

STAT 870: Applied Probability — Course notes

Richard A. Lockhart

September 4, 2006

Contents

1	Introduction	5
1.1	Probability Definitions	9
1.2	Random Variables	11
1.3	Appendix on Measurability	12
1.4	Appendix on Lebesgue Measure	13
2	Independence, conditional distributions	15
3	Expectation	23
3.1	Moments and independence	25
3.2	Appendix on Lebesgue Integration	27
4	The Strong Law of Large Numbers	31
4.1	Events in Set Notation	31
4.2	The Strong Law of Large Numbers with 4 moments	33
4.3	Proof without 4 finite moments	36
4.4	Consistency of MLE	42
5	Markov Chains	53
5.1	Chapman-Kolmogorov Equations	54
5.1.1	Extensions of the Markov Property	55
5.2	Classification of States	56
5.2.1	Conditional distributions and expectations	59
5.3	Initial Distributions	63
5.4	Recurrence and Transience	64
5.5	Mean return times	67
5.6	The ergodic theorem	69
5.7	Hitting Times	72
6	Continuous Time Markov Chains	89

Chapter 1

Introduction

Imagine tossing a coin n times. On trial k write down a 1 if the coin comes up heads and a 0 if it comes up tails. A typical outcome of the experiment will be a sequence, $\omega = (\omega_1, \dots, \omega_n)$ of zeros and ones. For $n = 3$ for instance there are eight possible outcomes, the elements of the set

$$\Omega = \{(0, 0, 0), (0, 0, 1), (0, 1, 0), (0, 1, 1), (1, 0, 0), (1, 0, 1), (1, 1, 0), (1, 1, 1)\}.$$

For the general case the set of all possible outcomes is $\Omega = \{0, 1\}^n$ which is just notation for the set of all sequences of n elements chose from $\{0, 1\}$; there are 2^n elements in Ω .

Before we do the experiment, we think of it as *random*. I am not going, in this course, to discuss the meaning of the word random, nor the various interpretations of the idea of probability. Instead I want to focus on the mathematics and the modelling. A probability measure on Ω is a function P defined on the set of all subsets of Ω with the following properties:

1. For each subset A of Ω , $P(A)$ is a real number in the set $[0, 1]$.
2. If A_1, \dots, A_k are *pairwise disjoint* (meaning that for $i \neq j$ the intersection $A_i \cap A_j$ which we usually write as $A_i A_j$ is the empty set \emptyset) then

$$P(\cup_1^k A_j) = \sum_1^k P(A_j)$$

3. $P(\Omega) = 1$.

Probability modelling is the process of selecting one (or a family) of possible probability measures which makes a good match between the mathematics and the real world phenomenon being described. What constitutes a good match depends on the interpretation of probability being used. In general in this course a long run limiting relative frequency interpretation of $P(A)$ will make the intuition underlying the modelling easiest.

There are actually quite a few possible probability measures on Ω in the examples given. How do we “select” one and how do we explain which one we have selected? Consider for

instance the case of $n = 3$. The set Ω has 8 elements and there are $2^8 = 256$ subsets of Ω . To fully specify a function P on the set of all subsets of Ω , we ought, it seems to specify 256 numbers.

That sort of specification is quite impractical. Instead we *model* by listing some assumptions about the function P . We then try to deduce what P is, or show how to calculate $P(A)$ for any A or at least for lots of interesting A .

Here are three approaches to modelling the coin tossing problem.

The first approach is, essentially, to provide a complete description of P , namely:

$$P(A) = \frac{\text{number of elements in } A}{\text{number of elements in } \Omega}. \quad (1.1)$$

This model has the disadvantage of not telling us how to approach very many other problems.

In fact, this first description is really a deduction from a fairly common model in games of chance, the equally likely outcomes model. In this case our model specifies that if A and B are two singleton sets in Ω , that is, $A = \{\omega_1\}$ and $B = \{\omega_2\}$, then $P(A) = P(B)$. If there are m elements in Ω , say $\Omega = (\omega_1, \dots, \omega_m)$ then we have

$$\begin{aligned} P(\Omega) &= P(\cup_1^m \{\omega_j\}) \\ &= \sum_1^m P(\{\omega_j\}) \\ &= mP(\{\omega_1\}). \end{aligned}$$

Since $P(\Omega) = 1$ we see that $P(\{\omega_i\}) = 1/m$ for each i and then a similar argument proves that this model is the same as the model (1.1).

These two approaches to the fair coin tossing model share a serious defect. They do not naturally make it easy to conceive of experiments in which the set of outcomes is infinite. Consider, for instance, the simple game of tossing until you see a head. It is natural to think of Ω as consisting of all sequences of some number k of zeros followed by a one. Or you could just think of Ω as the set $\{0, 1, 2, \dots\}$ by keeping track of how many tails came before the first head. Either way the natural space Ω is infinite and it is clearly not going to be adequate to assume that each element of Ω has the same probability.

Instead we will describe another method of modelling our coin tossing problem. Consider the model (1.1) for $n = 3$. Consider the event $A = \{\omega : \omega_1 = 1, \omega_2 = 0, \omega_3 = 1\}$. We can write A as the intersection of three events each defined in terms of a different co-ordinate in ω . Put

$$A_1 = \{\omega : \omega_1 = 1\}$$

$$A_2 = \{\omega : \omega_2 = 0\}$$

and

$$A_3 = \{\omega : \omega_3 = 1\}.$$

Then

$$A = A_1 \cap A_2 \cap A_3.$$

Note that $P(A) = 1/8$ and $P(A_i) = 1/2$ for each i so that

$$P(A) = \prod P(A_i).$$

This is a special case of a very general property of our equally likely outcomes model. Consider n tosses. For $i = 1, \dots, n$ let B_i be some subset of $\{0, 1\}$. (In our little example $B_1 = B_3 = \{1\}$ and $B_2 = \{0\}$.) Define

$$A_i = \{\omega : \omega_i \in B_i\}$$

and

$$A = \cap A_i.$$

It is then possible to prove that

$$P(A) = \prod P(A_i).$$

Later in the course we will describe this by saying that the random variables X_i defined by $X_i(\omega) = \omega_i$ are independent.

This property forms the basis of the most common modelling tactic for the coin tossing problem. We *assume*

$$P(\{\omega : \omega_i = 1\}) = P(\{\omega : \omega_i = 0\}) = 1/2 \tag{1.2}$$

and that for any set of events of the form given above

$$P(A) = \prod P(A_i). \tag{1.3}$$

The second assumption is a natural deduction of the long run relative frequency interpretation together with a physical assumption about cause and effect, namely, that the outcome of one toss of the coin is incapable of influencing the outcome of another toss.

The advantage of the second set of assumptions is that it permits generalization in many ways and, in particular, allows infinite spaces Ω . Most of our probability models are like this. We will assume that our probability satisfies identities such as (1.3) and then deduce other properties of P .

Here is the idea. Let us imagine tossing the coin an infinite number of times. Now $\Omega = \{\omega = (\omega_1, \omega_2, \dots)\}$ is an uncountably infinite set. We will assume that for any n and any event of the form $A = \cap_1^n A_i$ with each $A_i = \{\omega : \omega_i \in B_i\}$ we have

$$P(A) = \prod_1^n P(A_i). \tag{1.4}$$

This is the assumption that the outcomes of different tosses are independent. This assumption is already strong enough to deduce lots of results. To make it into the fair coin tossing model we add the assumption that

$$P(\{\omega : \omega_i = 1\}) = 1/2. \tag{1.5}$$

What values of $P(A)$ can be determined from these assumptions? Consider first, any set $A = \{\omega \in \Omega : (\omega_1, \dots, \omega_n) \in B\}$ where $B \subset \Omega_n = \{0, 1\}^n$. You can prove that the assumptions given guarantee that

$$P(A) = \frac{\text{number of elements in } B}{\text{number of elements in } \Omega_n}. \quad (1.6)$$

In other words, our model specifies that the first n of our infinite sequence of tosses behave like the equally likely outcomes model we started with. What about C_k which is the event that the first head occurs after k consecutive tails? You can write C_k as the intersection

$$C_k = A_1^c \cap A_2^c \cdots \cap A_k^c \cap A_{k+1}$$

where $A_i = \{\omega : \omega_i = 1\}$ and the notation A^c means the complement of A so that $A_i^c = \{\omega : \omega_i = 0\}$. Our assumption guarantees

$$\begin{aligned} P(C_k) &= P(A_1^c \cap A_2^c \cdots \cap A_k^c \cap A_{k+1}) \\ &= P(A_1^c) \cdots P(A_k^c) P(A_{k+1}) \\ &= 2^{-(k+1)}. \end{aligned}$$

It is possible to describe much more complicated events. Here are several which show up in the theorem known as the Strong Law of Large Numbers. See Homework set 1 for problems connected with these events.

$$A_1 \equiv \{\omega : \lim_{n \rightarrow \infty} (\omega_1 + \cdots + \omega_n)/n \text{ exists} \}$$

$$A_2 \equiv \{\omega : \lim_{n \rightarrow \infty} (\omega_1 + \cdots + \omega_n)/n = 1/2\}$$

$$A_3 \equiv \{\omega : \lim_{n \rightarrow \infty} \sum_1^n (2\omega_k - 1)/k \text{ exists} \}.$$

The Strong Law of Large Numbers says, among other things, that for our model $P(A_2) = 1$. In fact, it turns out that $A_3 \subset A_2 \subset A_1$. If $P(A_2) = 1$ then our assumptions guarantee that $P(A_1) = 1$ because $P(A_3)$ must be bigger than $P(A)$ for any subset A of A_1 . The proof of the strong law of large numbers actually proves that $P(A_3) = 1$ and so concludes also $P(A_2) = P(A_1) = 1$.

Before we get to that theorem we have some mathematical questions to answer:

1. Do our assumptions (1.4) and (1.5) determine the value of $P(A)$ for any subset A of Ω ? [The answer is no.]
2. Do our assumptions determine the value of $P(A_i)$ for $i = 1, 2, 3$? [The answer is yes.]
3. Are our assumptions about the probability P logically consistent? That is, is there any probability measure satisfying (1.4) and (1.5)? [The answer is yes.]

1.1 Probability Definitions

A **Probability Space** is an ordered triple (Ω, \mathcal{F}, P) . The idea is that Ω is the set of possible outcomes of a random experiment, \mathcal{F} is the set of those events, or subsets of Ω whose probability is defined and P is the rule for computing probabilities. Formally:

- Ω is a set.
- \mathcal{F} is a family of subsets of Ω with the property that \mathcal{F} is a σ -field (or σ -algebra):
 1. The empty set \emptyset and Ω are members of \mathcal{F} .
 2. \mathcal{F} is closed under complementation. That is, if A is in \mathcal{F} (meaning $P(A)$ is defined) then $A^c = \{\omega \in \Omega : \omega \notin A\}$ is in \mathcal{F} (because we want to be able to say $P(A^c) = 1 - P(A)$).
 3. If A_1, A_2, \dots are all in \mathcal{F} then so is $A = \cup_{i=1}^{\infty} A_i$. (A is the event that at least one of the A_i happens and we want to be sure that if each of the A_i has a probability then so does this event A .)
- P is a function whose domain is \mathcal{F} and whose range is a subset of $[0, 1]$ which satisfies the axioms for a probability:
 1. $P(\emptyset) = 0$ and $P(\Omega) = 1$.
 2. If A_1, A_2, \dots are **pairwise disjoint** (or **mutually exclusive**) (meaning for any $j \neq k$ $A_j \cap A_k = \emptyset$) then

$$P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i).$$

This property is called **countable additivity**.

These axioms guarantee that as we compute probabilities by the usual rules, including approximation of an event by a sequence of others we don't get caught in any logical contradictions. Here are some consequences of the axioms for which you should provide a proof.

1. If A_1, A_2, \dots, A_n are pairwise disjoint then

$$P(\cup_{i=1}^n A_i) = \sum_{i=1}^n P(A_i).$$

This property is called **finite additivity**.

2. For any event A

$$P(A^c) = 1 - P(A).$$

3. If $A_1 \subset A_2 \subset \dots$ are events then

$$P\left(\bigcup_1^\infty A_i\right) = \lim_{n \rightarrow \infty} P(A_n).$$

4. If $A_1 \subset A_2 \subset \dots$ then

$$P(\cap_1^\infty A_i) = \lim_{n \rightarrow \infty} P(A_n).$$

The most subtle point of the definitions is the σ -field, \mathcal{F} . This ingredient is needed to avoid some contradictions which arise if you try to define $P(A)$ for every subset A of Ω when Ω is a set with uncountably many elements.

A vector valued random variable is a function X whose domain is Ω and whose range is in some p dimensional Euclidean space, \mathbb{R}^p with the property that the events whose probabilities we would like to calculate from their definition in terms of X are in \mathcal{F} . We will write $X = (X_1, \dots, X_p)$. We will want to make sense of

$$P(X_1 \leq x_1, \dots, X_p \leq x_p)$$

for any constants (x_1, \dots, x_p) . In our formal framework the notation

$$X_1 \leq x_1, \dots, X_p \leq x_p$$

is just shorthand for an event, that is a subset of Ω , defined as

$$\{\omega \in \Omega : X_1(\omega) \leq x_1, \dots, X_p(\omega) \leq x_p\}$$

Remember that X is a function on Ω so that X_1 is also a function on Ω . In almost all of probability and statistics the dependence of a random variable on a point in the probability space is hidden! You almost always see X not $X(\omega)$.

Now for formal definitions:

The **Borel** σ -field in \mathbb{R}^p is the smallest σ -field in \mathbb{R}^p containing every open ball

$$B_y(r) = \{x \in \mathbb{R}^p : |x - y| < r\}.$$

(To see that there is, in fact, such a “smallest” σ -field you prove the following assertions:

1. The intersection of an arbitrary family of σ -fields is a σ -field. [Homework set 1.]
2. There is at least one σ -field of subsets of \mathbb{R}^p containing every open ball. [Homework set 1.]

Now define the Borel σ -field in \mathbb{R}^p to be

$$\mathcal{B}(\mathbb{R}^p) = \cap \mathcal{F}$$

where the intersection runs over all σ -fields \mathcal{F} which contain every open ball.)

Every common set is a Borel set, that is, in the Borel σ -field. For instance, let O be an open set. Then I will prove that O is Borel. For each x in O there is a point y all of whose co-ordinates are rational numbers and a rational number r such that

$$x \in B_y(r) \subset O$$

Now O is the union of all these $B_y(r)$. (Every $x \in O$ is in one of the $B_y(r)$ and every point in any $B_y(r)$ is in O .) But the union is countable because there are only countably many possible pairs (y, r) with all the co-ordinates rational numbers. [Homework set 1.]

1.2 Random Variables

Definition: An \mathbb{R}^p valued **random variable** is a map $X : \Omega \mapsto \mathbb{R}^p$ such that when A is Borel then $\{\omega \in \Omega : X(\omega) \in A\} \in \mathcal{F}$.

Fact: This is equivalent to

$$\{\omega \in \Omega : X_1(\omega) \leq x_1, \dots, X_p(\omega) \leq x_p\} \in \mathcal{F}$$

for all $(x_1, \dots, x_p) \in \mathbb{R}^p$. [Homework set 1.]

Jargon and notation: we write $P(X \in A)$ for $P(\{\omega \in \Omega : X(\omega) \in A\})$ and define the **distribution** of X to be the map

$$A \mapsto P(X \in A)$$

which is a probability on the set \mathbb{R}^p with the Borel σ -field rather than the original Ω and \mathcal{F} .

The **Cumulative Distribution Function** (or CDF) of X is the function F_X on \mathbb{R}^p defined by

$$F_X(x_1, \dots, x_p) = P(X_1 \leq x_1, \dots, X_p \leq x_p)$$

Properties of F_X (or just F when there's only one CDF under consideration):

1. $0 \leq F(x) \leq 1$.
2. $x > y \Rightarrow F(x) \geq F(y)$ (F is monotone non-decreasing).
3. $\lim_{x \rightarrow -\infty} F(x) = 0$
4. $\lim_{x \rightarrow \infty} F(x) = 1$
5. $\lim_{x \searrow y} F(x) = F(y)$ (F is right continuous).
6. $\lim_{x \nearrow y} F(x) \equiv F(y-)$ exists.
7. $F(x) - F(x-) = P(X = x)$.
8. $F_X(t) = F_Y(t)$ for all t implies that X and Y have the same distribution, that is, $P(X \in A) = P(Y \in A)$ for any (Borel) set A .

The distribution of a random variable X is **discrete** (we also call the random variable discrete) if there is a countable set x_1, x_2, \dots such that

$$P(X \in \{x_1, x_2, \dots\}) = 1 = \sum_i P(X = x_i).$$

In this case the **discrete density** or **probability mass function** of X is

$$f_X(x) = P(X = x).$$

The distribution of a random variable X is **absolutely continuous** if there is a function f such that

$$P(X \in A) = \int_A f(x) dx$$

for any (Borel) set A . This is a p dimensional integral in general. This condition is equivalent (when $p = 1$) to

$$F(x) = \int_{-\infty}^x f(y) dy$$

or for general p to

$$F(x_1, \dots, x_p) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_p} f(y_1, \dots, y_p) dy_1 \cdots dy_p.$$

We call f the **density** of X . For most values of x we then have F is differentiable at x and, for $p = 1$

$$F'(x) = f(x)$$

or in general

$$\frac{\partial^p}{\partial x_1 \cdots \partial x_p} F(x_1, \dots, x_p) = f(x_1, \dots, x_p).$$

1.3 Appendix on Measurability

1. A map f from \mathbb{R}^p to \mathbb{R}^q is Borel if for each Borel set $B \subset \mathbb{R}^q$ the inverse image

$$f^{-1}(B) = \{x \in \mathbb{R}^p : f(x) \in B\}$$

is a Borel set. Students are sometimes confused between inverse functions which exist only when f is bijective and inverse images like $f^{-1}(B)$ which are always defined.

2. If f from \mathbb{R}^p to \mathbb{R}^q and g from \mathbb{R}^q to \mathbb{R}^r are Borel then so is $h = g \circ f$, the composition of g and f from \mathbb{R}^p to \mathbb{R}^r .
3. Every continuous function is Borel.
4. If $f : \mathbb{R}^p \mapsto \mathbb{R}^{q_1}$ and $g : \mathbb{R}^p \mapsto \mathbb{R}^{q_2}$ are Borel then $h = (f, g) : \mathbb{R}^p \mapsto \mathbb{R}^{q_1+q_2}$ is Borel.
5. f is Borel if and only if $f^{-1}(O)$ is Borel for each open set O . The analogue for closed sets holds too.
6. If f_n is a sequence of Borel functions then

$$g \equiv \limsup f_n$$

is Borel. Similarly for $\liminf f_n$ and many other possibilities.

7. If X is an \mathbb{R}^p valued rv and f is a Borel map from \mathbb{R}^p to \mathbb{R}^q then $Y = f(X)$ is an \mathbb{R}^q valued random variable.

As an application suppose that $f_n, n = 1, 2, \dots$ is a sequence of non-negative real valued Borel functions. Define

$$g(x) = \sum_1^{\infty} f_n(x)$$

The sum might be infinite for some x values; all the concepts above can be adapted to permit functions to be $+\infty$ on a Borel set and $-\infty$ on some other Borel set and still use the Borel jargon but we will just assume the sum is finite for all x . Let g_n be the n th partial sum:

$$g_n = \sum_1^n f_m.$$

If we can prove each g_m is measurable then (6) above shows $g = \lim g_n = \limsup g_n$ is Borel. The function $h : \mathbb{R}^2 \mapsto \mathbb{R}$ defined by

$$h(x, y) = x + y$$

is continuous so Borel. If g_n is Borel for some n then $\phi = (g_n, f_{n+1})$ is Borel by (4) above and then $g_{n+1} = h \circ \phi$ is Borel by (2). Since $g_1 = f_1$ is Borel we conclude by induction that every g_n is Borel.

1.4 Appendix on Lebesgue Measure

There is a probability measure λ defined on the Borel subsets of $[0, 1]$ which agrees with length for intervals: $\lambda([a, b]) = b - a$ for $0 \leq a \leq b \leq 1$. The definition of λ can be extended to arbitrary Borel subsets of \mathbb{R} ; a p dimensional generalization of volume is also possible. This extension to \mathbb{R} has all the properties of a probability measure except that $\lambda(\mathbb{R}) = \infty$ not 1. Such an object is a **measure**. The measure λ is **translation invariant**:

$$\lambda(c + B) = \lambda(B)$$

for any Borel set B and real number c . The notation $c + B$ means

$$\{y : y = c + x \text{ for some } x \in B\}.$$

A Borel set B is called a Lebesgue null set if $\lambda(B) = 0$. Examples include:

1. The rational numbers \mathbb{Q} .
2. The Cantor set (all real numbers x in $[0, 1]$ whose expansions in base 3 may be written without the digit 1). The set shows up in lots of math courses as a set whose cardinality is the same as that of \mathbb{R} but whose Lebesgue measure is 0.

If $A \subset B$ and B is a Lebesgue null Borel set then it is natural to define $\lambda(A) = 0$ even if A is not Borel. We call all such A Lebesgue null sets. A set $A \subset \mathbb{R}$ is **Lebesgue measurable** if we can write $A = B \cup N$ with B Borel and N a Lebesgue null set. The family of all Lebesgue measurable sets is a σ -field (which is much larger than the Borel σ -field).

A property of a function $f(x)$ which holds except for a set N of x which is a Lebesgue null set is said to hold **almost everywhere** or *for almost all x* .

Chapter 2

Independence, conditional distributions

Definition: Events A and B are independent if

$$P(AB) = P(A)P(B).$$

(Note the notation: AB is the event that both A and B happen. It is also written $A \cap B$.)

Definition: Events A_i , $i = 1, \dots, p$ are **independent** if

$$P(A_{i_1} \cdots A_{i_r}) = \prod_{j=1}^r P(A_{i_j})$$

for any set of distinct indices i_1, \dots, i_r between 1 and p .

Example: $p = 3$

$$\begin{aligned} P(A_1 A_2 A_3) &= P(A_1)P(A_2)P(A_3) \\ P(A_1 A_2) &= P(A_1)P(A_2) \\ P(A_1 A_3) &= P(A_1)P(A_3) \\ P(A_2 A_3) &= P(A_2)P(A_3). \end{aligned}$$

You need all these equations to be true for independence!

Example: Toss a coin twice. If A_1 is the event that the first toss is a Head, A_2 is the event that the second toss is a Head and A_3 is the event that the first toss and the second toss are different then $P(A_i) = 1/2$ for each i and for $i \neq j$

$$P(A_i \cap A_j) = \frac{1}{4}$$

but

$$P(A_1 \cap A_2 \cap A_3) = 0 \neq P(A_1)P(A_2)P(A_3).$$

Definition: Random variables X and Y are **independent** if

$$P(X \in A; Y \in B) = P(X \in A)P(Y \in B)$$

for all A and B .

Definition: Random variables X_1, \dots, X_p are **independent** if

$$P(X_1 \in A_1, \dots, X_p \in A_p) = \prod P(X_i \in A_i)$$

for any choice of A_1, \dots, A_p .

Definition: σ -fields $\mathcal{F}_1, \dots, \mathcal{F}_p$ are **independent** if

$$P(A_1 \cap \dots \cap A_p) = \prod P(A_i)$$

for any choice of events $A_1 \in \mathcal{F}_1, \dots, A_p \in \mathcal{F}_p$.

Definition: If $X \in \mathbb{R}^p$ is a random variable we define its cumulative distribution F_X by

$$F_X(x) = P(X \leq x) = P(X_1 \leq x_1, \dots, X_p \leq x_p)$$

for any $x = (x_1, \dots, x_p) \in \mathbb{R}^p$.

Definition: We say that an \mathbb{R}^p valued random variable X has density f_X if f_X is a Borel measurable function defined on \mathbb{R}^p such that

$$P(X \in A) = \int_A f(x) dx$$

for every Borel set $A \subset \mathbb{R}^p$. (The integral is a Lebesgue integral; see the Appendix to this chapter.)

Theorem 1 *Suppose X and Y are two random variables. Assume X takes values in \mathbb{R}^p and Y takes values in \mathbb{R}^q .*

1. *If X and Y are independent then*

$$F_{X,Y}(x, y) = F_X(x)F_Y(y)$$

for all x, y . Here $F_X(x) = P(X \leq x)$ and inequalities are defined co-ordinatewise.

2. *If X and Y are independent and have joint density $f_{X,Y}(x, y)$ then X and Y have densities, say f_X and f_Y , and, for almost every pair (x, y)*

$$f_{X,Y}(x, y) = f_X(x)f_Y(y).$$

3. *If X and Y are independent and have marginal densities f_X and f_Y then (X, Y) has joint density $f_{X,Y}(x, y)$ given by*

$$f_{X,Y}(x, y) = f_X(x)f_Y(y).$$

4. *If*

$$F_{X,Y}(x, y) = F_X(x)F_Y(y)$$

for all x, y then X and Y are independent.

5. If (X, Y) has density $f(x, y)$ and there are functions $g(x)$ and $h(y)$ such that

$$f(x, y) = g(x)h(y)$$

for **almost all** (x, y) then X and Y are independent and they each have a density given by

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

and

$$f_Y(y) = h(y) / \int_{-\infty}^{\infty} h(u) du.$$

Proof:

1. Since X and Y are independent so are the events $X \leq x$ and $Y \leq y$; hence

$$P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y)$$

2. Suppose $X \in \mathbb{R}^p$ and $Y \in \mathbb{R}^q$. Let $B \subset \mathbb{R}^q$ be Borel. The event $Y \in B$ is the same event as $(X, Y) \in \mathbb{R}^p \times B$. Hence

$$\begin{aligned} P(Y \in B) &= P((X, Y) \in \mathbb{R}^p \times B) \\ &= \int_{\mathbb{R}^p \times B} f_{X,Y}(x, y) dx dy \\ &= \int_B g(y) dy \end{aligned}$$

where the function g is defined by

$$g(y) = \int_{\mathbb{R}^p} f_{X,Y}(x, y) dx$$

The function g satisfies the definition of density showing that Y has a density. Similarly X has a density. Then for any Borel sets $A \subset \mathbb{R}^p$ and $B \subset \mathbb{R}^q$

$$\begin{aligned} P(X \in A, Y \in B) &= \int_A \int_B f_{X,Y}(x, y) dy dx \\ P(X \in A)P(Y \in B) &= \int_A f_X(x) dx \int_B f_Y(y) dy \\ &= \int_A \int_B f_X(x) f_Y(y) dy dx \end{aligned}$$

See the appendix on Lebesgue integration for the theorems of Fubini and Tonelli which give conditions under which the double integral is equal to the iterated integral. These theorems are being used here.

Since $P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$ we see that for any sets A and B

$$\int_A \int_B [f_{X,Y}(x, y) - f_X(x)f_Y(y)] dy dx = 0$$

It follows (see the appendix) that the quantity in [] is 0 (for almost every pair (x, y)).

3. For any A and B we have

$$\begin{aligned} P(X \in A, Y \in B) &= P(X \in A)P(Y \in B) \\ &= \int_A f_X(x)dx \int_B f_Y(y)dy \\ &= \int_A \int_B f_X(x)f_Y(y)dydx \end{aligned}$$

If we **define** $g(x, y) = f_X(x)f_Y(y)$ then we have proved that for $C = A \times B$

$$P((X, Y) \in C) = \int_C g(x, y)dydx$$

To prove that g is the joint density of (X, Y) we need only prove that this integral formula is valid for an arbitrary Borel set C , not just a rectangle $A \times B$. This is proved via a *monotone class* argument. You prove that the collection of sets C for which the identity holds has closure properties which guarantee that this collection includes the Borel sets.

A **Monotone Class** is a collection \mathcal{C} of subsets of a given set which is closed under the operations of increasing countable unions and decreasing countable intersections. That is it has the properties:

- (a) If $A_1 \subset A_2 \subset \dots$ are in \mathcal{C} then $\cup_1^\infty A_i \in \mathcal{C}$.
- (b) If $A_1 \supset A_2 \supset \dots$ are in \mathcal{C} then $\cap_1^\infty A_i \in \mathcal{C}$.

Lemma 1 *The smallest monotone class containing a field \mathcal{F}_γ is a σ -field.*

A field is defined just the same way a σ -field is but with only finite unions and intersections. To prove the Lemma you let \mathcal{C} be the smallest monotone class containing \mathcal{F}_γ ; this is simply the intersection of all monotone classes containing \mathcal{F}_γ . Notice the similarity to the definition of the Borel σ -field.

Whenever you want to prove a fact about an object defined in this way you will need an indirect proof. Consider first the collection \mathcal{M} of all those sets $A \in \mathcal{C}$ for which $A^c \in \mathcal{C}$. You should prove that \mathcal{M} is a monotone class. Since \mathcal{F}_γ is a field and a subset of \mathcal{C} you see that \mathcal{M} contains \mathcal{F}_γ . Since \mathcal{C} is the smallest monotone class containing \mathcal{F}_γ and \mathcal{M} is another monotone class containing \mathcal{F}_γ we see that $\mathcal{C} \subset \mathcal{M}$. But this just means that every $A \in \mathcal{C}$ has the property $A^c \in \mathcal{C}$. In other words \mathcal{C} is closed under the operation of taking complements, one of the defining properties of a σ -field.

Next fix a set $A \in \mathcal{F}_\gamma$. Let \mathcal{M} be the collection of all $B \in \mathcal{C}$ for which $A \cup B \in \mathcal{C}$. Again you can check that \mathcal{M} is a monotone class and that \mathcal{M} contains \mathcal{F}_γ . As in the previous part this shows that $\mathcal{C} \subset \mathcal{M}$. In other words for every $A \in \mathcal{F}_\gamma$ and every $B \in \mathcal{C}$ we have $A \cup B \in \mathcal{C}$.

Finally fix a set A in \mathcal{C} . Let \mathcal{M} be the collection of all $B \in \mathcal{C}$ for which $A \cup B \in \mathcal{C}$. Again you can check that \mathcal{M} is a monotone class. The previous step showed that \mathcal{M}

contains \mathcal{F}_o . As in the previous part this shows that $\mathcal{C} \subset \mathcal{M}$. In other words for every $A \in \mathcal{C}$ and every $B \in \mathcal{C}$ we have $A \cup B \in \mathcal{C}$.

These three steps have shown that \mathcal{C} is a field. The fact that \mathcal{C} is a monotone class and a field allows us to rewrite a countable union of sets in \mathcal{C} as an increasing union of sets in \mathcal{C} . This proves \mathcal{C} is a σ -field. •

It remains to apply this lemma to the assertion. Call any set C of the form $A \times B$ where A and B are Borel a rectangle. You have proved that

$$P((X, Y) \in C) = \int_C f_X(x) f_Y(y) dx dy \quad (2.1)$$

for each rectangle C . If C is a finite union of rectangles then we may rewrite C as a finite disjoint union of rectangles to see that (2.1) holds for any such C . Similarly the complement of a rectangle may be rewritten as a finite disjoint union of rectangles. It follows that the collection \mathcal{M} of subsets of \mathbb{R}^{p+q} for which (2.1) holds contains the field formed by the collection of all finite disjoint unions of rectangles. Finally \mathcal{M} is a monotone class by the dominated convergence theorem. Here's precisely how the dominated convergence theorem applies.

The function $g(x, y) = f_X(x) f_Y(y)$ is positive and has integral equal to 1 by Tonelli's theorem. If $C_1 \subset C_2 \subset \dots$ then put

$$g_n(x, y) = g(x, y) 1(x \in C_n)$$

Note that $|g_n(x, y)| \leq |g(x, y)|$ and that g_n converges for almost all x to

$$g_\infty(x, y) \equiv g(x, y) 1(x \in \cup C_n)$$

Thus

$$\begin{aligned} \lim_{n \rightarrow \infty} \int g_n(x, y) dx dy &= \lim_{n \rightarrow \infty} \int_{C_n} f_X(x) f_Y(y) dx dy \\ &= \int g_\infty(x, y) dx dy \\ &= \int_{\cup C_n} f_X(x) f_Y(y) dx dy \end{aligned}$$

At the same time

$$\lim P((X, Y) \in C_n) = P((X, Y) \in \cup C_n)$$

so that

$$P((X, Y) \in \cup C_n) = \int_{\cup C_n} f_X(x) f_Y(y) dx dy$$

This and a corresponding argument for intersections show that \mathcal{M} is a monotone class containing the field of all finite unions of Borel rectangles. It thus contains the smallest σ field containing this field. But every open rectangle is in this field so in \mathcal{M} . Every open set is a countable union of open rectangles and so in \mathcal{M} . Thus \mathcal{M} is a σ field containing the open sets and so must contain the Borel σ field.

4. This is proved via another monotone class argument.
5.

$$\begin{aligned} P(X \in A, Y \in B) &= \int_A \int_B g(x)h(y)dydx \\ &= \int_A g(x)dx \int_B h(y)dy \end{aligned}$$

Take $B = \mathbb{R}^1$ to see that

$$P(X \in A) = c_1 \int_A g(x)dx$$

where $c_1 = \int h(y)dy$. From the definition of density we see that c_1g is the density of X . Since $\iint f_{X,Y}(xy)dxdy = 1$ we see that $\int g(x)dx \int h(y)dy = 1$ so that $c_1 = 1/\int g(x)dx$. A similar argument works for Y .

Theorem 2 *If X_1, \dots, X_p are independent and $Y_i = g_i(X_i)$ then Y_1, \dots, Y_p are independent. Moreover, (X_1, \dots, X_q) and (X_{q+1}, \dots, X_p) are independent.*

Conditional probability

Definition: $P(A|B) = P(AB)/P(B)$ provided $P(B) \neq 0$.

Definition: For discrete random variables X and Y the conditional probability mass function of Y given X is

$$\begin{aligned} f_{Y|X}(y|x) &= P(Y = y|X = x) \\ &= f_{X,Y}(x, y)/f_X(x) \\ &= f_{X,Y}(x, y)/\sum_t f_{X,Y}(x, t) \end{aligned}$$

For absolutely continuous X the problem is that $P(X = x) = 0$ for all x so how can we define $P(A|X = x)$ or $f_{Y|X}(y|x)$? The solution is to take a limit

$$P(A|X = x) = \lim_{\delta x \rightarrow 0} P(A|x \leq X \leq x + \delta x)$$

If, for instance, X, Y have joint density $f_{X,Y}$ then with $A = \{Y \leq y\}$ we have

$$\begin{aligned} P(A|x \leq X \leq x + \delta x) &= \frac{P(A \cap \{x \leq X \leq x + \delta x\})}{P(x \leq X \leq x + \delta x)} \\ &= \frac{\int_{-\infty}^y \int_x^{x+\delta x} f_{X,Y}(u, v)dudv}{\int_x^{x+\delta x} f_X(u)du} \end{aligned}$$

Divide the top and bottom by δx and let δx tend to 0. The denominator converges to $f_X(x)$ while the numerator converges to

$$\int_{-\infty}^y f_{X,Y}(x, v)dv$$

So we define the conditional CDF of Y given $X = x$ to be

$$P(Y \leq y | X = x) = \frac{\int_{-\infty}^y f_{X,Y}(x, v) dv}{f_X(x)}$$

Differentiate with respect to y to get the definition of the conditional density of Y given $X = x$ namely

$$f_{Y|X}(y|x) = f_{X,Y}(x, y) / f_X(x)$$

or in words “conditional = joint/marginal”.

Chapter 3

Expectation

In undergraduate courses we give two definitions of expected value, one for discrete random variables and one for absolutely continuous random variables. In general, there are random variables which are neither absolutely continuous nor discrete. Here's how probabilists define E in general.

Definition: A random variable X is simple if we can write

$$X(\omega) = \sum_1^n a_i 1(\omega \in A_i)$$

for some constants a_1, \dots, a_n and events A_i .

Definition: For a simple rv X we define

$$E(X) = \sum a_i P(A_i)$$

Remark: There are many ways to write a simple random variable as a sum $\sum a_i 1(\omega \in A_i)$. You should prove that if

$$P \left[\sum a_i 1(\omega \in A_i) = \sum b_j 1(\omega \in B_j) \right] = 1$$

then

$$\sum a_i P(A_i) = \sum b_j P(B_j)$$

so that our definition is really a definition.

For positive random variables which are not simple we extend our definition by approximation:

Definition: If $X \geq 0$ then

$$E(X) = \sup\{E(Y) : 0 \leq Y \leq X, Y \text{ simple}\}$$

Remark: For $X \geq 0$ we are permitting $+\infty$ as a value of E .

Remark: This definition *redefines* E for positive simple X . You must check that if X is simple, say

$$X = \sum a_i 1_{A_i}$$

and $X \geq 0$ then

$$\sup\{E(Y) : 0 \leq Y \leq X, Y \text{ simple}\} = \sum a_i P(A_i)$$

so that the definition has not been changed.

Definition: We call X **integrable** if

$$E(|X|) < \infty.$$

In this case we define

$$E(X) = E(\max(X, 0)) - E(\max(-X, 0))$$

Remark: Again you must check that you have not changed the definition of $E(X)$ for either of the previous categories of X .

Facts: E is a linear, monotone, positive operator:

1. **Linear:** $E(aX + bY) = aE(X) + bE(Y)$ provided X and Y are integrable.
2. **Positive:** $P(X \geq 0) = 1$ implies $E(X) \geq 0$.
3. **Monotone:** $P(X \geq Y) = 1$ and X, Y integrable implies $E(X) \geq E(Y)$.

Each of these facts is proved first for simple functions then for positive functions then for general integrable functions. You must follow the steps in the definition of E in any proof of general properties of E .

Major technical theorems:

Theorem 3 (Monotone Convergence) *If $0 \leq X_1 \leq X_2 \leq \dots$ (almost surely, that is, the probability of this event is 1) and $X = \lim X_n$ (which has to exist almost surely) then*

$$E(X) = \lim_{n \rightarrow \infty} E(X_n)$$

Theorem 4 (Dominated Convergence) *If $|X_n| \leq Y_n$ a.s. and there is a random variable X such that $X_n \rightarrow X$ almost surely and a random variable Y such that $Y_n \rightarrow Y$ almost surely with $E(Y_n) \rightarrow E(Y) < \infty$ then*

$$E(X_n) \rightarrow E(X)$$

This is often used with all Y_n the same random variable Y .

Theorem 5 (Fatou's Lemma) *If $X_n \geq 0$ a.s. then*

$$E(\limsup X_n) \leq \limsup E(X_n)$$

Theorem 6 *With this definition of E if X has density $f(x)$ (even in \mathbb{R}^p say) and $Y = g(X)$ then*

$$E(Y) = \int g(x)f(x)dx.$$

(This could be a multiple integral.) If X has pmf f then

$$E(Y) = \sum_x g(x)f(x).$$

This works for instance even if X has a density but Y doesn't.

Definition: The r^{th} moment (about the origin) of a real random variable X is $\mu'_r = E(X^r)$ (provided it exists). We generally use μ for $E(X)$. The r^{th} central moment is

$$\mu_r = E[(X - \mu)^r]$$

We call $\sigma^2 = \mu_2$ the variance.

Definition: For an \mathbb{R}^p valued random vector X we define $\mu_X = E(X)$ to be the vector whose i^{th} entry is $E(X_i)$ (provided all entries exist).

Definition: The $(p \times p)$ variance covariance matrix of X is

$$\text{Var}(X) = E[(X - \mu)(X - \mu)^t]$$

which exists provided each component X_i has a finite second moment. More generally if $X \in \mathbb{R}^p$ and $Y \in \mathbb{R}^q$ both have all components with finite second moments then

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)^T]$$

We have

$$\text{Cov}(AX + b, CY + d) = A\text{Cov}(X, Y)B^T$$

for general (conforming) matrices A, C and vectors b and d .

Moments and probabilities of rare events are closely connected as will be seen in a number of important probability theorems. Here is one version of Markov's inequality (one case is Chebyshev's inequality):

$$\begin{aligned} P(|X - \mu| \geq t) &= E[1(|X - \mu| \geq t)] \\ &\leq E\left[\frac{|X - \mu|^r}{t^r} 1(|X - \mu| \geq t)\right] \\ &\leq \frac{E[|X - \mu|^r]}{t^r} \end{aligned}$$

The intuition is that if moments are small then large deviations from average are unlikely.

3.1 Moments and independence

Theorem 7 *If X_1, \dots, X_p are independent and each X_i is integrable then $X = X_1 \cdots X_p$ is integrable and*

$$E(X_1 \cdots X_p) = E(X_1) \cdots E(X_p)$$

Proof: Suppose each X_i is simple: $X_i = \sum_j x_{ij} 1(X_i = x_{ij})$ where the x_{ij} are the possible

values of X_i . Then

$$\begin{aligned}
 E(X_1 \cdots X_p) &= \sum_{j_1 \cdots j_p} x_{1j_1} \cdots x_{pj_p} E(1(X_1 = x_{1j_1}) \cdots 1(X_p = x_{pj_p})) \\
 &= \sum_{j_1 \cdots j_p} x_{1j_1} \cdots x_{pj_p} P(X_1 = x_{1j_1} \cdots X_p = x_{pj_p}) \\
 &= \sum_{j_1 \cdots j_p} x_{1j_1} \cdots x_{pj_p} P(X_1 = x_{1j_1}) \cdots P(X_p = x_{pj_p}) \\
 &= \left[\sum_{j_1} x_{1j_1} P(X_1 = x_{1j_1}) \right] \cdots \left[\sum_{j_p} x_{pj_p} P(X_p = x_{pj_p}) \right] \\
 &= \prod E(X_i)
 \end{aligned}$$

For general $X_i \geq 0$ we create a sequence of simple approximations by rounding X_i down to the nearest multiple of 2^{-n} (to a maximum of n). Call this random variable $X_{i,n}$. Then each $X_{i,n}$ is simple and the variables $X_{1,n}, \dots, X_{p,n}$ are independent. Thus

$$E\left(\prod X_{j,n}\right) = \prod E(X_{j,n})$$

for each n . If

$$X_n^* = \prod X_{j,n}$$

then

$$0 \leq X_1^* \leq X_2^* \leq \cdots$$

and X_n^* converges to $X^* = \prod X_i$ so that

$$E(X^*) = \lim E(X_n^*)$$

by the monotone convergence theorem. Similarly by the monotone convergence theorem

$$\lim \prod E(X_{j,n}) = \prod E(X_j) < \infty$$

This shows both that X^* is integrable and that

$$E\left(\prod X_j\right) = \prod E(X_j)$$

The general case uses the fact that we can write each X_i as the difference of its positive and negative parts:

$$X_i = \max(X_i, 0) - \max(-X_i, 0)$$

Just expand out the product and use the previous case.

3.2 Appendix on Lebesgue Integration

Corresponding to Lebesgue measure is the Lebesgue integral which is defined in much the same way as E.

We call a Borel measurable function f simple if

$$f(x) = \sum_1^n a_i 1(x \in B_i)$$

for almost all $x \in \mathbb{R}^p$ and some constants a_i and Borel sets A_i with $\lambda(B_i) < \infty$. For such an f we define

$$\int f(x) dx = \sum a_i \lambda(A_i)$$

Again if

$$\sum a_i 1_{A_i} = \sum b_j 1_{B_j}$$

almost everywhere and all A_i and B_j have finite Lebesgue measure you must check that

$$\sum \lambda(A_i) = \sum b_j \lambda(B_j)$$

If $f \geq 0$ almost everywhere and f is Borel define

$$\int f(x) dx = \sup \left\{ \int g(y) dy \right\}$$

where the sup ranges over all simple functions g such that $0 \leq g(x) \leq f(x)$ for almost all x . Call $f \geq 0$ integrable if $\int f(x) dx < \infty$.

Call a general f integrable if $|f|$ is integrable and define for integrable f

$$\int f(x) dx = \int \max(f(x), 0) dx - \int \max(-f(x), 0) dx$$

Remark: Again you must check that you have not changed the definition of f for either of the previous categories of f .

Facts: \int is a linear, monotone, positive operator:

1. **Linear:** $\int a f(x) + b g(x) dx = a \int f(x) dx + b \int g(x) dx$ provided f and g are integrable.
2. **Positive:** If $f(x) \geq 0$ almost everywhere then $\int f(x) dx \geq 0$.
3. **Monotone:** If $f(x) > g(x)$ almost everywhere and f and g are integrable then $\int f(x) dx \geq \int g(x) dx$.

Each of these facts is proved first for simple functions then for positive functions then for general integrable functions.

Major technical theorems:

Theorem 8 (Monotone Convergence) *If $0 \leq f_1 \leq f_2 \leq \dots$ almost everywhere and $f = \lim f_n$ (which has to exist almost everywhere) then*

$$\int f(x)dx = \lim_{n \rightarrow \infty} \int f_n(x)dx$$

Theorem 9 (Dominated Convergence) *If $|f_n| \leq g_n$ and there is a Borel function f such that $f_n(x) \rightarrow f(x)$ for almost all x and a Borel function g such that $g_n(x) \rightarrow g(x)$ with $\int g_n(x)dx \rightarrow \int g(x)dx < \infty$ then f is integrable and*

$$\int f_n(x)dx \rightarrow \int f(x)dx$$

Theorem 10 (Fatou's Lemma) *If $f_n \geq 0$ almost everywhere then*

$$\int \limsup f_n(x)dx \leq \limsup \int f_n(x)dx$$

Notice the frequent of almost all or almost everywhere in the hypotheses. In our definition of E wherever we require a property of the function $X(\omega)$ we can require it to hold only for a set of ω whose complement has probability 0. In this case we say the property holds **almost surely**. For instance the dominated convergence theorem is usually written:

Dominated Convergence: If $|X_n| \leq Y_n$ almost surely (often abbreviated to a.s.) and there is a random variable X such that $X_n \rightarrow X$ a.s. and a random variable Y such that $Y_n \rightarrow Y$ almost surely with $E(Y_n) \rightarrow E(Y) < \infty$ then

$$E(X_n) \rightarrow E(X)$$

The hypothesis of almost sure convergence can be weakened; I hope to discuss this later in the course.

Multiple Integration

The previous appendix gave a definition of integrals over R^p in terms of Lebesgue measure on R^p . Students will be used to doing integrals one variable at a time. The fact that the multiple integral is equal to the iterated integral is contained in two theorems.

Theorem 11 (Tonelli) : *If $f : \mathbb{R}^{p+q} \mapsto \mathbb{R}$ is Borel and $f \geq 0$ almost everywhere then for almost every $x \in \mathbb{R}^p$ the integral*

$$g(x) \equiv \int f(x, y)dy$$

exists and

$$\int g(x)dx = \int f(x, y)dxdy$$

The right hand side of this formula denotes the $p + q$ dimensional integral defined in the previous appendix.

Theorem 12 (Fubini) *If $f : \mathbb{R}^{p+q} \mapsto \mathbb{R}$ is Borel and integrable then for almost every $x \in \mathbb{R}^p$ the integral*

$$g(x) \equiv \int f(x, y) dy$$

exists and is finite. Moreover g is integrable and

$$\int g(x) dx = \int f(x, y) dx dy .$$

Chapter 4

The Strong Law of Large Numbers

In this chapter I want to prove the strong law of large numbers. Along the way we will examine some standard approaches to probability theorems:

- Many events have sophisticated mathematical definitions (such as the event that a sequence has a limit, or the event that infinitely many times in an infinite sequence of coin tosses there will be more heads than tails). We will see how to convert the formal mathematical definitions into set notation.
- The resulting set theoretic expression may not be convenient. We use theorems of real analysis to re-express the events in more convenient ways.
- In order to prove $P(A) = 0$ it is sometimes convenient to find (guess or whatever) an event B which contains A and has $P(B) = 0$.
- When dealing with an event A defined in terms of infinitely many random variables it is helpful to write A as an increasing union or decreasing intersection of events each defined from only finitely many random variables.
- In the end inequalities play a central role. We will establish Kolmogorov's inequality – an extension of Chebyshev's inequality from one random variable to the maximum of several.

4.1 Events in Set Notation

The event that a sequence of random variables Y_n converges to 0 is

$$A \equiv \{\omega : \lim_{n \rightarrow \infty} Y_n(\omega) = 0\}$$

This event is not yet explicitly written in terms of simple events involving only a finite number of Y s at a time. To make the conversion we must recall the basic definition of limit: A sequence y_n converges to 0 if for every $\epsilon > 0$ there exists an N such that for every $n \geq N$ we have $|y_n| \leq \epsilon$.

Now to convert the definition into set theory notation we replace y_n by $Y_n(\omega)$, each *for every* by an intersection and each *there exists* with a union. We get

$$A = \bigcap_{\epsilon > 0} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{\omega : |Y_n(\omega)| \leq \epsilon\}$$

Is this an event? As it sits it is not obvious because the intersection over $\epsilon > 0$ is an uncountable intersection. In fact, however, the intersection is countable. Let

$$B_\epsilon \equiv \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{\omega : |Y_n(\omega)| \leq \epsilon\}$$

Notice that if $\epsilon' < \epsilon$ then $B_{\epsilon'} \subset B_\epsilon$. This means that

$$\bigcap_{\epsilon > 0} B_\epsilon = \bigcap_{m=1}^{\infty} B_{1/m}$$

This shows that A is a countable intersection of countable unions of countable intersections of events and so A is an event.

Here are some other events:

- The sequence S_n has a limit. Formally, a sequence s_n has a limit if there exists s_∞ such that for every $\epsilon > 0$ there exists an N such that for every $n \geq N$ we have $|s_n - s_\infty| \leq \epsilon$. Mechanically this leads to the event

$$\bigcup_s \bigcap_{\epsilon > 0} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{\omega : |S_n(\omega) - s| \leq \epsilon\}$$

Again the union over ϵ can be made countable. The union over s , however, is not easy to make countable. Instead we use a theorem of analysis to describe the existence of a limit in a different way. It is a theorem that a sequence s_n has a limit if and only if the sequence is Cauchy. Here is the definition of Cauchy: for every $\epsilon > 0$ there exists an N such that for every $n \geq N$ we have $|s_n - s_N| \leq \epsilon$. Our event that S_n has a limit is thus

$$\bigcap_{\epsilon > 0} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq \epsilon\}$$

Again the intersection over all $\epsilon > 0$ can be replaced with a countable intersection over $\epsilon = 1/r$ for positive integers r .

- The sequence Y_n is summable. This means that the series of partial sums $S_n = \sum_1^n Y_i$ has a limit so our event is

$$\bigcap_{r=1}^{\infty} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \bigcap_{m=n}^{\infty} \left\{ \left| \sum_{j=n+1}^m Y_j \right| \leq 1/r \right\}$$

- The sequence S_n is positive for infinitely many n . In other words for every N there is an $n \geq N$ such that $S_n > 0$. The event is

$$\bigcap_{N=1}^{\infty} \bigcup_{n=N}^{\infty} \{S_n > 0\}$$

- The limit superior of S_n is 1. It is easiest to write this as the intersection of two events, namely $\limsup S_n \leq 1$ and $\limsup S_n \geq 1$. The former means that for all $\epsilon > 0$ there is an N such that for all $n \geq N$ $S_n \leq 1 + \epsilon$. The latter means that for all $\epsilon > 0$ and all N there is an $n \geq N$ such that $S_n \geq 1 - \epsilon$. So our event is $A^* \cap A_*$ where

$$A^* = \bigcap_{r=1}^{\infty} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{S_n \leq 1 + 1/r\}$$

and

$$A_* = \bigcap_{r=1}^{\infty} \bigcap_{N=1}^{\infty} \bigcup_{n=N}^{\infty} \{S_n \geq 1 + 1/r\}$$

4.2 The Strong Law of Large Numbers with 4 moments

Theorem 13 *Suppose that $X_n; n \geq 1$ is a sequence of independent identically distributed random variables with $E(X_n) = 0$. Then $n^{-1} \sum_1^n X_k \rightarrow 0$ almost surely.*

We have seen in the previous section that we are supposed to prove that $P(A) = 1$ where

$$A = \bigcap_{r=1}^{\infty} \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} B_{n,r}$$

where $B_{n,r}$ is the event

$$\{|(X_1 + \cdots + X_n)/n| \leq 1/r\}$$

Let

$$C_r = \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} B_{n,r}$$

and notice that

$$C_1 \supset C_2 \supset \cdots$$

so that

$$P(C_1) \geq P(C_2) \geq \cdots \geq \lim P(C_r) = P(A)$$

It is thus necessary and sufficient to prove that $P(C_r) = 1$ for each r .

Remark: This process of recognizing decreasing or increasing sequences of sets plays a substantial role in many probability arguments. With several intersections and unions in a row it can be hard to figure out whether or not the objects really do increase or decrease. However many definitions which begin with “for every $\epsilon > 0$ ” will be decreasing intersections

because the condition in the definition is harder to satisfy for smaller ϵ . Similarly unions in definitions tend to be increasing.

Now each C_r is an increasing union. Let

$$C_{N,r} = \bigcap_{n=N}^{\infty} B_{n,r}$$

and notice that

$$C_{1,r} \subset C_{2,r} \subset \dots$$

We need to prove that

$$P(C_r) = \lim_N P(C_{N,r}) = 1$$

or equivalently that

$$\lim_N P(C_{N,r}^c) = 0$$

I will now show you one version of the strong law of large numbers which requires more moment conditions than the real thing but has a relatively direct proof. The proof uses a standard strategy worth knowing. It is based on the **Borel-Cantelli Lemma**.

Suppose A_n is a sequence of events. The event

$$B \equiv \{A_n \text{ infinitely often}\} = \{A_n \text{ i.o.}\}$$

is the set of ω belonging to infinitely many of the A_n . We have seen above that

$$B = \bigcap_{N=1}^{\infty} \bigcup_{n=N}^{\infty} A_n$$

Lemma 2 (Borel-Cantelli) *If $\sum P(A_n) < \infty$ then $P(A_n \text{ i.o.}) = 0$.*

To prove the lemma we need to prove that

$$\lim_{N \rightarrow \infty} P\left(\bigcup_{n=N}^{\infty} A_n\right) = 0$$

But (this next inequality is quite widely useful – it applies to an arbitrary sequence of events A_n)

$$P\left(\bigcup_{n=N}^{\infty} A_n\right) \leq \sum_{n=N}^{\infty} P(A_n)$$

It is a fact from real analysis that a convergent series has tail sums which converge to 0 so

$$\lim_{N \rightarrow \infty} \sum_{n=N}^{\infty} P(A_n) = 0$$

which proves the lemma.

Now to apply the lemma notice that the complement of

$$C_r = \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} B_{n,r}$$

is just

$$C_r^c = \bigcap_{N=1}^{\infty} \bigcup_{n=N}^{\infty} B_{n,r}^c$$

In words the opposite of the event that there is an N such that $B_{n,r}$ happens for every $n \geq N$ is that for every N there is at least one $n \geq N$ for which $B_{n,r}$ does not happen. That is

$$C_r^c = \{B_{n,r}^c \text{ i.o.}\}$$

Summary: if we prove

$$\sum_{n=1}^{\infty} P(B_{n,r}^c) < \infty$$

then we will conclude $P(C_r^c) = 0$ so $P(C_r) = 1$.

Eventually in most probability theorems you reach this stage; we need an upper bound to help us compute the sum. I will use one of the versions of Markov's inequality.

$$\begin{aligned} P(B_{n,r}^c) &= P(|X_1 + \cdots + X_n| > n/r) \\ &\leq \frac{E[(X_1 + \cdots + X_n)^4]}{(n/r)^4} \end{aligned}$$

The fourth moment is (after expanding everything out and collecting all the different kinds of terms

$$E[(X_1 + \cdots + X_n)^4] = nE(X_1^4) + n(n-1)[E(X_1^2)]^2$$

This produces the bound (using $n(n-1) \leq n^2$)

$$P(B_{n,r}^c) \leq \frac{r^4 E(X_1^4)}{n^3} + \frac{r^4 [E(X_1^2)]^2}{n^2}$$

Since

$$\sum_n n^{-p} < \infty$$

for any $p > 1$ the infinite sums converge proving the Strong Law of Large Numbers for random variables for which

$$E(X_1^4) < \infty \quad \text{and} \quad E(X_1^2) < \infty$$

In fact the inequality $X_1^2 \leq 1 + X_1^4$ proves that $E(X_1^4) < \infty$ implies $E(X_1^2) < \infty$. There are many inequalities concerning moments which are worth learning to use: Cauchy Schwarz, Minkowski, Holder, Jensen are the names of some of the most famous.

4.3 Proof without 4 finite moments

In probability (and mathematical statistics) considerable effort is devoted to eliminating the need for moment assumptions. The Strong Law of Large Numbers is an example. I am going to weaken the assumption of 4 finite moments to that of 1 finite moment. I will do this in two stages to illustrate the techniques. First I will suppose that each X_i has a finite variance (which is the same as $E(X_1^2)$ because the random variables have mean 0).

The first standard tactic of probabilists is to use some real variables result to replace the event whose probability you are computing by another whose probability is the same (or almost the same and then take limits). In this case we use Kronecker's Lemma.

Lemma 3 (Kronecker) *Suppose*

$$0 < b_1 < b_2 \cdots$$

and that x_n is a sequence of real numbers. If

$$s_n = \sum_{i=1}^n \frac{x_i}{b_i}$$

is a convergent sequence then

$$\lim_{n \rightarrow \infty} \frac{x_1 + \cdots + x_n}{b_n} = 0$$

I won't prove the lemma but observe this. The lemma implies that if

$$A^* = \left\{ \sum_1^n X_k/k \text{ is a convergent sequence} \right\}$$

then

$$A^* \subset A$$

so that we need only prove that $P(A^*) = 1$.

We have already seen that the event that a sequence is convergent is the same as the event that it is Cauchy so we must compute the probability of the event

$$\bigcap_r \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq 1/r\}$$

where $S_n = \sum 1^n X_i/i$. The key ingredient in the calculation is a lower bound on

$$P \left(\bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq \epsilon\} \right)$$

This is the same problem as finding an upper bound on the probability of the complement which is to bound

$$P \left(\bigcup_{n=N}^{\infty} \{|S_n - S_N| > \epsilon\} \right) = \lim_{m \rightarrow \infty} P \left(\bigcup_{n=N}^m \{|S_n - S_N| > \epsilon\} \right)$$

Now $S_n - S_N$ is a sum of $n - N$ independent random variables. Our general need is for a bound on

$$P\left(\bigcup_1^n \left\{ \left| \sum_1^k Z_i \right| > \epsilon \right\}\right)$$

where the Z_i are independent mean 0 random variables. We use Kolmogorov's inequality which is a strengthening of Chebyshev's:

Theorem 14 (Kolmogorov) *Suppose that Z_1, \dots, Z_n are independent random variables with finite variances and mean 0. Put $S_k = \sum_1^k Z_i$. Then*

$$P(\exists k, 1 \leq k \leq n : |S_k| > \epsilon) \leq \frac{\text{Var}(S_n)}{\epsilon^2}$$

The proof uses the idea of a stopping time. Define T to be the least value of $k \in \{1, \dots, n\}$ for which $|S_k| > \epsilon$ if such a k exists and $T = n$ if no such k exists. Notice that the event $T = k$ is determined by the variables Z_1, \dots, Z_k ; that is, there is a function $f_k(Z_1, \dots, Z_k)$ such that

$$1(T = k) = f_k(Z_1, \dots, Z_k)$$

Now

$$\begin{aligned} P(\exists k, 1 \leq k \leq n : |S_k| > \epsilon) &= P(|S_T| > \epsilon) \\ &\leq \frac{E(S_T^2)}{\epsilon^2} \end{aligned}$$

by Chebyshev's inequality.

Next we compare $E(S_n^2)$ to $E(S_T^2)$. We have

$$\begin{aligned} E(S_n^2) &= E[(S_n - S_T + S_T)^2] \\ &= E(S_T^2) + E[(S_n - S_T)^2] + 2E[S_T(S_n - S_T)] \end{aligned}$$

The first two terms are non-negative so if

$$E[S_T(S_n - S_T)] = 0$$

the theorem will be proved.

But

$$\begin{aligned} S_T(S_n - S_T) &= \sum_1^{n-1} S_k(S_n - S_k)1(T = k) \\ &= \sum_1^{n-1} S_k(S_n - S_k)f_k(Z_1, \dots, Z_k) \end{aligned}$$

Notice that $S_n - S_k$ and $S_k f_k(Z_1, \dots, Z_k)$ are independent so

$$\begin{aligned} E[S_T(S_n - S_T)] &= \sum_1^{n-1} E(S_n - S_k)E[S_k 1(T = k)] \\ &= 0 \end{aligned}$$

Remark: In fact all I needed to prove Kolmogorov's inequality was to know that $S_n - S_k$ was uncorrelated with any function of S_1, \dots, S_k (or equivalently uncorrelated with any function of Z_1, \dots, Z_k).

Definition: A sequence S_1, S_2, \dots is a **martingale** if, for each $k \geq 1$,

$$E[S_{k+1} | S_1, \dots, S_k] = S_k$$

(I have yet to define conditional expectation properly but another definition is that

$$E[(S_{k+1} - S_k)f_k(S_1, \dots, S_k)] = 0$$

for every bounded Borel function f_k .)

Theorem 15 (Kolmogorov) *Suppose that S_1, \dots, S_n is a martingale with finite variances and mean 0. Then*

$$P(\exists k, 1 \leq k \leq n : |S_k| > \epsilon) \leq \frac{\text{Var}(S_n)}{\epsilon^2}.$$

Now back to the strong law. Fix an $\epsilon > 0$ and an integer N . The events

$$B_{m,N} = \bigcap_{n=N}^m \{|S_n - S_N| \leq \epsilon\}$$

decrease, that is, $B_{1,N} \supset B_{2,N} \supset \dots$. So

$$P(B_{\infty,N}) \equiv P\left(\bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq \epsilon\}\right) = \lim_m P(B_{m,N})$$

Now

$$\begin{aligned} P(B_{m,N}) &= 1 - P(B_{m,N}^c) \\ &\geq 1 - \frac{\text{Var}(S_m - S_N)}{\epsilon^2} \\ &= 1 - \frac{\text{Var}(X_1) \sum_{N+1}^m i^{-2}}{\epsilon^2} \\ &\geq 1 - \frac{\text{Var}(X_1)}{\epsilon^2} \sum_{N+1}^{\infty} \frac{1}{i^2} \end{aligned}$$

This shows

$$P(B_{N,\infty}) \geq 1 - \frac{\text{Var}(X_1)}{\epsilon^2} \sum_{N+1}^{\infty} \frac{1}{i^2}$$

Since $\sum_1^{\infty} i^{-2} < \infty$ the right hand side of this inequality goes to 0 as N goes to ∞ . In other words

$$\lim_{N \rightarrow \infty} P\left(\bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq \epsilon\}\right) = 1$$

This proves

$$P\left(\bigcup_{m=1}^{\infty} \bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq \epsilon\}\right) = 1$$

because the union over m is increasing.

Finally we have shown for each r that $P(C_r) = 1$ where

$$C_r = \bigcup_{m=1}^{\infty} \bigcap_{n=N}^{\infty} \{|S_n - S_N| \leq 1/r\}$$

Hence

$$P\left(\bigcap_r C_r\right) = 1$$

This establishes the strong law of large numbers for variables with a finite variance.

The Strong Law with 1 Moment

Now suppose that we can't assume that $\text{Var}(X_1) < \infty$. In the arguments of the previous section we used this in Kolmogorov's inequality to try to verify that

$$P\left(\sum_1^n X_k/k \text{ converges}\right) = 1$$

The widely used probability trick here is called **truncation**. Let

$$X_k^* = X_k 1(|X_k| \leq k)$$

For most values of k , X_k and X_k^* are equal. Note that

$$\begin{aligned} P(X_k \neq X_k^*) &= P(|X_k| > k) \\ &= P(|X_1| > k) \end{aligned}$$

Note that

$$\begin{aligned} \sum_{k=1}^{\infty} P(|X| > k) &= \sum_{k=1}^{\infty} \sum_{j=k}^{\infty} P(j < |X| \leq j+1) \\ &= \sum_{j=1}^{\infty} \sum_{k=1}^j P(j < |X| \leq j+1) \\ &= \sum_{j=1}^{\infty} j P(j < |X| \leq j+1) \end{aligned}$$

Notice that

$$\begin{aligned} E\{|X| 1(j < |X| \leq j+1)\} &\geq j E\{1(j < |X| \leq j+1)\} \\ &= j P(j < |X| \leq j+1) \end{aligned}$$

Hence

$$\begin{aligned} \sum_{k=1}^{\infty} P(|X| > k) &\leq \sum_{j=1}^{\infty} E\{|X|1(j < |X| \leq j+1)\} \\ &= E\{|X|1(1 < |X|)\} \\ &\leq E(|X|) \\ &< \infty. \end{aligned}$$

Now apply the Borel Cantelli Lemma to prove that

$$P(X_k \neq X_k^* \text{ i.o.}) = 0$$

In other words except for a set N of ω with $P(N) = 0$ $X_k(\omega) = X_k^*(\omega)$ except for a finite number of k . If ω is in

$$A^* = \left\{ \sum_1^n X_k^*/n \rightarrow 0 \right\}$$

and not in N then ω is in

$$A = \left\{ \sum_1^n X_k/n \rightarrow 0 \right\}$$

Hence if $P(A^*) = 1$ then $P(A) = 1$.

We intend to mimic the proof given before but the following differences arise:

- Let

$$\mu_k = E(X_k^*)$$

Although $E(X_k) = 0$ the truncation may make the mean differ from 0. So Kolmogorov's inequality cannot be applied directly to the X_k^* variables. Instead we have to subtract the mean.

- The variable X_k^* is bounded by k so it does have a finite variance. But that variance depends on k so the sum we got above, of $\sum 1/i^2$ will be more complicated.

To deal with the first point we begin with an observation. If Y_n are any random variables and a_n a sequence of constants converging to 0 then the event $\{Y_n \rightarrow 0\}$ is almost surely the same as $\{Y_n + a_n \rightarrow 0\}$. For clarity the jargon is that two events are almost surely equal if their indicators are almost surely equal. Another way of saying this is

$$P(\{Y_n \rightarrow 0\} \Delta \{Y_n + a_n \rightarrow 0\}) = 0$$

where

$$A \Delta B = (AB^c) \cup (A^cB) = (A \cup B) \setminus (A \cap B)$$

is the symmetric difference of two sets.

So we will prove:

Claim 1: $\sum_1^n \mu_k/n$ converges to 0

and

Claim 2:

$$P \left\{ \sum_1^n (X_k^* - \mu_k)/n \rightarrow 0 \right\} = 1$$

For Claim 1 note

$$\begin{aligned} \mu_k &= \mathbf{E}(X_k^*) \\ &= \mathbf{E}(X_k 1(|X_k| \leq k)) \\ &= \mathbf{E}[X_1 1(|X_1| \leq k)] \end{aligned}$$

By the dominated convergence theorem

$$\mu_k \rightarrow \mathbf{E}(X_1) = 0.$$

Finally if $a_n \rightarrow 0$ then $(a_1 + \cdots + a_n)/n \rightarrow 0$. Hence

$$\sum_1^n \mu_k \rightarrow 0$$

which is Claim 1.

Turning to Claim 2 we define

$$A_3^* = \left\{ \sum_1^n (X_k^* - \mu_k)/k \text{ converges} \right\}$$

According to Kronecker's Lemma Claim 2 (and so the strong law) will follow if we prove $P(A_3^*) = 1$.

Let $\sigma_k^2 = \text{Var}(X_k^*)$ and

$$B_{m,N} = \bigcap_{n=N}^m \{|S_n^* - S_N^*| \leq \epsilon\}$$

where $S_n^* = \sum_1^n (X_k^* - \mu_k)/k$. As in the finite variance case:

$$\begin{aligned} P(B_{m,N}) &= 1 - P(B_{m,N}^c) \\ &\geq 1 - \frac{\text{Var}(S_m - S_N)}{\epsilon^2} \\ &= 1 - \frac{\sum_{N+1}^m \sigma_i^2/i^2}{\epsilon^2} \end{aligned}$$

As before we want to let $m \rightarrow \infty$ but need to know that

$$\sum_1^\infty \sigma_i^2/i^2 < \infty$$

But

$$\begin{aligned}\sigma_k^2 &\leq \mathbb{E}(X_k^2 1(|X_k| \leq k)) \\ &= \mathbb{E}(X_1^2 1(|X_1| \leq k)) \\ &= \sum_{j=1}^k \mathbb{E}(X_1^2 1(j-1 < |X_1| \leq j))\end{aligned}$$

Sum over k and switch the order of summation to get

$$\begin{aligned}\sum_1^n \sigma_k^2/k^2 &\leq \sum_{j=1}^n \sum_{k=j}^n \mathbb{E}(X_1^2 1(j-1 < |X_1| \leq j))/k^2 \\ &\leq \sum_{j=1}^n \mathbb{E}(X_1^2 1(j-1 < |X_1| \leq j)) \sum_{k=j}^{\infty} 1/k^2 \\ &\leq \sum_{j=1}^n \mathbb{E}(X_1^2 1(j-1 < |X_1| \leq j))/j \\ &\leq \sum_{j=1}^n \mathbb{E}(|X_1| 1(j-1 < |X_1| \leq j)) \\ &= \mathbb{E}(|X_1| 1(|X_1| \leq n)) \\ &\leq \mathbb{E}(|X_1|)\end{aligned}$$

This shows

$$P(B_{N,\infty}) \geq 1 - \frac{1}{\epsilon^2} \sum_{N+1}^{\infty} \sigma_k^2/k^2$$

Since $\sum_1^{\infty} \sigma_k^2/k^2 < \infty$ the right hand side of this inequality goes to 0 as N goes to ∞ . In other words

$$\lim_{N \rightarrow \infty} P\left(\bigcap_{n=N}^{\infty} \{|S_n^* - S_N^*| \leq \epsilon\}\right) = 1$$

Thus S_n^* is almost surely Cauchy and $P(A_3^*) = 1$.

4.4 Consistency of MLE

Suppose that X_1, X_2, \dots are independent and identically distributed with density $f(x, \theta_o)$ where

$$\{f(\cdot, \theta); \theta \in \Theta \subset \mathbb{R}\}$$

is a family of densities. In this section we investigate conditions under which the MLE of θ is almost surely consistent. Our focus is on techniques of making the assertion precise in probability theory terms.

Our goal is to find conditions under which we can prove

$$P(\hat{\theta}_n \rightarrow \theta_o) = 1$$

where $\hat{\theta}_n$ is the mle.

We face the following general technical problems:

- What is the precise definition of $\hat{\theta}_n$?
- Having settled on some definition is the resulting object a random variable?

We will focus on an example. The Cauchy(θ) density is

$$f(x, \theta) = \frac{1}{\pi \{1 + (x - \theta)^2\}}$$

For a sample X_1, \dots, X_n the likelihood is

$$\frac{1}{\pi^n \prod_{i=1}^n \{1 + (X_i - \theta)^2\}}$$

Normally we say that $\hat{\theta}$ is the value of θ which maximizes this function of θ .

To be precise this is supposed to define $\hat{\theta}$ as a function of X_1, \dots, X_n . In other words it is supposed that for each x_1, \dots, x_n there is a unique value $\hat{\theta}(x_1, \dots, x_n)$ which maximizes the likelihood. If this were so we would have a definition of a function from \mathbb{R}^n to \mathbb{R} .

In order to discuss this problem in general it is convenient to define the log-likelihood:

$$\ell(\theta|x_1, \dots, x_n) = -n \log(\pi) - \sum_{i=1}^n \log(1 + (x_i - \theta)^2)$$

It is all very well to say that $\hat{\theta}_n(x_1, \dots, x_n)$ maximizes this function of θ for each fixed value of (x_1, \dots, x_n) . But:

1. Is there, for every (x_1, \dots, x_n) a θ which maximizes ℓ ?
2. If so is the θ unique?
3. If so is $\hat{\theta}_n(x_1, \dots, x_n)$ a Borel function of x_1, \dots, x_n ?

Question 1: For the Cauchy density there is always a maximizer. Fix (x_1, \dots, x_n) . As $\theta \rightarrow \pm\infty$ it is easy to check that

$$\ell(\theta|x_1, \dots, x_n) \rightarrow -\infty$$

There is then a M such that $|\theta| > M$ implies

$$\begin{aligned} \sup\{\ell(\theta|x_1, \dots, x_n); |\theta| > M\} \\ &\leq \ell(0|x_1, \dots, x_n) \\ &\leq \sup\{\ell(\theta|x_1, \dots, x_n); |\theta| \leq M\} \end{aligned}$$

Now the function

$$\theta \mapsto \ell(\theta|x_1, \dots, x_n)$$

is continuous so that it assumes its maximum over $[-M, M]$. This shows the existence of at least one maximizing θ for any set of x values.

Question 2: Consider $n = 2$ and $x_1 = x = -x_2$. Then

$$\ell(\theta|x, -x)$$

is an even function of x . This function has derivative

$$\ell' = \frac{2(\theta - x)}{1 + (x - \theta)^2} + \frac{2(\theta + x)}{1 + (x + \theta)^2}$$

At $\theta = 0$ this evaluates to 0 so that $\theta = 0$ is a critical point of ℓ . The second derivative may be computed and evaluated at 0 to give

$$\ell''(0) = \frac{4(x^2 - 1)}{(1 + x^2)^2}$$

If $|x| < 1$ this is negative so that 0 is a local maximum but if $|x| > 1$ it is a local minimum. In this case, since ℓ is even there must be two maxima on either side of 0. Note that putting the two terms in ℓ' on a common denominator will give a numerator which is a multiple of

$$\theta(\theta^2 - (x^2 - 1))$$

Notice that there are exactly three roots if $x^2 > 1$.

Summary: defining $\hat{\theta}$ to be the maximizer of ℓ does not actually define a function.

Alternative strategies:

1: You might pick one of the maximizing θ values in an unequivocal way:

$$\hat{\theta} = \inf\{\theta : \ell(\theta) = \sup \ell\}$$

(The set of such θ is not empty and bounded so there is such a $\hat{\theta}$ and that $\hat{\theta}$ is finite. By continuity of ℓ

$$\ell(\hat{\theta}) = \sup \ell$$

2: You might try defining $\hat{\theta}$ to be a suitably chosen critical point of ℓ .

3: You might try to prove that

$$P(\text{card}(\{\theta : \ell(\theta|X_1, \dots, X_n) = \sup_{\phi} \ell(\phi|X_1, \dots, X_n)\}) = 1) = 1$$

In other words it might be true that the set of θ where ℓ achieves its maximum is almost surely a singleton when the x s are actually a data set.

I am going to follow method 2 since this is the one which works most generally.

For a vector (x_1, \dots, x_n) we define the order statistics

$$x_{(1)} \leq \dots \leq x_{(n)}$$

be the entries in the vector sorted into non-decreasing order. If $n = 2m - 1$ is odd then set

$$g_n(x_1, \dots, x_n) = x_{(m)}$$

If $n = 2m$ set

$$g_n(x_1, \dots, x_n) = (x_{(m)} + x_{(m+1)})/2.$$

Now define

$$\tilde{\theta}_n = g(X_1, \dots, X_n)$$

Lemma 4 *If X_1, \dots, X_n are iid from a distribution F with the properties:*

1. $F(0) = 1/2$.
2. For each $\epsilon > 0$

$$F(-\epsilon) < 1/2 < F(\epsilon)$$

Then $\tilde{\theta}_n$ converges almost surely to 0.

Remark: Part of the theorem is that

$$A \equiv \{\omega : \tilde{\theta}_n \rightarrow 0\}$$

is an event. Here is how we work our way back to primitive notions to prove this. Our formal description, in set notation, of this event shows that A is an event if each $\tilde{\theta}_n$ is a random variable. In turn, since the X_i are random variables, we need only check that g_n is Borel. This reduces to the assertion that

$$(x_1, \dots, x_n) \mapsto x_{(k)}$$

is Borel for each k and n . But

$$\{(x_1, \dots, x_n) : x_{(k)} < t\} = \{(x_1, \dots, x_n) : \sum_1^n 1(x_i < t) \geq k\}$$

Since a sum of Borel functions is Borel we are reduced to proving that the map

$$(x_1, \dots, x_n) \mapsto 1(x_i < t)$$

is Borel for each t and i . This is equivalent to

$$\{(x_1, \dots, x_n) : x_i < t\}$$

is Borel. Since this last set is actually open it is Borel.

Now to prove the lemma we begin by formalizing an argument we have used several times.

Lemma 5 *Suppose Y_n is a sequence of random variables. Then $Y_n \rightarrow 0$ almost surely is equivalent to $P(C_\epsilon) = 1$ for each $\epsilon > 0$ where*

$$C_\epsilon = \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{|Y_n| \leq \epsilon\}$$

So fix $\epsilon > 0$. For each x the random variables Y_1, Y_2, \dots defined by $Y_k = 1(X_k \leq x) - F(x)$ are independent and identically distributed with mean 0. According to the SLLN there is a null set N_x such that for all $\omega \notin N_x$ we have

$$\frac{1}{n} \sum_1^n Y_k \rightarrow 0.$$

Let $N = N_\epsilon \cup N_{-\epsilon}$. Then N is a null set. If $\omega \notin N$ then

$$\frac{1}{n} \sum_1^n 1(X_k \leq \epsilon) \rightarrow F(\epsilon) > 1/2$$

and

$$\frac{1}{n} \sum_1^n 1(X_k \leq -\epsilon) \rightarrow F(-\epsilon) < 1/2$$

For any such ω there is an M such that for all $n \geq M$ the number of X_i exceeding ϵ is less than $n/2$ and the number of X_i less than $-\epsilon$ is less than $n/2$. Thus for such ω , and $n \geq M$

$$-\epsilon \leq \tilde{\theta}_n \leq \epsilon$$

In other words the set $C_\epsilon^c \subset N$ so $P(C_\epsilon) = 1$. •

Now back to the Cauchy problem. For a vector (x_1, \dots, x_n) we define $h_n(x_1, \dots, x_n)$ to be that root of

$$\ell'(\theta|x_1, \dots, x_n) = \sum \frac{2(x_i - \theta)}{1 + (x_i - \theta)^2} = 0$$

which is closest to $g_n(x_1, \dots, x_n)$ (here g_n is the Borel function used in defining the median above). If $g_n(x_1, \dots, x_n)$ is exactly midway between two roots which are tied for closest we define h_n to be the root smaller than g_n .

It is possible to prove that this defines a Borel function from \mathbb{R}^n to \mathbb{R} . Now define

$$\hat{\theta}_n = h_n(X_1, \dots, X_n)$$

I claim that if X_1, X_2, \dots are iid Cauchy(0) then

$$\hat{\theta}_n \rightarrow 0 \tag{4.1}$$

almost surely.

To prove this fix $\epsilon > 0$ we prove that $P(C_\epsilon) = 1$ where

$$C_\epsilon = \bigcup_{N=1}^{\infty} \bigcap_{n=N}^{\infty} \{|\hat{\theta}_n| \leq \epsilon\}$$

In the discussion which follows we suppress dependence of ℓ on the data whenever possible. We use superscript \prime to denote differentiation with respect to θ . We will have use of the

notation:

$$\begin{aligned} L_i(\theta) &= -\log(f(X_i, \theta)) + \log(f(X_i, \theta_o)) \\ U_i(\theta) &= L'_i(\theta) \\ V_i(\theta) &= L''_i(\theta) \\ W_i(\theta) &= L'''_i(\theta) \end{aligned}$$

Here is how we proceed. We find an event D_ϵ inside C_ϵ with $P(D_\epsilon) = 1$. To define this new event note that if

1. ℓ' has a unique root over $[-3\epsilon, 3\epsilon]$ and
2. that root is actually in $[-\epsilon, \epsilon]$ and
3. $|\tilde{\theta}_n| \leq \epsilon$

then the root of ℓ' closest to $\tilde{\theta}_n$ is actually the root in points 1 and 2 and so

$$|\hat{\theta}_n| \leq \epsilon$$

Define $D_\epsilon^{(1)}$ to be the event that there is an N such that for all $n \geq N$ and all $|\theta| \leq 3\epsilon$ we have

$$\ell''(\theta|X_1, \dots, X_n) < 0$$

Define $D_\epsilon^{(2)}$ to be the event that there is an N such that for all $n \geq N$

$$\{\ell'(\epsilon|X_1, \dots, X_n) < 0$$

and

$$\{\ell'(-\epsilon|X_1, \dots, X_n) > 0$$

Finally define $D_\epsilon^{(3)}$ to be the event that there is an N such that for all $n \geq N$

$$|\tilde{\theta}_n| \leq \epsilon$$

We have already shown that $P(D_\epsilon^{(3)}) = 1$. Next I show that $P(D_\epsilon^{(2)}) = 1$. For the Cauchy case

$$U_k(\epsilon) = \frac{2(X_k - \epsilon)}{1 + (X_k - \epsilon)^2}$$

and

$$U_k(-\epsilon) = \frac{2(X_k + \epsilon)}{1 + (X_k + \epsilon)^2}$$

Then

$$\frac{1}{n} \ell''(\epsilon|X_1, \dots, X_n) = \overline{U_n(\epsilon)}$$

and

$$\frac{1}{n} \ell''(-\epsilon|X_1, \dots, X_n) = \overline{U_n(-\epsilon)}$$

Thus $D_\epsilon^{(2)}$ is the event that there is an N such that for all $n \geq N$

$$\overline{U_n(\epsilon)} < 0 \quad \text{and} \quad \overline{U_n(-\epsilon)} > 0$$

Since each of $\overline{U_n(-\epsilon)}$ and $\overline{U_n(\epsilon)}$ is an average of iid variates it suffices to show that

$$\begin{aligned} \mathbb{E}(U_k(\epsilon)) &< 0 \\ \mathbb{E}(U_k(-\epsilon)) &> 0. \end{aligned}$$

In fact

$$\begin{aligned} \mathbb{E}(U_k(\epsilon)) &= \frac{-2\pi\epsilon}{\epsilon^2 + 4} < 0 \\ \mathbb{E}(U_k(-\epsilon)) &= \frac{2\pi\epsilon}{\epsilon^2 + 4} > 0. \end{aligned}$$

This argument is not very easy to generalize since it hinged on an exact computation of a moment. A much more general tactic uses Jensen's inequality.

If l is smaller at $-\epsilon$ and at ϵ than it is at 0 then there must be a critical point in $[-\epsilon, \epsilon]$, that is, a root of l' . To deal with this recall the definition

$$L_i(\theta) = \log(1 + X_i^2) - \log(1 + (X_i - \theta)^2).$$

Define $D_\epsilon^{(4)}$ be the event that there is an N such that for all $n \geq N$

$$\left\{ \sum L_i(\epsilon) < 0, \sum L_i(-\epsilon) < 0 \right\}.$$

Define

$$\mu(\epsilon) = \mathbb{E}(L_i(\epsilon))$$

The strong law of large numbers shows that $P(D_\epsilon^{(5)}) = 1$ where

$$D_\epsilon^{(5)} = \left\{ \sum L_i(\epsilon)/n \rightarrow \mu(\epsilon), \sum L_i(-\epsilon)/n \rightarrow \mu(-\epsilon) \right\}$$

I claim that for all $\epsilon \neq 0$ $\mu(\epsilon) < 0$. If so then

$$D_\epsilon^{(5)} \subset D_\epsilon^{(4)}$$

and so $P(D_\epsilon^{(4)}) = 1$.

To prove the claim we apply Jensen's inequality:

Proposition 1 (Jensen) *Suppose Y is a random variable and ϕ is a function which is convex on an interval (a, b) with $P(a < Y < b) = 1$. Assume $\mathbb{E}(|Y|) < \infty$. Then*

$$\phi(\mathbb{E}(Y)) \leq \mathbb{E}(\phi(Y))$$

If ϕ is strictly convex then the inequality is strict unless $\text{Var}(Y) = 0$.

Jargon: ϕ is convex if for each x, y and $\lambda \in (0, 1)$

$$\phi(\lambda x + (1 - \lambda)y) \leq \lambda\phi(x) + (1 - \lambda)\phi(y)$$

We call ϕ strictly convex if the inequality is strict.

If ϕ is twice differentiable and $\phi'' \geq 0$ then ϕ is convex; a strict inequality shows ϕ is strictly convex. We apply Jensen's inequality with $\phi(x) = -\log(x)$. We apply it to $Y = g(X)/f(X)$ where X has density f and g is some other density. Then

$$E\{-\log(Y)\} > -\log\{E(Y)\}$$

But the latter is

$$\log \left\{ \int \frac{g(x)}{f(x)} f(x) dx \right\} = \log(1) = 0$$

Technically we ought to be careful. The interval (a, b) in the inequality is $(0, \infty)$. The assumption

$$P(0 < Y < \infty) = 1$$

deserves some discussion. If $f(x) = 0$ for some places where $g(x)$ is not 0 then

$$E \left\{ \frac{g(X)}{f(X)} \right\} = \int g(x) 1_{(f(X) > 0)} dx$$

which might be less than 1. This just makes the inequality stronger, however.

The other technical detail is that $g(x)$ might be 0 some places where $f(x)$ is not 0. This might mean $P(Y = 0) > 0$. On the event $Y = 0$ we will agree to take $-\log(Y) = \infty$ and conclude

$$E\{-\log(Y)\} = \infty$$

In any case we find

$$E\{-\log(Y)\} > 0$$

or

$$E[\log\{g(X)\} - \log\{f(X)\}] < 0.$$

Applied to our Cauchy problem we have shown $\mu(\theta) < 0$ for all $\theta \neq 0$. Hence $P(D_\epsilon^{(4)}) = 1$.

Finally we consider $D_\epsilon^{(1)}$. Up to now we have been able to make do with an arbitrary ϵ . In this case, however, the result holds only for small $\epsilon > 0$. First consider

$$\frac{1}{n} \ell''(0 | X_1, \dots, X_n)$$

According to the strong law of large numbers this converges almost surely to

$$E \left\{ \frac{2(1 - X^2)}{(1 + X^2)^2} \right\} = -\frac{1}{2} < 0$$

Now you can check that

$$\left| \frac{1}{n} \ell'''(\theta | X_1, \dots, X_n) \right| < 4.$$

(In fact each term in ℓ''' may be shown to be bounded by $3/2 + \sqrt{2}$.) As a result

$$\frac{1}{n} |\ell''(\theta|X_1, \dots, X_n) - \ell''(0|X_1, \dots, X_n)| \leq 4|\theta|$$

Pick $\epsilon > 0$ so that $4\epsilon < \pi/2$. If B is the event

$$\frac{1}{n} \ell''(0|X_1, \dots, X_n) \rightarrow -\frac{1}{2}$$

and ω is in B then there is an N such that for $n \geq N$ we have

$$\frac{1}{n} \ell''(\theta|X_1, \dots, X_n) < 0$$

for all $|\theta| < \epsilon$.

This proves that for all $3\epsilon < \pi/8$

$$P(D^{(1)}(\epsilon)) = 1.$$

We have now shown that for $\epsilon < \pi/24$

$$P(D^{(1)}(\epsilon) \cap D_\epsilon^{(2)} \cap D_\epsilon^{(3)}) = 1$$

For ω in this event we have that there is an N such that

$$|\hat{\theta}_n| \leq \epsilon$$

for all $n \geq N$. This establishes the result.

General Case

Now consider a parametric family

$$\{f(x|\theta); a < \theta < b\}$$

Let θ_o be the true value of θ , that is, the value used in making probability calculations. Let A_ϵ be the event that there is an N such that for all $n \geq N$ the log-likelihood has a local maximum on the interval $[\theta_o - \epsilon, \theta_o + \epsilon]$. We have proved quite generally that

$$P(A_\epsilon) = 1 \tag{4.2}$$

If we add the assumption

$$\ell \text{ has a continuous derivative} \tag{A}$$

then letting B_ϵ be the event that there is an N such that for all $n \geq N$ there is a critical point of ℓ in $(\theta_o - \epsilon, \theta_o + \epsilon)$ which is a local maximum of ℓ we have proved

$$P(B_\epsilon) = 1 \tag{4.3}$$

The event $B = \cap_{\epsilon} B_{\epsilon}$ then has probability 1. On this event there is a sequence of roots of the likelihood equations which is consistent.

The remaining question is whether or not we can prove, under reasonably general conditions, that there is probably only one root near θ_o .

Now consider the event that ℓ' is monotone on $[\theta_o - \epsilon, \theta_o + \epsilon]$. Our previous proof was based on showing that the next derivative was negative at θ_o and did not change much over a small enough interval.

Behaviour at θ_o is essentially the behaviour of

$$\frac{1}{n} \sum V_i(\theta_o)$$

which converges almost surely to

$$E(L_1''(\theta_o))$$

I claim this is negative for regular families.

Begin with

$$1 = \int f(x, \theta) dx$$

Differentiating with respect to θ gives

$$\begin{aligned} 0 &= \frac{d}{d\theta} \int f(x, \theta) dx \\ &= \lim_{\epsilon \rightarrow 0} \int \frac{f(x, \theta + \epsilon) - f(x, \theta)}{\epsilon} dx \end{aligned}$$

In order to take the limit inside the integral sign we must prove that for any sequence $\epsilon_n \rightarrow 0$

$$\lim \int \frac{f(x, \theta + \epsilon_n) - f(x, \theta)}{\epsilon_n} dx = \int \frac{\partial}{\partial \theta} f(x, \theta) dx$$

This is normally done by applying the dominated convergence theorem. If f is continuously differentiable with respect to θ then the difference quotient is exactly equal to

$$\frac{\partial}{\partial \theta} f(x, \theta_n^*)$$

The quantity θ_n^* depends on both n and x . One tactic is to try to compute

$$M(x, \epsilon) = \sup \left\{ \left| \frac{\partial}{\partial \theta} f(x, \theta) \right|; |\theta - \theta_o| \leq \epsilon \right\}$$

and show that

$$\int M(x, \epsilon) dx < \infty$$

which would permit application of dominated convergence.

Assuming that we can apply the dominated convergence theorem we get

$$\begin{aligned} 0 &= \int \frac{\partial}{\partial \theta} f(x, \theta) dx \\ &= \int \frac{\partial \log f(x, \theta)}{\partial \theta} f(x, \theta) dx \\ &= E(U_k(\theta_o)) \end{aligned}$$

Differentiating again and again passing limits through integrals gives

$$E(U_k^2(\theta_o)) = -E(V_k(\theta_o))$$

This shows that

$$\frac{1}{n} \ell''(\theta_o) \rightarrow E(V_k(\theta_o)) < 0$$

almost surely.

Next we consider

$$\frac{1}{n} \ell''(\theta) - \ell''(\theta_o)$$

For a three times continuously differentiable ℓ there is a θ_n^* (which is random but between θ_o and θ) such that this difference is

$$\frac{\theta - \theta_o}{n} \sum W_i(\theta_n^*)$$

Define

$$M_i(\epsilon) = \sup\{|W_i(\theta)| : |\theta - \theta_o| \leq \epsilon\}$$

The M_i are iid. If for some $\epsilon > 0$ the M_i are integrable then the strong law of large numbers shows that

$$\limsup \left| \frac{1}{n} \{\ell''(\theta) - \ell''(\theta_o)\} \right| \leq \epsilon E(M_1(\epsilon))$$

almost surely. The right hand side of this inequality can be made arbitrarily small by choosing ϵ small enough. Pick ϵ so small that the bound is strictly smaller than

$$I(\theta_o) \equiv -E(V_k(\theta_o))$$

we have the result that

$$P(E_\epsilon) = 1$$

where E_ϵ is the event that there is an N such that for all $n \geq N$ we have ℓ' is monotone decreasing on $[\theta_o - \epsilon, \theta_o + \epsilon]$.

Chapter 5

Markov Chains

Definition: A stochastic process is a family $\{X_i; i \in I\}$ of random variables where I is a set, called the index set. Often I is a subset of \mathbb{R} such as $[0, \infty)$, $[0, 1]$ \mathbb{Z} or \mathbb{N} . When I is an interval we speak of a continuous time stochastic process; often we use the letter t as a typical element of the index set. When $I \subset \mathbb{Z}$ we call X a discrete time stochastic process.

Generally we expect all the X_n to take values in the same **state space** S . In this chapter we will take S to be a finite or countable set so that each X_n is discrete. Usually we identify S with one of \mathbb{Z} , \mathbb{N} or $\{0, \dots, m\}$ for some finite m .

Definition: A Markov Chain is a stochastic process $X_n; n \in \mathbb{N}$. taking values in a finite or countable set S such that for every n and every event of the form

$$A = \{(X_0, \dots, X_{n-1}) \in B \subset S^{n+1}\}$$

we have

$$P(X_{n+1} = j | X_n = i, A) = P(X_1 = j | X_0 = i) \tag{5.1}$$

We then adopt the notation that \mathbf{P} is the (possibly infinite) array with elements

$$P_{ij} = P(X_1 = j | X_0 = i)$$

indexed by $i, j \in S$.

Conditional Probability: An Aside

If $P(A) > 0$ we define

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

In this chapter all random variables are discrete and so when we condition on values of a finite list of variables we will be conditioning on an event of positive probability. If we try to move to continuous state spaces we will have to improve our understanding of conditional probability and expectation.

WARNING: in (5.1) we require the condition to hold **only** when

$$P(X_n = i, A) > 0$$

End of Aside

Definition: \mathbf{P} is the (one step) **transition matrix** of the Markov Chain.

Evidently the entries in \mathbf{P} are non-negative and

$$\sum_j P_{ij} = 1 \tag{5.2}$$

for all $i \in S$. Any such matrix is called **stochastic**.

We define powers of \mathbf{P} by

$$(\mathbf{P}^n)_{ij} = \sum_k (\mathbf{P}^{n-1})_{ik} P_{kj}$$

Notice that even if S is infinite these sums converge absolutely.

5.1 Chapman-Kolmogorov Equations

A central tactic in Markov Chain calculations is conditioning on intermediate steps to compute probabilities of events. We do that now to compute

$$P(X_{l+n} = j | X_l = i)$$

We will condition on X_{l+n-1} :

$$\begin{aligned} P(X_{l+n} = j | X_l = i) &= \sum_k P(X_{l+n} = j, X_{l+n-1} = k | X_l = i) \\ &= \sum_k P(X_{l+n} = j | X_{l+n-1} = k, X_l = i) P(X_{l+n-1} = k | X_l = i) \\ &= \sum_k P(X_1 = j | X_0 = k) P(X_{l+n-1} = k | X_l = i) \\ &= \sum_k P(X_{l+n-1} = k | X_l = i) \mathbf{P}_{ik} \end{aligned}$$

Now condition on X_{l+n-2} to get

$$P(X_{l+n} = j | X_l = i) = \sum_{k_1 k_2} \mathbf{P}_{k_1 k_2} \mathbf{P}_{k_2 j} P(X_{l+n-2} = k_1 | X_l = i)$$

Notice that the sum over k_2 computes the $K_{1,j}$ entry in the matrix $\mathbf{P}\mathbf{P}$ which we denote \mathbf{P}^2 . In other words

$$P(X_{l+n} = j | X_l = i) = \sum_{k_1} (\mathbf{P}^2)_{k_1 j} P(X_{l+n-2} = k_1 | X_l = i)$$

We may now prove by induction on n that

$$P(X_{l+n} = j | X_l = i) = (\mathbf{P}^n)_{ij}.$$

Finally we have now proved the Chapman-Kolmogorov equations:

$$P(X_{l+m+n} = j | X_l = i) = \sum_k P(X_{l+m} = k | X_l = i) P(X_{l+m+n} = j | X_{l+m} = k)$$

These are simply a restatement of the identity

$$\mathbf{P}^{n+m} = \mathbf{P}^n \mathbf{P}^m .$$

Remark: It is important to notice that these probabilities depend on m and n but **not** on l . We say the chain has **stationary** transition probabilities. A more general definition of Markov chain than (5.1) is

$$P(X_{n+1} = j | X_n = i, A) = P(X_{n+1} = j | X_n = i) .$$

Notice that the right hand side of this formula is now permitted to depend on n . If we define $\mathbf{P}^{n,m}$ to be the matrix with i, j th entry

$$P(X_m = j | X_n = i)$$

for $m > n$ then we can prove

$$\mathbf{P}^{r,s} \mathbf{P}^{s,t} = \mathbf{P}^{r,t}$$

These equations are also called Chapman-Kolmogorov equations. This chain does not have stationary transitions.

Remark: The calculations above involve sums in which all terms are positive. They therefore apply even if the state space S is countably infinite.

5.1.1 Extensions of the Markov Property

Consider a function $f(x_0, x_1, \dots)$ defined on the set S^∞ of all infinite sequences of points in S . (This family of functions includes, in a natural way, all functions which depend on only finitely many of the co-ordinates.) Let B_n be the event

$$f(X_n, X_{n+1}, \dots) \in C$$

for some suitable C in the range space of f . Then

$$P(B_n | X_n = x, A) = P(B_0 | X_0 = x) \tag{5.3}$$

for any event A of the form

$$\{(X_0, \dots, X_{n-1}) \in D\}$$

Also

$$P(AB_n | X_n = x) = P(A | X_n = x) P(B_n | X_n = x) \tag{5.4}$$

The latter is usually expressed as “given the present the past and future are conditionally independent.”

To prove the first of these assertions you use monotone class arguments. The general result can be deduced from the special case where

$$B_n = \{(X_{n+1} = x_1, \dots, X_{n+m} = x_m)\}$$

For this B_n the left hand side of (5.3) may be evaluated by repeated conditioning as in Chapman-Kolmogorov to get

$$\mathbf{P}_{x,x_1} \mathbf{P}_{x_1,x_2} \cdots \mathbf{P}_{x_{m-1},x_m}$$

and so may the right. For general events defined only from X_n, \dots, X_{n+m} you sum over appropriate vectors x, x_1, \dots, x_m . The general case requires a measure theory argument via monotone class techniques.

To prove (5.4) write

$$\begin{aligned} P(AB_n|X_n = x) &= P(B_n|X_n = x, A)P(A|X_n = x) \\ &= P(B_n|X_n = x)P(A|X_n = x) \end{aligned}$$

using (5.3).

5.2 Classification of States

If an entry \mathbf{P}_{ij} is 0 it is not possible to go from state i to state j in one step. It may be possible to make the transition in some larger number of steps, however. We say i **leads to** j (or j is accessible from i) if there is an integer $n \geq 0$ such that

$$P(X_n = j|X_0 = i) > 0.$$

We use the notation $i \rightsquigarrow j$. By defining \mathbf{P}^0 to be the identity matrix \mathbf{I} we arrive at $i \rightsquigarrow j$ if there is an $n \geq 0$ for which $(\mathbf{P}^n)_{ij} > 0$.

We say states i and j **communicate** if $i \rightsquigarrow j$ and $j \rightsquigarrow i$. We write $i \leftrightarrow j$ if i and j communicate. Communication is an equivalence relation, that is, a reflexive, symmetric and transitive relation on the states of S . More precisely:

Reflexive: for all i we have $i \leftrightarrow j$.

Symmetric: if $i \leftrightarrow j$ then $j \leftrightarrow i$.

Transitive: if $i \leftrightarrow j$ and $j \leftrightarrow k$ then $i \leftrightarrow k$.

Proof: Reflexive follows from our inclusion of $n = 0$ in the definition of leads to. Symmetry is obvious. To check transitivity it suffices to check that $i \rightsquigarrow j$ and $j \rightsquigarrow k$ imply that $i \rightsquigarrow k$. But if $(\mathbf{P}^m)_{ij} > 0$ and $(\mathbf{P}^n)_{jk} > 0$ then

$$(\mathbf{P}^{m+n})_{ik} = \sum_l (\mathbf{P}^m)_{il} (\mathbf{P}^n)_{lk} \geq (\mathbf{P}^m)_{ij} (\mathbf{P}^n)_{jk} > 0$$

Any equivalence relation on a set partitions the set into **equivalence classes**; two elements are in the same equivalence class if and only if they are equivalent. So communication partitions the state space into equivalence classes called **communicating classes**.

Here is an example:

$$\mathbf{P} = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

To find communicating classes you start with say state 1 and see where it leads. We see $1 \rightsquigarrow 2$, $1 \rightsquigarrow 3$ and $1 \rightsquigarrow 4$ in row 1. Turning to row 4 we see that $4 \rightsquigarrow 1$. This proves (via transitivity) that 1, 2, 3 and 4 are all in the same communicating class. On the other hand I claim none of these states leads to any of 5, 6, 7 or 8. This is almost obvious because of the zeroes in the upper right hand corner of \mathbf{P} . For clarity, however, I will give a precise proof. So suppose $i \in \{1, 2, 3, 4\}$ and $j \in \{5, 6, 7, 8\}$. Then $(\mathbf{P}^n)_{ij}$ is a sum of products of terms of the form \mathbf{P}_{kl} . This cannot be positive unless there is a sequence $i_0 = i, i_1, \dots, i_n = j$ with $\mathbf{P}_{i_{k-1}, i_k} > 0$ for $k = 1, \dots, n$. Consider the first k for which $i_k \in \{5, 6, 7, 8\}$. Then $i_{k-1} \in \{1, 2, 3, 4\}$ and so $\mathbf{P}_{i_{k-1}, i_k} = 0$.

This proves $\{1, 2, 3, 4\}$ is a communicating class. Now look at state 5. Clearly $5 \rightsquigarrow 1$, $5 \rightsquigarrow 2$, $5 \rightsquigarrow 3$ and $5 \rightsquigarrow 4$. None of these lead to any of $\{5, 6, 7, 8\}$ so $\{5\}$ must be communicating class. Similarly you may check $\{6\}$ and $\{7, 8\}$ are communicating classes.

In the example the states 5 and 6 have a special property. Each time you are in either state you run a risk of going to one of the states 1, 2, 3 or 4. Eventually you will make such a transition and then never return to state 5 or 6. States 5 and 6 are transient. To make this precise we define hitting times:

$$T_k = \min\{n > 0 : X_n = k\}$$

We define

$$f_k = P(T_k < \infty | X_0 = k)$$

A state k is called transient if $f_k < 1$ and recurrent if $f_k = 1$.

Let N_k be the number of times the chain ever is in state k . I will show that if $f_i < 1$ then N_k has a Geometric distribution

$$P(N_k = r | X_0 = k) = f_k^{r-1}(1 - f_k)$$

for $r = 1, 2, \dots$. On the other hand if $f_i = 1$ then I will show

$$P(N_k = \infty | X_0 = k) = 1$$

To prove this theorem I am going to strengthen the Markov property to what is called the strong Markov Property. A stopping time for the Markov chain is a random variable T taking values in $\{0, 1, \dots\} \cup \{\infty\}$ such that for each finite k there is a function f_k such that

$$1(T = k) = f_k(X_0, \dots, X_k)$$

Notice that T_k is a stopping time.

It is now convenient to introduce some standard shorthand notation. By

$$P^x(A)$$

we mean

$$P(A|X_0 = x)$$

and similarly we define

$$E^x(Y) = E(Y|X_0 = x)$$

Our goal is to explain and prove the formula

$$E(f(X_T, \dots)|X_T, \dots, X_0) = E^{X_T}(f(X_0, \dots))$$

I will do a special case for illustration and then discuss conditional probabilities and expectations.

Consider first

$$P(X_{T+1} = j|X_T = i)$$

I will prove this is just

$$\mathbf{P}_{ij} = P^i(X_1 = j)$$

We evaluate

$$\begin{aligned} P(X_{T+1} = j|X_T = i) &= \frac{P(X_{T+1} = j, X_T = i)}{P(X_T = i)} \\ &= \frac{\sum_k P(X_{T+1} = j, X_T = i, T = k)}{\sum_k P(X_T = i, T = k)} \\ &= \frac{\sum_k P(X_{k+1} = j, X_k = i, T = k)}{\sum_k P(X_k = i, T = k)} \\ &= \frac{\sum_k P(X_{k+1} = j|X_k = i, T = k)P(X_k = i, T = k)}{\sum_k P(X_k = i, T = k)} \\ &= \frac{\sum_k P(X_1 = j|X_0 = i)P(X_k = i, T = k)}{\sum_k P(X_k = i, T = k)} \end{aligned}$$

$\mathbf{P}_{i,j}$

Notice the use in deleting $T = k$ from the condition of the fact that $T = k$ is an event defined in terms of X_0, \dots, X_k .

There are some technical problems with this proof, however:

- It might be that $P(T = \infty) > 0$. In that case what are X_T and X_{T+1} on the event $T = \infty$. The answer is that the formula must be revised to condition also on $T < \infty$.
- It is not obvious that the event being conditioned on has positive probability. The answer is that we must seek to prove the formula only for those stopping times where $\{T < \infty\} \cap \{X_T = i\}$ has positive probability.

We will now fix up these technical details.

5.2.1 Conditional distributions and expectations

When X and Y are discrete we have

$$E(Y|X = x) = \sum_y P(Y = y|X = x)$$

for any x for which $P(X = x)$ is positive. This defines a function of x . When this function is evaluated at X we get a random variable which is a function of X denoted

$$E(Y|X)$$

Here are some properties of that function.

1. Suppose A is a function defined on the range of X . Then

$$E(A(X)Y|X = x) = A(x)E(Y|X = x)$$

and so

$$E(A(X)Y|X) = A(X)E(Y|X)$$

The second assertion follows by definition from the first. To prove the first note that if $Z = A(X)Y$ then Z is discrete and

$$P(Z = z) = \sum_{x,y} P(Y = y, X = x)1(z = A(x)y)$$

Also

$$\begin{aligned} P(Z = z|X = x) &= \frac{\sum_y P(Y = y, X = x)1(z = A(x)y)}{P(X = x)} \\ &= \sum_y P(Y = y|X = x)1(z = A(x)y) \end{aligned}$$

Thus

$$\begin{aligned} E(Z|X = x) &= \sum_z zP(Z = z|X = x) \\ &= \sum_z \sum_y zP(Y = y|X = x)1(z = A(x)y) \\ &= \sum_z \sum_y A(x)yP(Y = y|X = x)1(z = A(x)y) \\ &= A(x) \sum_y yP(Y = y|X = x) \sum_z 1(z = A(x)y) \\ &= A(x) \sum_y yP(Y = y|X = x) \end{aligned}$$

2. Repeated conditioning: if X, Y and Z discrete then

$$\begin{aligned} E\{E(Z|X, Y)|X\} &= E(Z|X) \\ E\{E(Y|X)\} &= E(Y) \end{aligned}$$

3. Additivity

$$E(Y + Z|X) = E(Y|X) + E(Z|X)$$

4. Putting the first two items together gives

$$E\{A(X)E(Y|X)\} = E(A(X)Y) \quad (5.5)$$

In fact the general definition of $E(Y|X)$ when X and Y are not assumed to be discrete is that it is a random variable which is a measurable function of X satisfying (5.5). The existence of this function is a measure theory problem.

Now suppose that $f(x_0, x_1, \dots)$ is a (measurable) function on $S^{\mathbb{N}}$. Put

$$Y_n = f(X_n, X_{n+1}, \dots).$$

Assume that $E(|Y_0| | X_0 = x) < \infty$ for all x . Then I claim

$$E(Y_n | X_n, A) = E^{X_n}(Y_0) \quad (5.6)$$

whenever A is any event defined in terms of X_0, \dots, X_n .

Proof: : The family of functions f for which the claim is true includes all indicators since then the expectations are just probabilities and (5.6) reduces to (5.3). The family of functions for which the claim is true is a vector space (the point being that if f and g are in the family then so is $af + bg$ for any constants a and b). Thus the family of functions f for which the claim is true includes all simple functions. The family of functions f for which the claim is true is closed under the action of taking monotone increasing limits (of non-negative f_n) by the Monotone Convergence theorem. Hence the claim is true for every non-negative integrable f . Finally the claim follows for every integrable f by linearity.

Aside on “measurable”: what sorts of events can be defined in terms of a family $\{Y_i : i \text{ in } I\}$? It seems natural to regard any event of the form $(Y_{i_1}, \dots, Y_{i_k}) \in C$ as a definition in terms of the family for any finite set i_1, \dots, i_k and any (Borel) set C in S^k . When the state space is countable, incidentally, the natural notion of Borel is to make each single point $(s_1, \dots, s_k) \in S^k$ Borel. This makes every subset of S^k Borel.

It is also natural to agree that if you can define each of a sequence of events A_n in terms of the Y s then the definition “there exists an n such that (definition of A_n) ...” defines $\cup A_n$. Similarly if A is definable in terms of the Y s then A^c can be defined from the Y s by just inserting the phrase “It is not true that” in front of the definition of A . In other words the family of events definable in terms of the family $\{Y_i : i \text{ in } I\}$ is a σ -field which includes every event of the form $(Y_{i_1}, \dots, Y_{i_k}) \in C$. We call the smallest such σ -field, $\mathcal{F}(\{Y_i : i \text{ in } I\})$, the σ -field generated by the family $\{Y_i : i \text{ in } I\}$.

This discussion permits some shorthand. We define

$$\mathcal{F}_n = \mathcal{F}(\{X_0, \dots, X_n\})$$

and

$$\mathcal{F}_\infty = \mathcal{F}(\{X_0, X_1, \dots\})$$

Finally I want to discuss conditioning on X_0, \dots, X_T where T is a stopping time. There are several ways to deal with this problem. One is to simply say that the vector X_0, \dots, X_T takes values in a somewhat exotic space, namely:

$$\left(\bigcup_{n=1}^{\infty} S^n \right) \cup S^{\mathbb{N}}$$

We define

$$\mathcal{F}_T = \mathcal{F}(X_0, \dots, X_T)$$

The random vector X_0, \dots, X_T is discrete on the event $T < \infty$ but not discrete on the event $T = \infty$. Our interest, however, is really on the event $T < \infty$.

Suppose X is discrete and $X^* = g(X)$ is a one to one transformation of X . Since $X = x$ is the same event as $X^* = g(x)$ we find

$$E(Y|X = x) = E(Y|X^* = g(x))$$

Let $h^*(u)$ denote the function $E(Y|X^* = u)$ and $h(u) = E(Y|X = u)$. Then

$$h(x) = h^*(g(x))$$

Thus

$$h(X) = h^*(g(X)) = h^*(X^*)$$

This just means

$$E(Y|X) = E(Y|X^*)$$

This formula deserves some interpretation. First of all it says something obvious. Here's an example. I toss a coin $n = 20$ times. Let Y denote the indicator that the first toss is a heads. Let X be the number of heads and X^* be the number of tails. Then what is being said is

$$E(Y|X = 17) = E(Y|X^* = 3)$$

In fact for a general k and n

$$E(Y|X = k) = \frac{k}{n}$$

so

$$E(Y|X) = \frac{X}{n}$$

At the same time

$$E(Y|X^* = j) = \frac{n-j}{n}$$

so

$$E(Y|X^*) = \frac{n - X^*}{n}$$

But of course $X = n - X^*$ so these are just two ways of describing the same random variable.

Another interpretation is this. The random variable X partitions Ω into a countable set of events of the form $X = x$. For any other random variable X^* which partitions Ω into the same events the values of $E(Y|X^* = x^*)$ are just the same as the values of $E(Y|X = x)$ but labelled. When we form $E(Y|X)$ we are simply taking the value ω , computing the value of $X(\omega)$ to see which member A of the partition we are talking about and then writing down the corresponding $E(Y|A)$. This should make it apparent that we would get the same answer if we used X^* as long as X and X^* define the same partition of Ω .

For random variables which are not discrete the notion of a partition must be replaced by a σ -field. In fact suppose X and X^* are two random variables such that

$$\mathcal{F}(X) = \mathcal{F}(X^*)$$

Then:

- There is a Borel g which is one to one and has a one to one Borel inverse such that $X^* = g(X)$.
- $E(Y|X) = E(Y|X^*)$ almost surely.

In other words $E(Y|X)$ depends *only* on the σ -field generated by X . We write

$$E(Y|\mathcal{F}(X)) = E(Y|X)$$

Definition: Suppose \mathcal{G} is a sub- σ -field of \mathcal{F} . We say that X is \mathcal{G} measurable if, for every Borel B

$$\{\omega : X(\omega) \in B\} \in \mathcal{G}.$$

Definition: $E(Y|\mathcal{G})$ is any \mathcal{G} measurable random variable such that for every \mathcal{G} measurable random variable A we have

$$E(AY) = E\{AE(Y|\mathcal{G})\}.$$

Again the existence of such a random variable is a measure theory problem.

Now we define \mathcal{F}_T to be the collection of all events B such that

$$B \cap \{T \leq n\} \in \mathcal{F}_n$$

for each n . It is easy to prove that \mathcal{F}_T is a σ -field and that T and X_T are \mathcal{F}_T measurable.

Finally I rewrite the Strong Markov Property: if

- $f : S^{\mathbb{N}} \mapsto \mathbb{R}$ is Borel and
- $Y_T = f(X_T, X_{T+1}, \dots)1(T < \infty)$ and
- $Y_0 = f(X_0, X_1, \dots)$ and

- $E^x(|Y_0|) < \infty$ for all x .

then

$$1(T < \infty)E(f(X_T, X_{T+1}, \dots)|\mathcal{F}_T) = 1(T < \infty)E^{X_T}(f(X_0, X_1, \dots)) \quad (5.7)$$

Proof of the Strong Markov Property.

Since T is \mathcal{F}_T measurable we have

$$1(T < \infty)E\{f(X_T, X_{T+1}, \dots)|\mathcal{F}_T\} = E\{1(T < \infty)f(X_T, X_{T+1}, \dots)|\mathcal{F}_T\}$$

Now suppose A is \mathcal{F}_T measurable. Then we need only prove that

$$E[A1(T < \infty)E^{X_T}\{f(X_0, X_1, \dots)\}] = E\{A1(T < \infty)f(X_T, X_{T+1}, \dots)\}$$

We evaluate these by breaking up the events into the various possible values of T :

$$\begin{aligned} A1(T < \infty)E^{X_T}\{f(X_0, X_1, \dots)\} &= \sum_k A1(T = k)E^{X_T}\{f(X_0, X_1, \dots)\} \\ &= \sum_k A1(T = k)E^{X_k}\{f(X_0, X_1, \dots)\} \end{aligned} \quad (5.8)$$

On the other hand

$$\begin{aligned} A1(T < \infty)f(X_T, X_{T+1}, \dots) &= \sum_k A1(T = k)f(X_T, X_{T+1}, \dots) \\ &= \sum_k A1(T = k)f(X_k, X_{k+1}, \dots) \end{aligned}$$

But

$$\begin{aligned} E\{A1(T = k)f(X_k, X_{k+1}, \dots)\} &= E[E\{A1(T = k)f(X_k, X_{k+1}, \dots)|\mathcal{F}_k\}] \\ &= E[A1(T = k)E\{f(X_k, X_{k+1}, \dots)|\mathcal{F}_k\}] \\ &= E[A1(T = k)E^{X_k}\{f(X_0, X_1, \dots)\}] \end{aligned}$$

by (5.3). But this is exactly the expectation of the k th term in (5.8). This proves the Strong Markov Property.

5.3 Initial Distributions

Up to now I have hidden the nature of the unconditional expected values. The Markov property (5.1) specifies only conditional probabilities. There is no way to deduce the marginal distributions. Instead, for every possible distribution π on S and each possible transition matrix \mathbf{P} there is a stochastic process X_0, X_1, \dots with

$$P(X_0 = k) = \pi_k$$

and which is a Markov Chain with transition matrix \mathbf{P} .

In our proof of the strong Markov property we used unconditional expectations. Part of the point of the calculation is that it works for any initial distribution. Notation: if π is an initial distribution (a probability measure on S) then E^π denotes expected values for chains with initial distribution π ; we use \mathbf{P}^π for the corresponding probability measure.

Summary of easily verified facts:

- For any sequence of states i_0, \dots, i_k

$$P(X_0 = i_0, \dots, X_k = i_k) = \pi_{i_0} \mathbf{P}_{i_0 i_1} \cdots \mathbf{P}_{i_{k-1} i_k}$$

- For any event A :

$$\mathbf{P}^\pi(A) = \sum_k \pi_k \mathbf{P}^k(A)$$

- For any bounded random variable $Y = f(X_0, \dots)$

$$E^\pi(Y) = \sum_k \pi_k E^k(A)$$

5.4 Recurrence and Transience

Now consider a transient state k , that is, a state for which

$$f_k = P^k(T_k < \infty) < 1$$

Note that $T_k = \min\{n > 0 : X_n = k\}$ is a stopping time. Let N_k be the number of visits to state k . That is

$$N_k = \sum_{n=0}^{\infty} 1(X_n = k)$$

Notice that if we define the function

$$f(x_0, x_1, \dots) = \sum_{n=0}^{\infty} 1(x_n = k)$$

then

$$N_k = f(X_0, X_1, \dots)$$

Notice, also, that on the event $T_k < \infty$

$$N_k = 1 + f(X_{T_k}, X_{T_k+1}, \dots)$$

and on the event $T_k = \infty$ we have

$$N_k = 1$$

In short:

$$N_k = 1 + f(X_{T_k}, X_{T_k+1}, \dots)1(T_k < \infty)$$

Hence

$$\begin{aligned}
\mathbf{P}^k(N_k = r) &= \mathbf{E}^k \{P(N_k = r | \text{cal}F_T)\} \\
&= \mathbf{E}^k \{\mathbf{P}(1 + f(X_{T_k}, X_{T_k+1}, \dots)1(T_k < \infty) = r | \text{cal}F_T)\} \\
&= \mathbf{E}^k [1(T_k < \infty)P^{X_{T_k}} \{f(X_0, X_1, \dots) = r - 1\}] \\
&= \mathbf{E}^k \{1(T_k < \infty)P^k(N_k = r - 1)\} \\
&= \mathbf{E}^k \{1(T_k < \infty)\} P^k(N_k = r - 1) \\
&= f_k P^k(N_k = r - 1)
\end{aligned}$$

It is easily verified by induction, then, that

$$\mathbf{P}^k(N_k = r) = f_k^{r-1} P^k(N_k = 1)$$

But $N_k = 1$ if and only if $T_k = \infty$ so

$$\mathbf{P}^k(N_k = r) = f_k^{r-1} (1 - f_k)$$

which shows that N_k has (if the chain starts from state k) a Geometric distribution with mean $1/(1 - f_k)$. Notice that the argument also shows that if $f_k = 1$ then

$$P^k(N_k = 1) = P^k(N_k = 2) = \dots$$

which can only happen if all these probabilities are 0. Thus

$$P(N_k = \infty) = 1$$

if $f_k = 1$.

Now

$$N_k = \sum_{n=0}^{\infty} 1(X_n = k)$$

Taking expectations gives

$$\mathbf{E}^k(N_k) = \sum_{n=0}^{\infty} (\mathbf{P}^n)_{kk}$$

We therefore learn:

State k is transient if and only if

$$\sum_{n=0}^{\infty} (\mathbf{P}^n)_{kk} < \infty$$

and that this sum is $1/(1 - f_k)$.

Proposition 2 *Recurrence (or transience) is a class property. That is, if i and j are in the same communicating class then i is recurrent (respectively transient) if and only if j is recurrent (respectively transient).*

Proof: Suppose i is recurrent and $i \leftrightarrow j$. There are integers m and n such that

$$(\mathbf{P}^m)_{ji} > 0 \quad \text{and} \quad (\mathbf{P}^n)_{ij} > 0$$

Then

$$\begin{aligned} \sum_k (\mathbf{P}^k)_{jj} &\geq \sum_{k \geq 0} (\mathbf{P}^{m+k+n})_{jj} \\ &\geq \sum_{k \geq 0} (\mathbf{P}^m)_{ji} (\mathbf{P}^k)_{ii} (\mathbf{P}^n)_{ij} \\ &= (\mathbf{P}^m)_{ji} \left\{ \sum_{k \geq 0} (\mathbf{P}^k)_{ii} \right\} (\mathbf{P}^n)_{ij} \end{aligned}$$

The middle term is infinite and the two outside terms positive so

$$\sum_k (\mathbf{P}^k)_{jj} = \infty$$

which shows j is recurrent. •

In a finite state space chain there is at least one recurrent state. If all states were transient we would have for each k $P(N_k < \infty) = 1$. This would mean $P(\forall k . N_k < \infty) = 1$. But for any ω there must be at least one k for which $N_k = \infty$ (the total of a finite list of finite numbers is finite).

For an infinite state space chain this argument does not work. The chain X_n satisfying $X_{n+1} = X_n + 1$ on the integers has all states transient.

A more interesting example is this. Toss a coin repeatedly. Let X_n be X_0 plus the number of heads minus the number of tails in the first n tosses. Let p denote the probability of heads on an individual trial.

Notice that $X_n - X_0$ is a sum of n iid random variables Y_i where $P(Y_i = 1) = p$ and $P(Y_i = -1) = 1 - p$. The Strong Law of Large Numbers shows X_n/n converges almost surely to $2p - 1$. If $p \neq 1/2$ this is not 0. In order for X_n/n to have a positive limit we must have $X_n \rightarrow \infty$ almost surely so all states are visited only finitely many times. That is, all states are transient. Similarly for $p < 1/2$ $X_n \rightarrow -\infty$ almost surely and all states are transient.

Now look at $p = 1/2$. The law of large numbers argument no longer shows anything. I will show that all states are recurrent.

Proof: We show $\sum_n (\mathbf{P}^n)_{00}$ and show the sum is infinite. If n is odd then $(p_n)_{00} = 0$ so we evaluate

$$\sum_m (\mathbf{P}^{2m})_{00}$$

Now

$$(\mathbf{P}^{2m})_{00} = \binom{2m}{m} 2^{-2m}$$

According to Stirling's approximation

$$\lim_{m \rightarrow \infty} \frac{m!}{m^{m+1/2} e^{-m} \sqrt{2\pi}} = 1$$

Hence

$$\lim_{m \rightarrow \infty} \sqrt{m}(\mathbf{P}^{2m})_00 = \frac{1}{\sqrt{\pi}}$$

Since

$$\sum \frac{1}{\sqrt{m}} = \infty$$

we are done.

5.5 Mean return times

There is a substantial difference between the recurrent states in a finite state space Markov Chain and the recurrent states in our last example. The first step in examining the difference is to compute expected times to return. Again, for $x \in S$ let T_x denote the hitting time for x . Consider a recurrent state x . Since the chain, once having entered the class containing x cannot leave it we will simply assume that there is only one communicating class for the chain. We now derive a set of linear equations between the various expected values of different T_x . Note that each T_x is a certain function f_x applied to X_1, \dots . Setting $\mu_{ij} = E^i(T_j)$ we find

$$\mu_{ij} = \sum_k E^i(T_j 1(X_1 = k))$$

Note that if $X_1 = x$ then $T_x = 1$ so

$$E^i(T_j 1(X_1 = j)) = \mathbf{P}_{ij}$$

For $k \neq j$

$$T_x = 1 + f_x(X_2, X_3, \dots)$$

and, by conditioning on $X_1 = k$ we find

$$E^i(T_j 1(X_1 = k)) = \mathbf{P}_{ik} \{1 + E^k(T_j)\}$$

This gives

$$\mu_{ij} = 1 + \sum_{k \neq j} \mathbf{P}_{ik} \mu_{kj} \tag{5.9}$$

Technically, I should check that the expectations in (5.9) are finite. All the random variables involved are non-negative, however, and the equation actually makes sense even if some terms are infinite. (To prove this you actually study

$$T_{x,n} = \min(T_x, n)$$

deriving an identity for a fixed n , letting $n \rightarrow \infty$ and applying the monotone convergence theorem.)

Here is a simple example:

$$\mathbf{P} = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

The identity (5.9) becomes

$$\begin{aligned}\mu_{1,1} &= 1 + \frac{1}{2}\mu_{2,1} + \frac{1}{2}\mu_{3,1} \\ \mu_{1,2} &= 1 + \frac{1}{2}\mu_{3,1} \\ \mu_{1,3} &= 1 + \frac{1}{2}\mu_{2,1} \\ \mu_{2,1} &= 1 + \frac{1}{2}\mu_{3,1} \\ \mu_{2,2} &= 1 + \frac{1}{2}\mu_{1,2} + \frac{1}{2}\mu_{3,2} \\ \mu_{2,3} &= 1 + \frac{1}{2}\mu_{1,3} \\ \mu_{3,1} &= 1 + \frac{1}{2}\mu_{2,1} \\ \mu_{3,2} &= 1 + \frac{1}{2}\mu_{1,2} \\ \mu_{3,3} &= 1 + \frac{1}{2}\mu_{1,3} + \frac{1}{2}\mu_{2,3}\end{aligned}$$

The seventh and fourth together give

$$\mu_{2,1} = \mu_{3,1}$$

Similar calculations lead to $\mu_{ii} = 3$ and, for $i \neq j$ $\mu_{i,j} = 2$.

Example: The coin tossing Markov Chain with $p = 1/2$ shows that the situation can be different when S is infinite. The equations above become:

$$\begin{aligned}m_{0,0} &= 1 + \frac{1}{2}m_{1,0} + \frac{1}{2}m_{-1,0} \\ m_{1,0} &= 1 + \frac{1}{2}m_{2,0}\end{aligned}$$

and many more.

Here are some observations:

- You have to go through 1 to get to 0 from 2 so

$$m_{2,0} = m_{2,1} + m_{1,0}.$$

- The chain which counts tails minus heads has exactly the same distribution as X_n by symmetry. So:

$$m_{1,0} = m_{-1,0}$$

- The transition probabilities are **spatially homogeneous**, meaning \mathbf{P}_{ij} is a function of $j - i$ only. Hence

$$m_{2,1} = m_{1,0}$$

Conclusion:

$$\begin{aligned} m_{0,0} &= 1 + m_{1,0} \\ &= 1 + 1 + \frac{1}{2}m_{2,0} \\ &= 2 + m_{1,0} \end{aligned}$$

Notice that there are **no** finite solutions!

Summary of the situation:

- Every state is recurrent.
- All the expected hitting times m_{ij} are infinite.
- All entries \mathbf{P}_{ij}^n converge to 0.

Jargon: The states in this chain are **null recurrent**.

Definition: A recurrent state x is positive recurrent if $E^x(T_x) < \infty$ otherwise x is null recurrent.

5.6 The ergodic theorem

Consider a finite state space chain. If x is a vector then the i th entry in $\mathbf{P}x$ is

$$\sum_j \mathbf{P}_{ij}x_j$$

This is, since the rows of \mathbf{P} are probability vectors, a weighted average of the entries in x . If the weights are strictly between 0 and 1 and the largest and smallest entries in x are not the same then $\sum_j \mathbf{P}_{ij}x_j$ is strictly between the largest and smallest entries in x . In fact

$$\sum_j \mathbf{P}_{ij}x_j - \min(x_k) \geq \min_j \{p_{ij}\} (\max\{x_k\} - \min\{x_k\})$$

and

$$\max\{x_j\} - \sum_j \mathbf{P}_{ij}x_j \geq \min_j \{p_{ij}\} (\max\{x_k\} - \min\{x_k\})$$

Now consider what happens when we multiply \mathbf{P}^r by \mathbf{P}^m . The ij th entry in \mathbf{P}^{r+m} is a weighted average of the j th column of \mathbf{P}^m and so, if all the entries in row i of \mathbf{P}^r are positive and the j th column of \mathbf{P}^m is not constant, the i th entry in the j th column of \mathbf{P}^{r+m} must be strictly between the minimum and maximum entries of the j th column of \mathbf{P}^m . In fact, fix a j and let \bar{x}_m be the maximum entry in column j of \mathbf{P}^m and \underline{x}_m the minimum entry. Suppose all entries of \mathbf{P}^r are positive. Let $\delta > 0$ be the smallest entry in \mathbf{P}^r . Our argument above shows that

$$\bar{x}_{m+r} \leq \bar{x}_m - \delta(\bar{x}_m - \underline{x}_m)$$

and

$$\underline{x}_{m+r} \geq \underline{x}_m + \delta(\bar{x}_m - \underline{x}_m)$$

Putting these together gives

$$(\bar{x}_{m+r} - \underline{x}_{m+r}) \leq (1 - 2\delta)(\bar{x}_m - \underline{x}_m)$$

In summary the column maximum decreases, the column minimum increases and the gap between the two decreases exponentially along the sequence $m, m+r, m+2r, \dots$. This idea can be used to prove

Proposition 3 *Suppose X_n is a finite state space Markov Chain with stationary transition matrix \mathbf{P} . Assume that there is a power r such that all entries in \mathbf{P}^r are positive. Then for \mathbf{P}^k has all entries positive for all $k \geq r$ and \mathbf{P}^n converges, as $n \rightarrow \infty$ to a matrix \mathbf{P}^∞ . Moreover,*

$$(\mathbf{P}^\infty)_{ij} = \pi_j$$

where π is the unique row vector satisfying

$$\pi = \pi \mathbf{P}$$

whose entries sum to 1.

Proof: First for $k > r$

$$(\mathbf{P}^k)_{ij} = \sum_k (\mathbf{P}^{k-r})_{ik} (\mathbf{P}^r)_{kj}$$

For each i there is a k for which $(\mathbf{P}^{k-r})_{ik} > 0$ and since $(\mathbf{P}^r)_{kj} > 0$ we see $(\mathbf{P}^k)_{ij} > 0$. The argument before the proposition shows that

$$\lim_{j \rightarrow \infty} \mathbf{P}^{m+jk}$$

exists for each m and $k \geq r$. This proves \mathbf{P}^n has a limit which we call \mathbf{P}^∞ . Since \mathbf{P}^{n-1} also converges to \mathbf{P}^∞ we find

$$\mathbf{P}^\infty = \mathbf{P}^\infty \mathbf{P}$$

Hence each row of \mathbf{P}^∞ is a solution of $x\mathbf{P} = x$. The argument before the statement of the proposition shows all rows of \mathbf{P}^∞ are equal. Let π be this common row.

Now if α is any vector whose entries sum to 1 then $\alpha\mathbf{P}^n$ converges to

$$\alpha\mathbf{P}^\infty = \pi$$

If α is any solution of $x = x\mathbf{P}$ we have by induction $\alpha\mathbf{P}^n = \alpha$ so $\alpha\mathbf{P}^\infty = \alpha$ so $\alpha = \pi$. That is exactly one vector whose entries sum to 1 satisfies $x = x\mathbf{P}$. \bullet

The proposition has the condition that there be an r for which all entries in \mathbf{P}^r are positive and that the chain have a finite state space. Consider first the finite state space case. To see that it is possible for \mathbf{P}^n not to have a limit just let

$$\mathbf{P} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

and note that \mathbf{P}^{2n} is the identity while $\mathbf{P}^{2n+1} = \mathbf{P}$. Note, too, that

$$\frac{\mathbf{P}^0 + \cdots + \mathbf{P}^n}{n+1} \rightarrow \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

Consider the equations $\pi = \pi\mathbf{P}$ with $\pi_1 + \pi_2 = 1$. We get

$$\pi_1 = \frac{1}{2}\pi_1 + \frac{1}{2}(1 - \pi_1) = \frac{1}{2}$$

so that the solution to $\pi = \pi\mathbf{P}$ is again unique.

Definition: The period d of a state i is the greatest common divisor of

$$\{n : (\mathbf{P}^n)_{ii} > 0\}$$

Lemma 6 *If $i \leftrightarrow j$ then i and j have the same period.*

Definition: A state is **aperiodic** if its period is 1.

Proof: I do the case $d = 1$. Fix i . Let

$$G = \{k : (\mathbf{P}^k)_{ii} > 0\}$$

If $k_1, k_2 \in G$ then $k_1 + k_2 \in G$.

This (and aperiodic) implies by a number theory argument that there is an r such that $k \geq r$ implies $k \in G$.

Now find m and n so that

$$(\mathbf{P}^m)_{ij} > 0 \text{ and } (\mathbf{P}^n)_{ji} > 0$$

For $k > r + m + n$ we see $(\mathbf{P}^k)_{jj} > 0$ so the gcd of the set of k such that $(\mathbf{P}^k)_{jj} > 0$ is 1. •

The case of period $d > 1$ can be dealt with by considering \mathbf{P}^d .

$$\mathbf{P} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

For this example $\{1, 2, 3\}$ is a class of period 3 states and $\{4, 5\}$ a class of period 2 states.

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

has a single communicating class of period 2.

A chain is **aperiodic** if all its states are aperiodic.

5.7 Hitting Times

Start irreducible recurrent chain X_n in state i . Let T_j be first $n > 0$ such that $X_n = j$. Define

$$m_{ij} = E(T_j | X_0 = i)$$

First step analysis:

$$\begin{aligned} m_{ij} &= 1 \cdot P(X_1 = j | X_0 = i) \\ &\quad + \sum_{k \neq j} (1 + E(T_j | X_0 = k)) P_{ik} \\ &= \sum_j P_{ij} + \sum_{k \neq j} P_{ik} m_{kj} \\ &= 1 + \sum_{k \neq j} P_{ik} m_{kj} \end{aligned}$$

Example

$$\mathbf{P} = \begin{bmatrix} \frac{3}{5} & \frac{2}{5} \\ \frac{1}{5} & \frac{4}{5} \end{bmatrix}$$

The equations are

$$\begin{aligned} m_{11} &= 1 + \frac{2}{5} m_{21} \\ m_{12} &= 1 + \frac{3}{5} m_{12} \\ m_{21} &= 1 + \frac{4}{5} m_{21} \\ m_{22} &= 1 + \frac{1}{5} m_{12} \end{aligned}$$

The second and third equations give immediately

$$\begin{aligned} m_{12} &= \frac{5}{2} \\ m_{21} &= 5 \end{aligned}$$

Then plug in to the others to get

$$\begin{aligned} m_{11} &= 3 \\ m_{22} &= \frac{3}{2} \end{aligned}$$

Notice stationary initial distribution is

$$\left(\frac{1}{m_{11}}, \frac{1}{m_{22}} \right)$$

Consider fraction of time spent in state j :

$$\frac{1(X_0 = j) + \cdots + 1(X_n = j)}{n + 1}$$

Imagine chain starts in chain i ; take expected value.

$$\frac{\sum_{r=1}^n \mathbf{P}_{ij}^r + 1(i = j)}{n + 1}$$

If rows of \mathbf{P} converge to π then fraction converges to π_j ; i.e. limiting fraction of time in state j is π_j .

Heuristic: start chain in i . Expect to return to i every m_{ii} time units. So are in state i about once every m_{ii} time units; i.e. limiting fraction of time in state i is $1/m_{ii}$.

Conclusion: for an irreducible recurrent finite state space Markov chain

$$\pi_i = \frac{1}{m_{ii}}.$$

Real proof: Renewal theorem or variant.

Idea: $S_1 < S_2 < \dots$ are times of visits to i . Segment j :

$$X_{S_{j-1}+1}, \dots, X_{S_j}.$$

Segments are iid by Strong Markov.

Number of visits to i by time S_k is exactly k .

Total elapsed time is $S_k = T_1 + \cdots + T_k$ where T_j are iid.

Fraction of time in state i by time S_k is

$$\frac{k}{S_k} \rightarrow \frac{1}{m_{ii}}$$

by SLLN. So if fraction converges to π_i must have

$$\pi_i = \frac{1}{m_{ii}}.$$

Summary of Theoretical Results:

For an irreducible aperiodic positive recurrent Markov Chain:

1. \mathbf{P}^n converges to a stochastic matrix \mathbf{P}^∞ .
2. Each row of \mathbf{P}^∞ is π the unique stationary initial distribution.
3. The stationary initial distribution is given by

$$\pi_i = 1/m_i$$

where m_i is the mean return time to state i from state i .

If the state space is finite an irreducible chain is positive recurrent.

Ergodic Theorem

Notice slight of hand: I showed

$$\frac{E \{ \sum_{i=0}^n 1(X_i = k) \}}{n} \rightarrow \pi_k$$

but claimed

$$\frac{\sum_{i=0}^n 1(X_i = k)}{n} \rightarrow \pi_k$$

almost surely which is also true. This is a step in the proof of the ergodic theorem. For an irreducible positive recurrent Markov chain and any f on S such that $E^\pi(f(X_0)) < \infty$:

$$\frac{\sum_0^n f(X_i)}{n} \rightarrow \sum \pi_j f(j)$$

almost surely. The limit works in other senses, too. You also get

$$\frac{\sum_0^n f(X_i, \dots, X_{i+k})}{n} \rightarrow E^\pi \{f(X_0, \dots, X_k)\}$$

E.g. fraction of transitions from i to j goes to

$$\pi_i \mathbf{P}^{ij}$$

For an irreducible positive recurrent chain of period d :

1. \mathbf{P}^d has d communicating classes each of which forms an irreducible aperiodic positive recurrent chain.
2. $(\mathbf{P}^{n+1} + \dots + \mathbf{P}^{n+d})/d$ has a limit \mathbf{P}^∞ .
3. Each row of \mathbf{P}^∞ is π the unique stationary initial distribution.
4. Stationary initial distribution places probability $1/d$ on each of the communicating classes in 1.

For an irreducible null recurrent chain:

1. \mathbf{P}^n converges to 0 (pointwise).
2. there is no stationary initial distribution.

For an irreducible transient chain:

1. \mathbf{P}^n converges to 0 (pointwise).
2. there is no stationary initial distribution.

For a chain with more than 1 communicating class:

1. If \mathcal{C} is a recurrent class the submatrix $\mathbf{P}_{\mathcal{C}}$ of \mathbf{P} made by picking out rows i and columns j for which $i, j \in \mathcal{C}$ is a stochastic matrix. The corresponding entries in \mathbf{P}^n are just $(\mathbf{P}_{\mathcal{C}})^n$ so you can apply the conclusions above.
2. For any transient or null recurrent class the corresponding columns in \mathbf{P}^n converge to 0.
3. If there are multiple positive recurrent communicating classes then the stationary initial distribution is not unique.

Poisson Processes

Particles arriving over time at a particle detector. Several ways to describe most common model.

Approach 1: numbers of particles arriving in an interval has Poisson distribution, mean proportional to length of interval, numbers in several non-overlapping intervals independent.

For $s < t$, denote number of arrivals in $(s, t]$ by $N(s, t)$. Jargon: $N(A)$ = number of points in A is a **counting process**. Model is

1. $N(s, t)$ has a $\text{Poisson}(\lambda(t - s))$ distribution.
2. For $0 \leq s_1 < t_1 \leq s_2 < t_2 \cdots \leq s_k < t_k$ the variables $N(s_i, t_i); i = 1, \dots, k$ are independent.

Approach 2:

Let $0 < S_1 < S_2 < \cdots$ be the times at which the particles arrive.

Let $T_i = S_i - S_{i-1}$ with $S_0 = 0$ by convention. T_i are called **interarrival times**.

Then T_1, T_2, \dots are independent Exponential random variables with mean $1/\lambda$.

Note $P(T_i > x) = e^{-\lambda x}$ is called **survival function** of T_i .

Approaches are equivalent. Both are deductions of a model based on **local** behaviour of process.

Approach 3: Assume:

1. given all the points in $[0, t]$ the probability of 1 point in the interval $(t, t + h]$ is of the form

$$\lambda h + o(h)$$

2. given all the points in $[0, t]$ the probability of 2 or more points in interval $(t, t + h]$ is of the form

$$o(h)$$

Notation: given functions f and g we write

$$f(h) = g(h) + o(h)$$

provided

$$\lim_{h \rightarrow 0} \frac{f(h) - g(h)}{h} = 0$$

[Aside: if there is a constant M such that

$$\limsup_{h \rightarrow 0} \left| \frac{f(h) - g(h)}{h} \right| \leq M$$

we say

$$f(h) = g(h) + O(h)$$

Notation due to Landau. Another form is

$$f(h) = g(h) + O(h)$$

means there is $\delta > 0$ and M s.t. for all $|h| < \delta$

$$|f(h) - g(h)| \leq Mh$$

Idea: $o(h)$ is tiny compared to h while $O(h)$ is (very) roughly the same size as h .]

Generalizations:

1. First (Poisson) model generalizes to $N(s, t]$ having a Poisson distribution with parameter $\Lambda(t) - \Lambda(s)$ for some non-decreasing non-negative function Λ (called **cumulative intensity**). Result called **inhomogeneous** Poisson process.
2. Exponential interarrival model generalizes to independent non-exponential interarrival times. Result is **renewal process** or **semi-Markov** process.
3. Infinitesimal probability model generalizes to other infinitesimal jump rates. Model specifies **infinitesimal generator**. Yields other **continuous time Markov Chains**.

However, all 3 approaches to Poisson process are equivalent. I show: 3 implies 1, 1 implies 2 and 2 implies 3. First explain o , O .

Model 3 implies 1: Fix t , define $f_t(s)$ to be conditional probability of 0 points in $(t, t + s]$ given value of process on $[0, t]$.

Derive differential equation for f . Given process on $[0, t]$ and 0 points in $(t, t + s]$ probability of no points in $(t, t + s + h]$ is

$$f_{t+s}(h) = 1 - \lambda h + o(h)$$

Given the process on $[0, t]$ the probability of no points in $(t, t + s]$ is $f_t(s)$. Using $P(AB|C) = P(A|BC)P(B|C)$ gives

$$\begin{aligned} f_t(s + h) &= f_t(s)f_{t+s}(h) \\ &= f_t(s)(1 - \lambda h + o(h)) \end{aligned}$$

Now rearrange, divide by h to get

$$\frac{f_t(s+h) - f_t(s)}{h} = \lambda f_t(s) + \frac{o(h)}{h}$$

Let $h \rightarrow 0$ and find

$$\frac{\partial f_t(s)}{\partial s} = -\lambda f_t(s)$$

Differential equation has solution

$$f_t(s) = f_t(0) \exp(-\lambda s) = \exp(-\lambda s).$$

Things to notice:

- $f_t(s) = e^{-\lambda s}$ is survival function of exponential rv..
- We had suppressed dependence of $f_t(s)$ on $N(u); 0 \leq u \leq t$ but solution does not depend on condition.
- So: the event of getting 0 points in $(t, t+s]$ is independent of $N(u); 0 \leq u \leq t$.
- We used: $f_t(s)o(h) = o(h)$. Other rules:

$$\begin{aligned} o(h) + o(h) &= o(h) \\ O(h) + O(h) &= O(h) \\ O(h) + o(h) &= O(h) \\ o(h^r)O(h^s) &= o(h^{r+s}) \\ O(o(h)) &= o(h) \end{aligned}$$

General case:

Notation: $N(t) = N(0, t)$.

$N(t)$ is a non-decreasing function of t . Let

$$P_k(t) = P(N(t) = k).$$

Evaluate $P_k(t+h)$ by conditioning on $N(s); 0 \leq s < t$ and $N(t) = j$.

Given $N(t) = j$ probability that $N(t+h) = k$ is conditional probability of $k-j$ points in $(t, t+h]$.

So, for $j \leq k-2$:

$$P(N(t+h) = k | N(t) = j, N(s), 0 \leq s < t) = o(h).$$

For $j = k-1$ we have

$$P(N(t+h) = k | N(t) = k-1, N(s), 0 \leq s < t) = \lambda h + o(h)$$

For $j = k$ we have

$$P(N(t+h) = k | N(t) = k, N(s), 0 \leq s < t) = 1 - \lambda h + o(h).$$

N is increasing so only consider $j \leq k$.

$$\begin{aligned} P_k(t+h) &= \sum_{j=0}^k P(N(t+h) = k | N(t) = j) P_j(t) \\ &= P_k(t)(1 - \lambda h) + \lambda h P_{k-1}(t) + o(h) \end{aligned}$$

Rearrange, divide by h and let $h \rightarrow 0$ to get

$$P'_k(t) = -\lambda P_k(t) + \lambda P_{k-1}(t)$$

For $k = 0$ the term P_{k-1} is dropped and

$$P'_0(t) = \lambda P_0(t)$$

Using $P_0(0) = 1$ we get

$$P_0(t) = e^{-\lambda t}$$

Put this into the equation for $k = 1$ to get

$$P'_1(t) = -\lambda P_1(t) + \lambda e^{-\lambda t}$$

Multiply by $e^{\lambda t}$ to see

$$(e^{\lambda t} P_1(t))' = \lambda$$

With $P_1(0) = 0$ we get

$$P_1(t) = \lambda t e^{-\lambda t}$$

For general k we have $P_k(0) = 0$ and

$$(e^{\lambda t} P_k(t))' = \lambda e^{\lambda t} P_{k-1}(t)$$

Check by induction that

$$e^{\lambda t} P_k(t) = (\lambda t)^k / k!$$

Hence: $N(t)$ has Poisson(λt) distribution.

Similar ideas permit proof of

$$P(N(s, t) = k | N(u); 0 \leq u \leq s) = \frac{\{\lambda(t-s)\}^k e^{-\lambda}}{k!}$$

From which (by induction) we can prove that N has independent Poisson increments.

Exponential Interarrival Times

If N is a Poisson Process we define T_1, T_2, \dots to be the times between 0 and the first point, the first point and the second and so on.

Fact: T_1, T_2, \dots are iid exponential rvs with mean $1/\lambda$.

We already did T_1 rigorously. The event $T > t$ is exactly the event $N(t) = 0$. So

$$P(T > t) = \exp(-\lambda t)$$

which is the survival function of an exponential rv.

I do case of T_1, T_2 . Let t_1, t_2 be two positive numbers and $s_1 = t_1, s_2 = t_1 + t_2$. The event

$$\{t_1 < T_1 \leq t_1 + \delta_1\} \cap \{t_2 < T_2 \leq t_2 + \delta_2\}.$$

This is almost the same as the intersection of the four events:

$$\begin{aligned} N(0, t_1] &= 0 \\ N(t_1, t_1 + \delta_1] &= 1 \\ N(t_1 + \delta_1, t_1 + \delta_1 + t_2] &= 0 \\ N(s_2 + \delta_1, s_2 + \delta_1 + \delta_2] &= 1 \end{aligned}$$

which has probability

$$e^{-\lambda t_1} \times \lambda \delta_1 e^{-\lambda \delta_1} \times e^{-\lambda t_2} \times \lambda \delta_2 e^{-\lambda \delta_2}$$

Divide by $\delta_1 \delta_2$ and let δ_1 and δ_2 go to 0 to get joint density of T_1, T_2 is

$$e^{-\lambda t_1} e^{-\lambda t_2}$$

which is the joint density of two independent exponential variates.

Here is a more rigorous argument. It has the following steps:

- Find joint density of S_1, \dots, S_k .
- Use **change of variables** to find joint density of T_1, \dots, T_k .

First step: Compute

$$P(0 < S_1 \leq s_1 < S_2 \leq s_2 \cdots < S_k \leq s_k)$$

This is just the event of exactly 1 point in each interval $(s_{i-1}, s_i]$ for $i = 1, \dots, k-1$ ($s_0 = 0$) and at least one point in $(s_{k-1}, s_k]$ which has probability

$$\prod_1^{k-1} \{\lambda(s_i - s_{i-1}) e^{-\lambda(s_i - s_{i-1})}\} (1 - e^{-\lambda(s_k - s_{k-1})})$$

Second step: write this in terms of joint cdf of S_1, \dots, S_k . I do $k = 2$:

$$P(0 < S_1 \leq s_1 < S_2 \leq s_2) = F_{S_1, S_2}(s_1, s_2) - F_{S_1, S_2}(s_1, s_1)$$

Notice tacit assumption $s_1 < s_2$.

Differentiate twice, that is, take

$$\frac{\partial^2}{\partial s_1 \partial s_2}$$

to get

$$f_{S_1, S_2}(s_1, s_2) = \frac{\partial^2}{\partial s_1 \partial s_2} \lambda s_1 e^{-\lambda s_1} (1 - e^{-\lambda(s_2 - s_1)})$$

Simplify to

$$\lambda^2 e^{-\lambda s_2}$$

Recall tacit assumption to get

$$f_{S_1, S_2}(s_1, s_2) = \lambda^2 e^{-\lambda s_2} \mathbf{1}(0 < s_1 < s_2)$$

That completes the first part.

Now compute the joint cdf of T_1, T_2 by

$$F_{T_1, T_2}(t_1, t_2) = P(S_1 < t_1, S_2 - S_1 < t_2)$$

This is

$$\begin{aligned} P(S_1 < t_1, S_2 - S_1 < t_2) &= \int_0^{t_1} \int_{s_1}^{s_1+t_2} \lambda^2 e^{-\lambda s_2} ds_2 ds_1 \\ &= \lambda \int_0^{t_1} (e^{-\lambda s_1} - e^{-\lambda(s_1+t_2)}) ds_1 \\ &= 1 - e^{-\lambda t_1} - e^{-\lambda t_2} + e^{-\lambda(t_1+t_2)} \end{aligned}$$

Differentiate twice to get

$$f_{T_1, T_2}(t_1, t_2) = \lambda e^{-\lambda t_1} \lambda e^{-\lambda t_2}$$

which is the joint density of two independent exponential random variables.

Summary so far:

Have shown:

Instantaneous rates model implies independent Poisson increments model implies independent exponential interarrivals.

Next: show independent exponential interarrivals implies the instantaneous rates model.

Suppose T_1, \dots iid exponential rvs with means $1/\lambda$. Define N_t by $N_t = k$ if and only if

$$T_1 + \dots + T_k \leq t \leq T_1 + \dots + T_{k+1}$$

Let A be the event $N(s) = n(s); 0 < s \leq t$. We are to show

$$P(N(t, t+h) = 1 | N(t) = k, A) = \lambda h + o(h)$$

and

$$P(N(t, t+h) \geq 2 | N(t) = k, A) = o(h)$$

If $n(s)$ is a possible trajectory consistent with $N(t) = k$ then n has jumps at points

$$t_1, t_1 + t_2, \dots, s_k \equiv t_1 + \dots + t_k < t$$

and at no other points in $(0, t]$.

So given $N(s) = n(s); 0 < s \leq t$ with $n(t) = k$ we are essentially being given

$$T_1 = t_1, \dots, T_k = t_k, T_{k+1} > t - s_k$$

and asked the conditional probability in the first case of the event B given by

$$t - s_k < T_{k+1} \leq t - s_k + h < T_{k+2} + T_{k+1}.$$

Conditioning on T_1, \dots, T_k irrelevant (independence).

$$\begin{aligned} P(N(t, t+h] = 1 | N(t) = k, A) / h \\ &= P(B | T_{k+1} > t - s_k) / h \\ &= \frac{P(B)}{h e^{-\lambda(t-s_k)}} \end{aligned}$$

The numerator may be evaluated by integration:

$$P(B) = \int_{t-s_k}^{t-s_k+h} \int_{t-s_k+h-u_1}^{\infty} \lambda^2 e^{-\lambda(u_1+u_2)} du_2 du_1$$

Let $h \rightarrow 0$ to get the limit

$$P(N(t, t+h] = 1 | N(t) = k, A) / h \rightarrow \lambda$$

as required.

The computation of

$$\lim_{h \rightarrow 0} P(N(t, t+h] \geq 2 | N(t) = k, A) / h$$

is similar.

Properties of exponential rvs

Convolution: If X and Y independent rvs with densities f and g respectively and $Z = X + Y$ then

$$P(Z \leq z) = \int_{-\infty}^{\infty} \int_{-\infty}^{z-x} f(x)g(y)dydx$$

Differentiating with respect to z we get

$$f_Z(z) = \int_{-\infty}^{\infty} f(x)g(z-x)dx$$

This integral is called the **convolution** of densities f and g .

If T_1, \dots, T_n iid Exponential(λ) then $S_n = T_1 + \dots + T_n$ has a Gamma(n, λ) distribution. Density of S_n is

$$f_{S_n}(s) = \lambda(\lambda s)^{n-1} e^{-\lambda s} / (n-1)!$$

for $s > 0$.

Proof:

$$\begin{aligned} P(S_n > s) &= P(N(0, s] < n) \\ &= \sum_{j=0}^{n-1} (\lambda s)^j e^{-\lambda s} / j! \end{aligned}$$

Then

$$\begin{aligned} f_{S_n}(s) &= \frac{d}{ds} P(S_n \leq s) \\ &= \frac{d}{ds} \{1 - P(S_n > s)\} \\ &= -\lambda \sum_{j=1}^{n-1} n-1 \{j(\lambda s)^{j-1} - (\lambda s)^j\} \frac{e^{-\lambda s}}{j!} \\ &\quad + \lambda e^{-\lambda s} \\ &= \lambda e^{-\lambda s} \sum_{j=1}^{n-1} \left\{ \frac{(\lambda s)^j}{j!} - \frac{(\lambda s)^{j-1}}{(j-1)!} \right\} \\ &\quad + \lambda e^{-\lambda s} \end{aligned}$$

This telescopes to

$$f_{S_n}(s) = \lambda(\lambda s)^{n-1} e^{-\lambda s} / (n-1)!$$

Extreme Values: If X_1, \dots, X_n are independent exponential rvs with means $1/\lambda_1, \dots, 1/\lambda_n$ then $Y = \min\{X_1, \dots, X_n\}$ has an exponential distribution with mean

$$\frac{1}{\lambda_1 + \dots + \lambda_n}$$

Proof: :

$$\begin{aligned} P(Y > y) &= P(\forall k X_k > y) \\ &= \prod e^{-\lambda_k y} \\ &= e^{-\sum \lambda_k y} \end{aligned}$$

Memoryless Property: conditional distribution of $X - x$ given $X \geq x$ is exponential if X has an exponential distribution.

Proof: :

$$\begin{aligned}
 P(X - x > y | X \geq x) &= \frac{P(X > x + y, X \geq x)}{P(X > x)} \\
 &= \frac{P(X > x + y)}{P(X \geq x)} \\
 &= \frac{e^{-\lambda(x+y)}}{e^{-\lambda x}} \\
 &= e^{-\lambda y}
 \end{aligned}$$

Hazard Rates

The hazard rate, or instantaneous failure rate for a positive random variable T with density f and cdf F is

$$r(t) = \lim_{\delta \rightarrow 0} \frac{P(t < T \leq t + \delta | T \geq t)}{\delta}$$

This is just

$$r(t) = \frac{f(t)}{1 - F(t)}$$

For an exponential random variable with mean $1/\lambda$ this is

$$h(t) = \frac{\lambda e^{-\lambda t}}{e^{-\lambda t}} = \lambda$$

The exponential distribution has constant failure rate.

Weibull random variables have density

$$f(t|\lambda, \alpha) = \lambda(\lambda t)^{\alpha-1} e^{-(\lambda t)^\alpha}$$

for $t > 0$. The corresponding survival function is

$$1 - F(t) = e^{-(\lambda t)^\alpha}$$

and the hazard rate is

$$r(t) = \lambda(\lambda t)^{\alpha-1}$$

which is increasing for $\alpha > 1$, decreasing for $\alpha < 1$. For $\alpha = 1$ this is the exponential distribution.

Since

$$r(t) = \frac{dF(t)/dt}{1 - F(t)} = -\frac{d \log(1 - F(t))}{dt}$$

we can integrate to find

$$1 - F(t) = \exp\left\{-\int_0^t r(s)ds\right\}$$

so that r determines F and f .

Properties of Poisson Processes

- 1) If N_1 and N_2 are independent Poisson processes with rates λ_1 and λ_2 , respectively, then $N = N_1 + N_2$ is a Poisson processes with rate $\lambda_1 + \lambda_2$.
- 2) Suppose N is a Poisson process with rate λ . Suppose each point is marked with a label, say one of L_1, \dots, L_r , independently of all other occurrences. Suppose p_i is the probability that a given point receives label L_i . Let N_i count the points with label i (so that $N = N_1 + \dots + N_r$). Then N_1, \dots, N_r are independent Poisson processes with rates $p_i\lambda$.
- 3) Suppose U_1, U_2, \dots independent rvs, each uniformly distributed on $[0, T]$. Suppose M is a Poisson(λT) random variable independent of the U 's. Let

$$N(t) = \sum_1^M 1(U_i \leq t)$$

Then N is a Poisson process on $[0, T]$ with rate λ .

- 4) Suppose N is a Poisson process with rate λ . Let $S_1 < S_2 < \dots$ be the times at which points arrive. Given $N(T) = n$, S_1, \dots, S_n have the same distribution as the order statistics of a sample of size n from the uniform distribution on $[0, T]$.
- 5) Given $S_{n+1} = t$, S_1, \dots, S_n have the same distribution as the order statistics of a sample of size n from the uniform distribution on $[0, T]$.

Indications of some proofs:

1) N_1, \dots, N_r independent Poisson processes rates λ_i , $N = \sum N_i$. Let A_h be the event of 2 or more points in N in the time interval $(t, t+h]$, B_h , the event of exactly one point in N in the time interval $(t, t+h]$.

Let A_{ih} and B_{ih} be the corresponding events for N_i .

Let H_t denote the history of the processes up to time t ; we condition on H_t . Technically, H_t is the σ -field generated by

$$\{N_i(s); 0 \leq s \leq t, i = 1, \dots, r\}$$

We are given:

$$P(A_{ih}|H_t) = o(h)$$

and

$$P(B_{ih}|H_t) = \lambda_i h + o(h).$$

Note that

$$A_h \subset \bigcup_{i=1}^r A_{ih} \cup \bigcup_{i \neq j} (B_{ih} \cap B_{jh})$$

Since

$$\begin{aligned} P(B_{ih} \cap B_{jh} | H_t) &= P(B_{ih} | H_t) P(B_{jh} | H_t) \\ &= (\lambda_i h + o(h)) (\lambda_j h + o(h)) \\ &= O(h^2) \\ &= o(h) \end{aligned}$$

and

$$P(A_{ih} | H_t) = o(h)$$

we have checked one of the two infinitesimal conditions for a Poisson process.

Next let C_h be the event of no points in N in the time interval $(t, t+h]$ and C_{ih} the same for N_i . Then

$$\begin{aligned} P(C_h | H_t) &= P(\cap C_{ih} | H_t) \\ &= \prod P(C_{ih} | H_t) \\ &= \prod (1 - \lambda_i h + o(h)) \\ &= 1 - (\sum \lambda_i) h + o(h) \end{aligned}$$

shows

$$\begin{aligned} P(B_h | H_t) &= 1 - P(C_h | H_t) - P(A_h | H_t) \\ &= (\sum \lambda_i) h + o(h) \end{aligned}$$

Hence N is a Poisson process with rate $\sum \lambda_i$.

2) The infinitesimal approach used for 1 can do part of this. See text for rest. Events defined as in **1)**: The event B_{ih} that there is one point in N_i in $(t, t+h]$ is the event, B_h that there is exactly one point in any of the r processes together with a subset of A_h where there are two or more points in N in $(t, t+h]$ but exactly one is labeled i . Since $P(A_h | H_t) = o(h)$

$$\begin{aligned} P(B_{ih} | H_t) &= p_i P(B_h | H_t) + o(h) \\ &= p_i (\lambda h + o(h)) + o(h) \\ &= p_i \lambda h + o(h) \end{aligned}$$

Similarly, A_{ih} is a subset of A_h so

$$P(A_{ih} | H_t) = o(h)$$

This shows each N_i is Poisson with rate λp_i . To get independence requires more work; see the homework for an easier algebraic method.

3) Fix $s < t$. Let $N(s, t)$ be the number of points in $(s, t]$. Given $N = n$ the conditional distribution of $N(s, t)$ is Binomial(n, p) with $p = (t - s)/T$. So

$$\begin{aligned}
P(N(s, t) = k) &= \sum_{n=k}^{\infty} P(N(s, t) = k, N = n) \\
&= \sum_{n=k}^{\infty} P(N(s, t) = k | N = n) P(N = n) \\
&= \sum_{n=k}^{\infty} \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} \frac{(\lambda T)^n}{n!} e^{-\lambda T} \\
&= \frac{e^{-\lambda T}}{k!} (\lambda T p)^k \sum_{n=k}^{\infty} \frac{(1-p)^{n-k} (\lambda T)^{n-k}}{(n-k)!} \\
&= \frac{e^{-\lambda T}}{k!} (\lambda T p)^k \sum_{m=0}^{\infty} \frac{(1-p)^m (\lambda T)^m}{m!} \\
&= \frac{e^{-\lambda T}}{k!} (\lambda T p)^k e^{-\lambda T(1-p)} \\
&= \frac{e^{-\lambda(t-s)} (\lambda(t-s))^k}{k!}
\end{aligned}$$

4): Fix s_i, h_i for $i = 1, \dots, n$ such that

$$0 < s_1 < s_1 + h_1 < s_2 < \dots < s_n < s_n + h_n < T$$

Given $N(T) = n$ we compute the probability of the event

$$A = \bigcap_{i=1}^n \{s_i < S_i < s_i + h_i\}$$

Intersection of A , $N(T) = n$ is ($s_0 = h_0 = 0$):

$$B \equiv \left[\bigcap_{i=1}^n \{N(s_{i-1} + h_{i-1}, s_i] = 0, N(s_i, s_i + h_i] = 1\} \right] \cap \{N(s_n + h_n, T] = 0\}$$

whose probability is

$$\left(\prod \lambda h_i \right) e^{-\lambda T}$$

So

$$\begin{aligned}
P(A | N(t) = n) &= \frac{P(A, N(T) = n)}{P(N(T) = n)} \\
&= \frac{\lambda^n e^{-\lambda T} \prod h_i}{(\lambda T)^n e^{-\lambda T} / n!} \\
&= \frac{n! \prod h_i}{T^n}
\end{aligned}$$

Divide by $\prod h_i$ and let all h_i go to 0 to get joint density of S_1, \dots, S_n is

$$\frac{n!}{T^n} 1(0 < s_1 < \dots < s_n < T)$$

which is the density of order statistics from a Uniform $[0, T]$ sample of size n .

5) Replace the event $S_{n+1} = T$ with $T < S_{n+1} < T + h$. With A as before we want

$$P(A|T < S_{n+1} < T + h) = \frac{P(B, N(T, T + h] \geq 1)}{P(T < S_{n+1} < T + h)}$$

Note that B is independent of $\{N(T, T + h] \geq 1\}$ and that we have already found the limit

$$\frac{P(B)}{\prod h_i} \rightarrow \lambda^n e^{-\lambda T}$$

We are left to compute the limit of

$$\frac{P(N(T, T + h] \geq 1)}{P(T < S_{n+1} < T + h)}$$

The denominator is

$$\sum_{k=0}^n P(N(0, T] = k, N(T, T + h] = n + 1 - k) = P(N(0, T] = n) \lambda h + o(h)$$

Thus

$$\begin{aligned} \frac{P(N(T, T + h] \geq 1)}{P(T < S_{n+1} < T + h)} &= \frac{\lambda h + o(h)}{\frac{(\lambda T)^n}{n!} e^{-\lambda T} \lambda h + o(h)} \\ &\rightarrow \frac{n!}{(\lambda T)^n e^{-\lambda T}} \end{aligned}$$

This gives the conditional density of S_1, \dots, S_n given $S_{n+1} = T$ as in 4).

Inhomogeneous Poisson Processes

The idea of hazard rate can be used to extend the notion of Poisson Process. Suppose $\lambda(t) \geq 0$ is a function of t . Suppose N is a counting process such that

$$P(N(t + h) = k + 1 | N(t) = k, H_t) = \lambda(t)h + o(h)$$

and

$$P(N(t + h) \geq k + 2 | N(t) = k, H_t) = o(h)$$

Then N has independent increments and $N(t + s) - N(t)$ has a Poisson distribution with mean

$$\int_t^{t+s} \lambda(u) du$$

If we put

$$\Lambda(t) = \int_0^t \lambda(u) du$$

then mean of $N(t+s) - N(t)$ is $\Lambda(t+s) - \Lambda(t)$.

Jargon: λ is the **intensity** or **instantaneous intensity** and Λ the **cumulative intensity**.

Can use the model with Λ any non-decreasing right continuous function, possibly without a derivative. This allows ties.

Compound Poisson Processes

Imagine insurance claims arise at times of a Poisson process, $N(t)$, (more likely for an inhomogeneous process).

Let Y_i be the value of the i th claim associated with the point whose time is S_i .

Assume that the Y 's are independent of each other and of N .

Let

$$\mu = E(Y_i) \text{ and } \sigma^2 = \text{var}(Y_i)$$

Let

$$X(t) = \sum_{i=1}^{N(t)} Y_i$$

be the total claim up to time t . We call X a compound Poisson Process.

Useful properties:

$$\begin{aligned} E\{X(t)|N(t)\} &= N(t)\mu \\ \text{Var}\{X(t)|N(t)\} &= N(t)\sigma^2 \\ E\{X(t)\} &= E\{N(t)\} \\ &= \lambda t \\ \text{Var}\{X(t)\} &= \text{Var}[E\{X(t)|N(t)\}] \\ &\quad + E[\text{Var}\{X(t)|N(t)\}] \\ &= \lambda t\mu^2 + \lambda t\sigma^2 \end{aligned}$$

Space Time Poisson Processes

Suppose at each time S_i of a Poisson Process, N , we have rv Y_i with the Y_i iid and independent of the Poisson process. Let M be the counting process on $[0, \infty) \times \mathcal{Y}$ (where \mathcal{Y} is the range space of the Y s) defined by

$$M(A) = \#\{(S_i, Y_i) \in A\}$$

Then M is an inhomogeneous Poisson process with mean function μ a measure extending

$$\mu([a, b] \times C) = \lambda(b-a)P(Y \in C)$$

This means that each $M(A)$ has a Poisson distribution with mean $\mu(A)$ and if A_1, \dots, A_r are disjoint then $M(A_1), \dots, M(A_r)$ are independent. The proof in general is a monotone class argument. The first step is: if $(a_i, b_i), i = 1, \dots, r$ are disjoint intervals and C_1, \dots, C_s disjoint subsets of \mathcal{Y} then the rs rvs $M((a_i, b_i) \times C_j)$ are independent Poisson random variables. See the homework for proof of a special case.

Chapter 6

Continuous Time Markov Chains

Consider a population of single celled organisms in a stable environment. In a short time interval each organism might be regarded as having some probability of dividing to produce two organisms and some other probability of dying. We might suppose:

- Different organisms behave independently.
- If the time interval has length say h then the probability of division is λh plus a quantity small compared to h .
- The probability of death is say μh plus a quantity small compared to h .
- The probability that an organism divides twice (or divides once and dies) in the interval of length h is $o(h)$.

Notice the tacit assumption that the constants of proportionality do not depend on time (that is our interpretation of “stable environment”). Notice too that we have taken the constants not to depend on which organism we are talking about. We are really assuming that the organisms are all similar and live in similar environments.

Let $Y(t)$ be the total population at time t . Let \mathcal{H}_t be the history of the process up to time t . (We generally take

$$\mathcal{H}_t = \sigma\{Y(s); 0 \leq s \leq t\}$$

but, in fact, the general definition of a history is any family of σ -fields indexed by t satisfying:

- $s < t$ implies $\mathcal{H}_s \subset \mathcal{H}_t$.
- $Y(t)$ is a \mathcal{H}_t measurable random variable.
- $\mathcal{H}_t \cap (s > t)\mathcal{H}_s$.

The last assumption is a technical detail we will ignore from now on.

Condition on the event $Y(t) = n$. Then the probability of two or more divisions (either more than one division by a single organism or two or more organisms dividing) is $o(h)$ by

our assumptions. Similarly the probability of both a division and a death or of two or more deaths is $o(h)$. We deduce:

$$\begin{aligned}P(Y(t_h) = n + 1 | Y(t) = n, \mathcal{H}_t) &= n\lambda h + o(h) \\P(Y(t_h) = n - 1 | Y(t) = n, \mathcal{H}_t) &= n\mu h + o(h) \\P(Y(t_h) = n | Y(t) = n, \mathcal{H}_t) &= 1 - n(\lambda + \mu)h + o(h) \\P(Y(t_h) \in \{n - 1, n, n + 1\} | Y(t) = n, \mathcal{H}_t) &= o(h)\end{aligned}$$