}X_i = 0,$$

and covariance function

$$r(i, j) = r((x^i, y^i), (x^j, y^j)) = \text{Cov}(X_{(x^i, y^i)}, X_{(x^j, y^j)}) = \begin{cases} 1 & \text{if } (x^i, y^i) = (x^j, y^j) \\ 0 & \text{otherwise} \end{cases}.$$

We can simulate this in Matlab by simple generating a matrix of standard

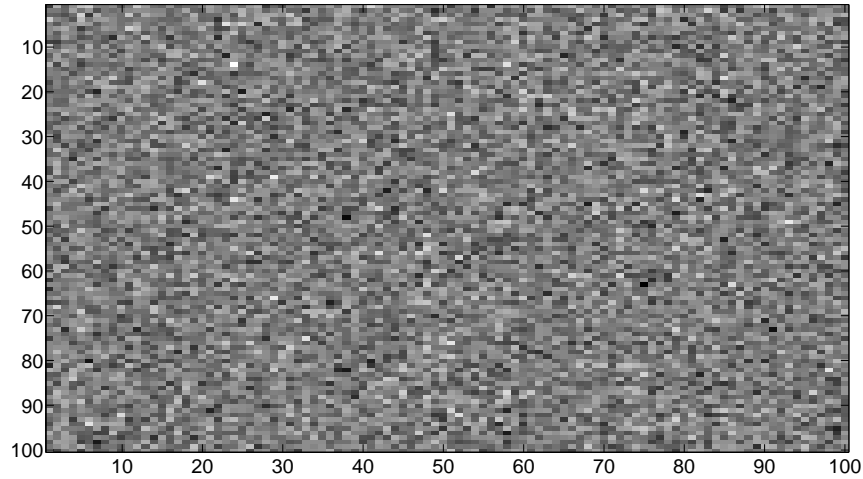


Figure 7.1.1: A realization of Gaussian white noise

normal random variables.

Listing 7.1: Matlab code

```

1 m = 100;
2 Z = randn(m,m);
3 imagesc(Z);
4 colormap(gray);

```

Example 7.1.3 (Brownian Sheet). The *Brownian sheet*, $\{W_{(x,y)}\}_{(x,y) \in [0,1]^2}$ is a natural random field extension of Brownian motion on $[0, 1]$. It has mean function

$$\mu(i) = \mathbb{E}X_i = 0,$$

and covariance function

$$r(i, j) = r((x^i, y^i), (x^j, y^j)) = \text{Cov}(W_{(x^i, y^i)}, W_{(x^j, y^j)}) = \min(x^i, x^j) \min(y^i, y^j).$$

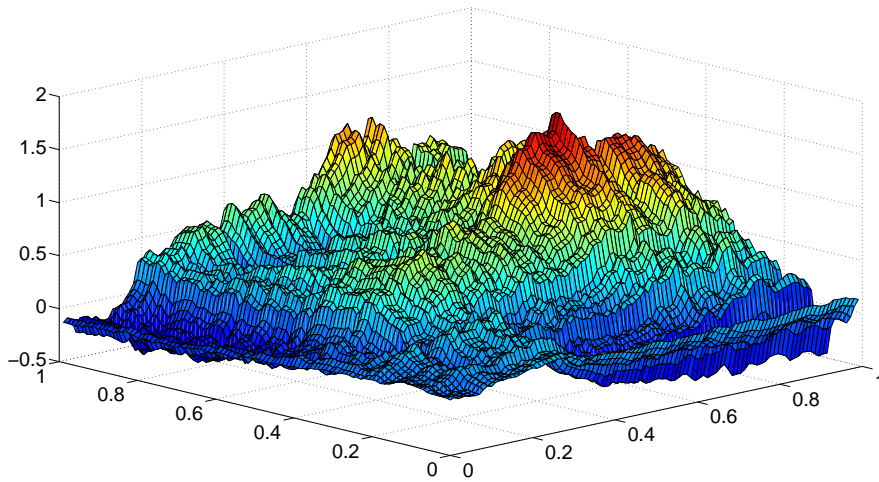


Figure 7.1.2: A realization of a Brownian sheet

We can easily simulate this in Matlab.

Listing 7.2: Matlab code

```

1 alpha = 1.5; x = .01:.01:1; y = .01:.01:1;
2
3 n1 = length(x); n2 = length(y);
4
5 Sigma = zeros(n1*n2,n1*n2);
6
7 for i = 1:n2
8     for j = 1:n1
9         row_state = (i-1)*n2 + j;
10        for k = 1:n2
11            for l = 1:n1
12                column_state = (k-1)*n2 + l;
13                Sigma(row_state,column_state) = ...

```

```

14         min(x(j),x(1))*min(y(i),y(k));
15     end
16 end
17 end
18 end
19
20 A = chol(Sigma);
21 Y = randn(1,n1*n2) * A;
22
23 X = zeros(n1,n2);
24
25 for i = 1:n2
26     for j = 1:n1
27         X(i,j) = Y((i-1)*n2 + j);
28     end
29 end
30
31 surf(x,y,X)

```

Random fields are, in general, difficult to work with. For that reason, it is nice to identify classes of random fields that are more tractable.

Definition 7.1.4 (Stationarity). We say a random field, $\{X_i\}_{i \in I}$, is stationary if $(X_{i_1}, \dots, X_{i_n})$ has the same distribution as $(X_{i_1+s}, \dots, X_{i_n+s})$ for all $n \geq 1$, $(i_1, \dots, i_n) \in I^n$ and $s \in I$.

Weak stationarity implies that the first two moments of the distributions do not change.

Definition 7.1.5 (Wide sense (weak) stationarity). A random field, $\{X_i\}_{i \in I}$ is wide-sense stationary if $\mu(i) = c$ for all $i \in I$ and $r(i, j) = r(|i - j|)$. That is, the covariance function is simply a function of the displacement vector between the points i and j .

Analogously to the case of Gaussian processes, Gaussian fields are stationary if they are weak stationary. When considering random fields, a different type of invariance property is also useful. This is called *isotropy*. Basically, this means the distribution of the object does not change when the object is rotated. We will only consider a very special case of isotropy.

Definition 7.1.6 (Stationary wide sense isotropy). A stationary random field is wide sense isotropic if $r(i, j) = r(\|i - j\|)$. That is, the covariance between points i and j only depends on the distance between them.

Note that non-stationary random fields can be isotropic.

7.2 Markov Random Fields

When considering stochastic processes, the Markov property simplified things considerably. Markov random fields are the random field analogues of Markov stochastic processes. To define a meaningful version of the Markov property, we exploit the idea of conditional independence.

Definition 7.2.1 (Conditional Independence for Events). We say two events A and B are said to be *conditionally independent* given C if

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

Clearly, a Markov process is conditional independent of its past given its present.

Definition 7.2.2 (Conditional Independence for Random Variables). We can easily extend this definition to random vectors \mathbf{X}, \mathbf{Y} and \mathbf{Z} with joint density $\pi(\mathbf{x}, \mathbf{y}, \mathbf{z})$. We say \mathbf{X} and \mathbf{Y} are conditionally independent given \mathbf{Z} (denoted $\mathbf{x} \perp \mathbf{y} | \mathbf{z}$) if

$$\pi(\mathbf{x}, \mathbf{y} | \mathbf{z}) = \pi(\mathbf{x} | \mathbf{z})\pi(\mathbf{y} | \mathbf{z}).$$

There is a nice criterion for determining conditional independence given a joint density.

Theorem 7.2.3 (Factorization Theorem). Given random vectors \mathbf{X}, \mathbf{Y} and \mathbf{Z} with joint density $\pi(\mathbf{x}, \mathbf{y}, \mathbf{z})$, we say $\mathbf{x} \perp \mathbf{y} | \mathbf{z}$ if and only if

$$\pi(\mathbf{x}, \mathbf{y}, \mathbf{z}) = f(\mathbf{x}, \mathbf{z})g(\mathbf{y}, \mathbf{z})$$

for some functions f and g and all \mathbf{z} with $\pi(\mathbf{z}) > 0$.

Example 7.2.4. Take

$$\pi(x, y, z) \propto \exp\{x + xz + yz\}$$

on some bounded region. We can write

$$\pi(x, y, z) \propto \exp\{x + xz\}\exp\{yz\} = f(x, z)g(y, z)$$

so $x \perp y | z$.

Example 7.2.5. Take

$$\pi(x, y, z) \propto \exp\{xyz\}$$

on some bounded region. We cannot factorize this into a function of x and z and a function of y and z , so x and y are not conditionally independent given z here.

We can represent conditional dependence relations using graphs called *graphical models*.

Example 7.2.6.

PICTURE HERE

This example encodes the dependency structure

$$\pi(a, b, c, d) = \pi(a)\pi(b)\pi(c | b, d)\pi(d | a, b, c).$$

Markov random fields depict a specific dependency structure using undirected graphs $G = (V, E)$. Alternatively, you can think of these as graphs with arrows on both ends of the edges.

PICTURE HERE

In Markov random fields, random variables are conditionally independent of the rest of the graph given their immediate neighbors.

Definition 7.2.7 (Neighbors of a Vertex). The neighbors of the vertex i are the members of the set

$$\mathcal{N}_i = \{j \in V : (i, j) \in E\}.$$

For example, in the graph above, $\mathcal{N}_1 = \{2, 3\}$, $\mathcal{N}_2 = \{1, 3, 4\}$, $\mathcal{N}_3 = \{1, 2, 4\}$ and $\mathcal{N}_4 = \{2, 3\}$.

We can now define a Markov random field. It will be convenient to use the following notation. For $C \subset V$, let $\mathbf{X}_C = \{\mathbf{X}_i : i \in C\}$ and $\mathbf{X}_{-C} = \{\mathbf{X}_i : i \in V \setminus C\}$.

Definition 7.2.8 (Markov Random Field). We say a sequence of random variables, $\mathbf{X} = \{X_i\}_{i \in V}$, indexed by the vertices of $G = (V, E)$, is a Markov random field if $X_i | \mathbf{X}_{- \{i\}}$ has the same distribution as $X_i | \mathbf{X}_{\mathcal{N}_i}$ for all $i \in V$.

7.3 Gaussian Random Markov Fields

We will consider a special class of Markov random fields, called *Gaussian random Markov fields*.

Definition 7.3.1 (Gaussian Random Markov Field). A Gaussian random Markov field (GRMF) is a Markov random field which is also a Gaussian field.

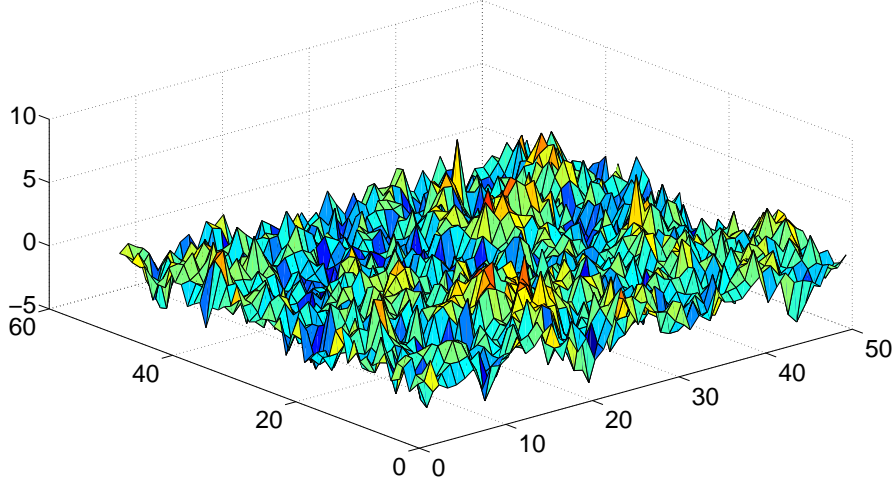


Figure 7.3.1: A realization of a Gaussian Markov random field on a 50×50 lattice

We will assume that the GRMFs we consider have densities. That is, the covariance matrix of a finite number of points, Σ , will always be positive definite. Recall that the pdf of a normal vector \mathbf{X} is given by

$$f(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = (2\pi)^{-n/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

In general, the covariance matrix is a natural object to consider when talking about Gaussian processes and fields. However, in the case of Gaussian random fields, it turns out that the *precision matrix* is a more natural object to work with.

Definition 7.3.2 (Precision Matrix). The precision matrix, Q , of a covariance matrix, Σ , is its inverse. That is,

$$Q = \Sigma^{-1}.$$

We can rewrite the pdf of a normal vector in terms of its precision matrix as

$$f(\mathbf{x}; \boldsymbol{\mu}, Q) = (2\pi)^{-n/2} |Q|^{1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top Q (\mathbf{x} - \boldsymbol{\mu}) \right\}.$$

The following result reveals the reason why precision matrices are so appealing when working with GRMFs.

Theorem 7.3.3. Let \mathbf{X} be normally distributed with mean $\boldsymbol{\mu}$ and precision matrix Q . Then, for $i \neq j$,

$$X_i \perp X_j \mid \mathbf{X}_{-\{i,j\}} \iff Q_{i,j} = 0.$$

Proof. Assume without loss of generality that $\boldsymbol{\mu} = \mathbf{0}$. We have

$$\begin{aligned} \pi(\mathbf{x}) &\propto \exp \left\{ -\frac{1}{2} \mathbf{x}^\top Q \mathbf{x} \right\} = \exp \left\{ -\frac{1}{2} \sum_{l \in V} \sum_{k \in V} x_k Q_{k,l} x_l \right\} \\ &= \exp \left\{ -\frac{1}{2} x_i x_j (Q_{i,j} + Q_{j,i}) - \frac{1}{2} \sum_{\{k,l\} \neq \{i,j\}} x_k Q_{k,l} x_l \right\} \end{aligned}$$

If you consider examples 7.2.5 and 7.2.6, then it is clear that this pdf can only be factorized in the necessary way if $Q_{i,j} = Q_{j,i} = 0$. \square

Theorem 7.3.4. Let \mathbf{X} be a GRMF with respect to $G = (V, E)$ with mean $\boldsymbol{\mu}$ and precision matrix Q . Then,

- (i) $\mathbb{E}[X_i \mid \mathbf{X}_{-i}] = \mu_i - \frac{1}{Q_{i,i}} \sum_{j \in \mathcal{N}_i} Q_{i,j} (X_j - \mu_j)$.
- (ii) $\text{Prec}(X_i \mid \mathbf{X}_{-\{i\}}) = Q_{i,i}$.
- (iii) $\text{Corr}(X_i, X_j \mid \mathbf{X}_{i,j}) = \frac{-Q_{i,j}}{\sqrt{Q_{i,i} Q_{j,j}}}$.

Theorem 7.3.5. Let \mathbf{X} be a GRMF with respect to $G = (V, E)$. Then, the following are equivalent:

- (i) The *pairwise Markov property*

$$X_i \perp X_j \mid \mathbf{X}_{-\{i,j\}} \quad (i, j) \notin E, i \neq j.$$

- (ii) The *local Markov property*

$$X_i \perp \mathbf{X}_{-\{i\} \cup \mathcal{N}_i} \mid \mathbf{X}_{\mathcal{N}_i} \quad \forall i \in V.$$

- (iii) The *global Markov property*

$$\mathbf{X}_A \perp \mathbf{X}_B \mid \mathbf{X}_C$$

for all disjoint A, B and C , where C separates A and B .

GRMFs are very appealing from a computational perspective as the precision matrix Q is sparse (that is, it is mostly zeros). In general, matrix operations are much faster for sparse matrices. In addition, far less memory is required in order to store sparse matrices.

Remember that the inverse of a SPD matrix is also SPD. This implies that Q has a Cholesky decomposition. We can use this to generate normal random vectors with the desired distributions.

Theorem 7.3.6. Given a SPD covariance matrix Σ and a vector $\boldsymbol{\mu}$, the random vector

$$\mathbf{X} = \boldsymbol{\mu} + (C^\top)^{-1} \mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma),$$

where C is such that $\Sigma^{-1} = Q = CC^\top$ and $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, I)$.

Proof. We know \mathbf{X} is multivariate normal, so we just need to check that it has the correct mean and variance. Now,

$$\mathbb{E} [\boldsymbol{\mu} + (C^\top)^{-1} \mathbf{Z}] = \boldsymbol{\mu} + (C^\top)^{-1} \mathbb{E} \mathbf{Z} = \boldsymbol{\mu},$$

and

$$\text{Var} (\boldsymbol{\mu} + (C^\top)^{-1} \mathbf{Z}) = (C^\top)^{-1} \text{Var}(\mathbf{Z}) (C)^{-1} = (CC^\top)^{-1} = Q^{-1} = \Sigma.$$

□

Example 7.3.7 (A Gaussian random Markov field on a Lattice). We simulate a zero mean GRMF on the 200×200 square lattice with $Q_{i,i} = 1$ and

$$Q_{i,j} = \begin{cases} -0.25 & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}.$$

Listing 7.3: Matlab code

```

1 m = 200; d1 = 1; d2 = -0.25;
2 nels = m*(5*m-4);
3
4 a = zeros(1, nels); b = zeros(1,nels); q = zeros(1,nels);
5 %compute the links and weights for the precision matrix
6 k=0;
7 for i=1:m
8     for j=1:m
9         A = findneigh(i,j,m);
10        number_neighbours = size(A,1);
11        for h=1:number_neighbours
12            a(k+h)= ij2k(i,j,m);

```

```

13         b(k+h)= ij2k(A(h,1),A(h,2),m);
14         if h==1
15             q(k+h) = d1;
16         else
17             q(k+h) = d2;
18         end
19     end
20     k = k+number_neighbours;
21 end
22 end
23 %construct the precision matrix
24 Q = sparse(a,b,q,m^2,m^2);
25 %calculate the Cholesky matrix
26 C = chol(Q,'lower');
27 Z = randn(m^2,1);
28 % generate the Gaussian process
29 x = C\Z;
30 colormap gray, brighten(-0.2)
31 imagesc(reshape(x,m,m)) % plot the result

```

This uses the following functions.

Listing 7.4: Matlab code

```

1 function A = findneigh(i,j,m)
2 % find neighbors of the (i,j)-th site of an m by m square lattice
3 if i==1
4     if j==1
5         A = [1,1;1,2;2,1];
6     elseif j==m
7         A = [1,m;1,m-1;2,m];
8     else
9         A = [1,j;1,j-1;1,j+1;2,j];
10    end
11 elseif i==m
12    if j==1
13        A = [m,1;m,2;m-1,1];
14    elseif j==m
15        A = [m,m;m,m-1;m-1,m];
16    else
17        A = [m,j;m,j-1;m,j+1;m-1,j];
18    end
19 else
20    if j==1
21        A = [i,1;i,2;i-1,1;i+1,1];

```

```
22     elseif j==m
23         A = [i,m;i,m-1;i+1,m;i-1,m];
24     else
25         A = [i,j;i,j-1;i,j+1;i+1,j;i-1,j];
26     end
27 end
28 end
```

Listing 7.5: Matlab code

```
1 function k = ij2k(i,j,m)
2 k = (i-1)*m + j;
```

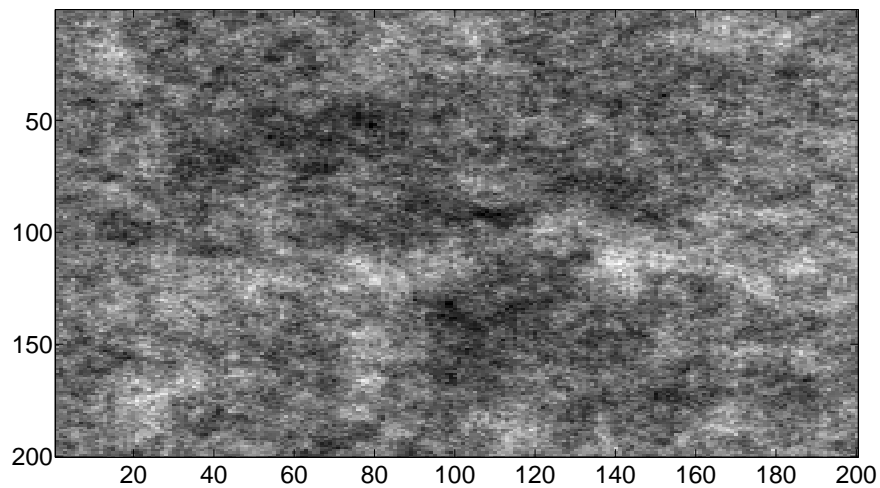


Figure 7.3.2: A realization of the random field

Bibliography

- [1] J. Norris. *Markov Chains*. Cambridge University Press, Cambridge, 1997.