

|   |   |
|---|---|
| <p>APTS-ASP 1</p> <h2 style="text-align: center;">APTS Applied Stochastic Processes<br/>Preliminary material</h2> <p style="text-align: center;">Nicholas Georgiou<sup>1</sup> &amp; Matt Roberts<sup>2</sup><br/> nicholas.georgiou@durham.ac.uk and mattiroberts@gmail.com</p> <p style="text-align: center;">(Notes originally produced by Wilfrid Kendall; some material due to<br/> Stephen Connor, Christina Goldschmidt and Amanda Turner)</p> <p style="text-align: center;"><sup>1</sup>Department of Mathematical Sciences, Durham University<br/> <sup>2</sup>Probability Laboratory, University of Bath</p> <p style="text-align: center;">8th March 2021</p>                            | <p>APTS-ASP 2</p> <p>Introduction</p> <p>Preliminary material<br/> Expectation and probability<br/> Markov chains</p> <p>Some useful texts</p>   |
| <p>APTS-ASP 3<br/> └ Introduction</p> <h2 style="text-align: center;">Introduction</h2> <p>This module will introduce students to two important notions in stochastic processes — reversibility and martingales — identifying the basic ideas, outlining the main results and giving a flavour of some of the important ways in which these notions are used in statistics.</p>    | <p>2021-03-08</p> <p>APTS-ASP<br/> └ Introduction</p> <p style="text-align: right;">Introduction</p> <p style="text-align: right;"><small>This module will introduce students to two important notions in stochastic processes — reversibility and martingales — identifying the basic ideas, outlining the main results and giving a flavour of some of the important ways in which these notions are used in statistics.</small></p> <p>Probability provides one of the major underlying languages of statistics, and purely probabilistic concepts often cross over into the statistical world. So statisticians need to acquire some fluency in the general language of probability . . .</p>   |
| <p>APTS-ASP 5<br/> └ Introduction</p> <h2 style="text-align: center;">Learning Outcomes</h2> <p>After successfully completing this module an APTS student will be able to:</p> <ul style="list-style-type: none"> <li>▶ describe and calculate with the notion of a reversible Markov chain, both in discrete and continuous time;</li> <li>▶ describe the basic properties of discrete-parameter martingales and check whether the martingale property holds;</li> <li>▶ recall and apply some significant concepts from martingale theory;</li> <li>▶ explain how to use Foster-Lyapunov criteria to establish recurrence and speed of convergence to equilibrium for Markov chains.</li> </ul>  | <p>2021-03-08</p> <p>APTS-ASP<br/> └ Introduction</p> <p style="text-align: right;">Learning Outcomes</p> <p style="text-align: right;"><small>After successfully completing this module an APTS student will be able to:</small></p> <ul style="list-style-type: none"> <li>• describe and calculate with the notion of a reversible Markov chain, both in discrete and continuous time;</li> <li>• describe the basic properties of discrete-parameter martingales and check whether the martingale property holds;</li> <li>• recall and apply some significant concepts from martingale theory;</li> <li>• explain how to use Foster-Lyapunov criteria to establish recurrence and speed of convergence to equilibrium for Markov chains.</li> </ul> <p>These outcomes interact interestingly with various topics in applied statistics. However the most important aim of this module is to help students to acquire general awareness of further ideas from probability as and when that might be useful in their further research.</p> |

Preliminary material  
 Expectation and probability

For most APTS students most of this material should be well-known:

- ▶ Probability and conditional probability;
- ▶ Basic expectation and conditional expectation;
- ▶ discrete versus continuous (sums and integrals);
- ▶ limits versus expectations.

It is set out here, describing key ideas rather than details, in order to establish a sound common basis for the module.



Probability

1. Sample space  $\Omega$  of possible outcomes;
2. Probability  $\mathbb{P}$  assigns a number between 0 and 1 inclusive (the *probability*) to each (sensible) subset  $A \subseteq \Omega$  (we say  $A$  is an *event*);
3. Advanced (measure-theoretic) probability takes great care to specify what *sensible* means:  $A$  has to belong to a pre-determined  $\sigma$ -algebra  $\mathcal{F}$ , a family of subsets closed under countable unions and complements, often generated by open sets. We shall avoid these technicalities, though it will later be convenient to speak of  $\sigma$ -algebras  $\mathcal{F}_t$  as a shorthand for "information available by time  $t$ ".
4. Rules of probability:

**Normalization:**  $\mathbb{P}(\Omega) = 1$ ;

**$\sigma$ -additivity:** if  $A_1, A_2, \dots$  form a disjoint sequence of events then

$$\mathbb{P}(A_1 \cup A_2 \cup \dots) = \sum_i \mathbb{P}(A_i).$$



Conditional probability

1. We declare the *conditional probability* of  $A$  given  $B$  to be  $\mathbb{P}(A|B) = \mathbb{P}(A \cap B) / \mathbb{P}(B)$ , and declare the case when  $\mathbb{P}(B) = 0$  as undefined.
2. **Bayes:** if  $B_1, B_2, \dots$  is an exhaustive disjoint partition of  $\Omega$  then

$$\mathbb{P}(B_i|A) = \frac{\mathbb{P}(A|B_i) \mathbb{P}(B_i)}{\sum_j \mathbb{P}(A|B_j) \mathbb{P}(B_j)}.$$

3. Conditional probabilities are clandestine random variables! Let  $X$  be the Bernoulli<sup>1</sup> random variable which indicates<sup>2</sup> event  $B$ . Consider the conditional probability of  $A$  given information of whether or not  $B$  occurs: it is random, being  $\mathbb{P}(A|B)$  if  $X = 1$  and  $\mathbb{P}(A|B^c)$  if  $X = 0$ .

<sup>1</sup>Taking values only 0 or 1.

<sup>2</sup> $X = 1$  exactly when  $B$  happens.



**Preliminary material**  
 Introduction to probability

For most APTS students most of this material should be well-known:

- Probability and conditional probability.
- Basic expectation and conditional expectation.
- discrete versus continuous (sums and integrals).
- limits versus expectations.

It is set out here, describing key ideas rather than details, in order to establish a sound common basis for the module.

This material uses a two-panel format. Left-hand panels present the theory, often using itemized lists. Right-hand panels present commentary and useful exercises (announced by "Test understanding"). It is likely that you will have seen most, if not all, of the preliminary material at undergraduate level. However syllabi are not uniform across universities; if some of this material is not well-known to you then:

- read through it to pick up the general sense and notation;
- supplement by reading (for example) the first five chapters of [Grimmett and Stirzaker \(2001\)](#).

**Probability**

1. Sample space  $\Omega$  of possible outcomes.
2. Probability  $\mathbb{P}$  assigns a number between 0 and 1 inclusive (the *probability*) to each (sensible) subset  $A \subseteq \Omega$  (we say  $A$  is an *event*).
3. Advanced (measure-theoretic) probability takes great care to specify what *sensible* means:  $A$  has to belong to a pre-determined  $\sigma$ -algebra  $\mathcal{F}$ , a family of subsets closed under countable unions and complements, often generated by open sets. We shall avoid these technicalities, though it will later be convenient to speak of  $\sigma$ -algebras  $\mathcal{F}_t$  as a shorthand for "information available by time  $t$ ".
4. Rules of probability:

**Normalization:**  $\mathbb{P}(\Omega) = 1$ ;

**$\sigma$ -additivity:** if  $A_1, A_2, \dots$  form a disjoint sequence of events then

$$\mathbb{P}(A_1 \cup A_2 \cup \dots) = \sum_i \mathbb{P}(A_i).$$

1. Example:  $\Omega = (-\infty, \infty)$ .
2. We could for example start with  $\mathbb{P}((a, b)) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-u^2/2} du$  and then use the rules of probability to determine probabilities for all manner of sensible subsets of  $(-\infty, \infty)$ .
3. In our example a "natural" choice for  $\mathcal{F}$  is the family of all sets generated from intervals by indefinitely complicated countably infinite combinations of countable unions and complements.
4. **Test understanding:** use these rules to explain
  - (a) why  $\mathbb{P}(\emptyset) = 0$ ,
  - (b) why  $\mathbb{P}(A^c) = 1 - \mathbb{P}(A)$  if  $A^c = \Omega \setminus A$ , and
  - (c) why it makes no sense in general to try to extend  $\sigma$ -additivity to uncountable unions such as  $(-\infty, \infty) = \bigcup_x \{x\}$ .

**Conditional probability**

1. We declare the *conditional probability* of  $A$  given  $B$  to be  $\mathbb{P}(A|B) = \mathbb{P}(A \cap B) / \mathbb{P}(B)$ , and declare the case when  $\mathbb{P}(B) = 0$  as undefined.
2. **Bayes:** if  $B_1, B_2, \dots$  is an exhaustive disjoint partition of  $\Omega$  then

$$\mathbb{P}(B_i|A) = \frac{\mathbb{P}(A|B_i) \mathbb{P}(B_i)}{\sum_j \mathbb{P}(A|B_j) \mathbb{P}(B_j)}.$$

3. Conditional probabilities are clandestine random variables! Let  $X$  be the Bernoulli random variable which indicates<sup>2</sup> event  $B$ . Consider the conditional probability of  $A$  given information of whether or not  $B$  occurs: it is random, being  $\mathbb{P}(A|B)$  if  $X = 1$  and  $\mathbb{P}(A|B^c)$  if  $X = 0$ .

<sup>1</sup>Taking values only 0 or 1.  
<sup>2</sup> $X = 1$  exactly when  $B$  happens.

1. Actually we *often* use limiting arguments to make sense of cases when  $\mathbb{P}(B) = 0$ .
2. Hence all of Bayesian statistics ...  
**Test understanding:** write out an explanation of why Bayes' theorem is a completely obvious consequence of the definitions of probability and conditional probability.
3. The idea of conditioning is developed in probability theory to the point where this notion (that conditional probabilities are random variables) becomes entirely natural not artificial.  
**Test understanding:** establish the law of inclusion and exclusion: if  $A_1, \dots, A_n$  are potentially overlapping events then

$$\begin{aligned} \mathbb{P}(A_1 \cup \dots \cup A_n) &= \mathbb{P}(A_1) + \dots + \mathbb{P}(A_n) \\ &\quad - (\mathbb{P}(A_1 \cap A_2) + \dots + \mathbb{P}(A_i \cap A_j) + \dots + \mathbb{P}(A_{n-1} \cap A_n)) \\ &\quad + \dots + (-1)^n \mathbb{P}(A_1 \cap \dots \cap A_n). \end{aligned}$$

Hint: represent RHS as expectation of expansion of  $1 - (1 - X_1) \dots (1 - X_n)$  for suitable Bernoulli random variables  $X_i$  indicating various  $A_i$ .

## Expectation

Statistical intuition about expectation is based on *properties*:

1. If  $X \geq 0$  is a non-negative random variable then we can define its (possibly infinite) *expectation*  $\mathbb{E}[X]$ .
2. If  $X = X^+ - X^- = \max\{X, 0\} - \max\{-X, 0\}$  is such that  $\mathbb{E}[X^\pm]$  are both finite<sup>3</sup> then set  $\mathbb{E}[X] = \mathbb{E}[X^+] - \mathbb{E}[X^-]$ .
3. Familiar properties of expectation follow from **linearity** ( $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$ ) and **monotonicity** ( $\mathbb{P}(X \geq a) = 1$  implies  $\mathbb{E}[X] \geq a$ ) for constants  $a, b$ .
4. Useful notation: for an event  $A$  write  $\mathbb{E}[X; A] = \mathbb{E}[X \mathbb{1}_A]$ , where  $\mathbb{1}_A$  is the Bernoulli random variable indicating  $A$ .
5. If  $X$  has countable range then  $\mathbb{E}[X] = \sum_x x \mathbb{P}(X = x)$ .
6. If  $X$  has *density*  $f_X$  then  $\mathbb{E}[X] = \int x f_X(x) dx$ .

<sup>3</sup>We wish to avoid having to make sense of  $\infty - \infty$ !



**Expectation**  
 Statistical intuition about expectation is based on properties:  
 1. If  $X \geq 0$  is a non-negative random variable then we can define its (possibly infinite) expectation  $\mathbb{E}[X]$ .  
 2. If  $X = X^+ - X^- = \max\{X, 0\} - \max\{-X, 0\}$  is such that  $\mathbb{E}[X^\pm]$  are both finite<sup>3</sup> then we set  $\mathbb{E}[X] = \mathbb{E}[X^+] - \mathbb{E}[X^-]$ .  
 3. Familiar properties of expectation follow from **linearity** ( $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$ ) and **monotonicity** ( $\mathbb{P}(X \geq a) = 1$  implies  $\mathbb{E}[X] \geq a$ ) for constants  $a, b$ .  
 4. Useful notation: for an event  $A$  write  $\mathbb{E}[X; A] = \mathbb{E}[X \mathbb{1}_A]$ , where  $\mathbb{1}_A$  is the Bernoulli random variable indicating  $A$ .  
 5. If  $X$  has countable range then  $\mathbb{E}[X] = \sum_x x \mathbb{P}(X = x)$ .  
 6. If  $X$  has density  $f_X$  then  $\mathbb{E}[X] = \int x f_X(x) dx$ .  
<sup>3</sup>We wish to avoid having to make sense of  $\infty - \infty$ !

1. Full definition of expectation takes 3 steps: obvious definition for Bernoulli random variables, finite range random variables by linearity, general case by monotonic limits  $X_n \uparrow X$ . The hard work lies in proving this is all consistent . . . .
2. Any decomposition as difference of *integrable* random variables will do.
3. **Test understanding**: using these properties
  - deduce  $\mathbb{E}[a] = a$  for constant  $a$ .
  - show *Markov's inequality*  $\mathbb{P}(X \geq a) \leq \frac{1}{a} \mathbb{E}[X]$  for  $X \geq 0, a > 0$ .
4. So in absolutely continuous case  $\mathbb{E}[X; A] = \int_A x f_X(x) dx$  and in discrete case  $\mathbb{E}[X; X = k] = k \mathbb{P}(X = k)$ .
5. Countable [=discrete] case: expectation defined exactly when sum converges absolutely.
6. Density [=absolutely continuous] case: expectation defined exactly when integral converges absolutely.

## Independence

Events  $A$  and  $B$  are *independent* if  $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$ .

1. This definition can be extended to more than two events:  $A_1, A_2, \dots, A_n$  are independent if for any set  $J \subseteq \{1, \dots, n\}$

$$\mathbb{P}(\cap_{j \in J} A_j) = \prod_{j \in J} \mathbb{P}(A_j).$$

2. If  $A$  and  $B$  are independent, with  $\mathbb{P}(B) > 0$ , then  $\mathbb{P}(A|B) = \mathbb{P}(A)$ .

Random variables  $X$  and  $Y$  are independent if for all  $x, y \in \mathbb{R}$

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

In this case,  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ .



**Independence**  
 Events  $A$  and  $B$  are independent if  $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$ .  
 This definition can be extended to more than two events:  $A_1, A_2, \dots, A_n$  are independent if for any set  $J \subseteq \{1, \dots, n\}$   
 $\mathbb{P}(\cap_{j \in J} A_j) = \prod_{j \in J} \mathbb{P}(A_j)$ .  
 If  $A$  and  $B$  are independent, with  $\mathbb{P}(B) > 0$ , then  $\mathbb{P}(A|B) = \mathbb{P}(A)$ .  
 Random variables  $X$  and  $Y$  are independent if for all  $x, y \in \mathbb{R}$   
 $\mathbb{P}(X \leq x, Y \leq y) = \mathbb{P}(X \leq x)\mathbb{P}(Y \leq y)$ .  
 In this case,  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ .

1. Note that it's *not* enough to simply ask for any two events  $A_i$  and  $A_j$  to be independent (i.e. pairwise independence)!  
**Test understanding**: find a set of three events which are *pairwise* independent, but for which  $\mathbb{P}(A_1 \cap A_2 \cap A_3) \neq \mathbb{P}(A_1)\mathbb{P}(A_2)\mathbb{P}(A_3)$ .
2. **Test understanding**: Show that:
  - if  $A$  and  $B$  are independent, then events  $A^c$  and  $B^c$  are independent;
  - if  $A_1, A_2, \dots, A_n$  are independent then
$$\mathbb{P}(A_1 \cup A_2 \cup \dots \cup A_n) = 1 - \prod_{i=1}^n \mathbb{P}(A_i^c).$$
3. This is equivalent to  $\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A)\mathbb{P}(Y \in B)$  for any (Borel) sets  $A, B \subseteq \mathbb{R}$ .
4. **Test understanding**: Suppose that  $X$  and  $Y$  are independent, and are either both discrete or both absolutely continuous; show that  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ .

## Generating functions

We're often interested in expectations of *functions* of random variables (e.g. recall that in the discrete case  $\mathbb{E}[g(X)] = \sum_x g(x)\mathbb{P}(X = x)$ ).

Some functions are particularly useful:

1. when  $g(x) = z^x$  for some  $z \geq 0$  we obtain the *probability generating function (pgf)* of  $X$ ,  $G_X(z) = \mathbb{E}[z^X]$ ;
2. when  $g(x) = e^{tx}$  we get the *moment generating function (mgf)*,  $m_X(t) = \mathbb{E}[e^{tX}]$ ;
3. when  $g(x) = e^{itx}$ , where  $i = \sqrt{-1}$ , we get the *characteristic function* of  $X$ ,  $\phi_X(t)$ .



**Generating functions**  
 We're often interested in expectations of functions of random variables (e.g. recall that in the discrete case  $\mathbb{E}[g(X)] = \sum_x g(x)\mathbb{P}(X = x)$ ).  
 Some functions are particularly useful:  
 1. when  $g(x) = z^x$  for some  $z \geq 0$  we obtain the *probability generating function (pgf)* of  $X$ ,  $G_X(z) = \mathbb{E}[z^X]$ .  
 2. when  $g(x) = e^{tx}$  we get the *moment generating function (mgf)*,  $m_X(t) = \mathbb{E}[e^{tX}]$ .  
 3. when  $g(x) = e^{itx}$ , where  $i = \sqrt{-1}$ , we get the *characteristic function* of  $X$ ,  $\phi_X(t)$ .

1. The pgf is only really useful when  $X$  takes values in  $\{0, 1, 2, \dots\}$ . The mgf and characteristic function are more generally applicable.
2. **Test understanding**: Show that  $\mathbb{E}[X] = G'_X(1)$  and  $\mathbb{P}(X = k) = G_X^{(k)}(0)/k!$  (where  $G_X^{(k)}(0)$  means the  $k^{\text{th}}$  derivative of  $G_X$ , evaluated at  $z = 0$ ).
3. **Test understanding**: Show that  $\mathbb{E}[X] = m'_X(0)$  and

$$m_X(t) = \sum_k \frac{\mathbb{E}[X^k]}{k!} t^k.$$

## Uses of generating functions

Generating functions are helpful in many ways. In particular:

1. They can be used to determine distributions
2. They can often provide an easy route to finding e.g. moments of a distribution
3. They're useful when working with sums of independent random variables, since the generating function of a *convolution* of distributions is the product of their generating functions. So

$$G_{X+Y}(z) = G_X(z)G_Y(z) \text{ etc.}$$



Uses of generating functions  
 Generating functions are helpful in many ways. In particular:  
 1. They can be used to determine distributions  
 2. They can often provide an easy route to finding e.g. moments of a distribution  
 3. They're useful when working with sums of independent random variables, since the generating function of a convolution of distributions is the product of their generating functions. So  
 $G_{X+Y}(z) = G_X(z)G_Y(z)$  etc.

1. Characteristic functions always uniquely determine distributions (i.e. there is a one-to-one correspondence between a distribution and its characteristic function); the same is true of pgfs and distributions on  $\{0, 1, \dots\}$ ; mgfs are slightly more complicated, but *mostly* they can be used to identify a distribution. See [Grimmett and Stirzaker \(2001\)](#) for more on this.
2. See the two exercises on the previous notes slide!
3. **Test understanding:** show that if  $X$  and  $Y$  are independent random variables then  $G_{X+Y}(z) = G_X(z)G_Y(z)$ ,  $m_{X+Y}(t) = m_X(t)m_Y(t)$  and  $\phi_{X+Y}(t) = \phi_X(t)\phi_Y(t)$ . (Only one argument is needed to see all three results!) Use the first of these as a quick method of proving that the sum of two independent Poisson random variables is itself Poisson.

## Conditional Expectation (I): property-based definition

1. Conventional definitions treat two separate cases (discrete and absolutely continuous):
  - ▶  $\mathbb{E}[X|Y=y] = \sum_x x \mathbb{P}(X=x|Y=y)$ ,
  - ▶  $\mathbb{E}[X|Y=y] = \int x f_{X|Y=y}(x) dx$ .
  - ... but what if  $X$  is mixed discrete/continuous? or worse?

Focus on *properties* to get unified approach:

2. If  $\mathbb{E}[X] < \infty$ , we say  $Z = \mathbb{E}[X|Y]$  if:
  - (a)  $\mathbb{E}[Z] < \infty$ ;
  - (b)  $Z$  is a function of  $Y$ ;
  - (c)  $\mathbb{E}[Z; A] = \mathbb{E}[X; A]$  for events  $A$  defined in terms of  $Y$ .
 This defines  $\mathbb{E}[X|Y]$  uniquely, up to events of prob 0.
3. We can now define  $\mathbb{E}[X|Y_1, Y_2, \dots]$  simply by using "is a function of  $Y_1, Y_2, \dots$ " and "defined in terms of  $Y_1, Y_2, \dots$ ", etc. Indeed we often write  $\mathbb{E}[X|\mathcal{G}]$ , where ( $\sigma$ -algebra)  $\mathcal{G}$  represents information conveyed by a specified set of random variables and events.



Conditional Expectation (I): property-based definition  
 1. Conventional definitions treat two separate cases (discrete and absolutely continuous)  
 +  $\mathbb{E}[X|Y=y] = \sum_x x \mathbb{P}(X=x|Y=y)$   
 +  $\mathbb{E}[X|Y=y] = \int x f_{X|Y=y}(x) dx$   
 but what if  $X$  is mixed discrete/continuous? or worse?  
 Focus on properties to get unified approach  
 2. If  $\mathbb{E}[X] < \infty$ , we say  $Z = \mathbb{E}[X|Y]$  if:  
 (a)  $\mathbb{E}[Z] < \infty$   
 (b)  $Z$  is a function of  $Y$   
 (c)  $\mathbb{E}[Z; A] = \mathbb{E}[X; A]$  for events  $A$  defined in terms of  $Y$   
 This defines  $\mathbb{E}[X|Y]$  uniquely, up to events of prob 0.  
 3. We can now define  $\mathbb{E}[X|Y_1, Y_2, \dots]$  simply by using "is a function of  $Y_1, Y_2, \dots$ " and "defined in terms of  $Y_1, Y_2, \dots$ ", etc. Indeed we often write  $\mathbb{E}[X|\mathcal{G}]$ , where ( $\sigma$ -algebra)  $\mathcal{G}$  represents information conveyed by a specified set of random variables and events

Conditional expectation needs careful definition to capture all cases. But focus on *properties* to build intuitive understanding.

1. Notice that conditional expectation is also properly viewed as a random variable.
2. - " $\mathbb{E}[Z] < \infty$ " is needed to get a good definition of *any* kind of expectation;  
 - We could express " $Z$  is a function of  $Y$ " etc more formally using measure theory if we had to;  
 - We need (b) to rule out  $Z = X$ , for example.  
**Test understanding:** verify that the discrete definition of conditional expectation satisfies the three properties (a), (b), (c). Hint: use  $A$  running through events  $A = \{Y = y\}$  for  $y$  in the range of  $Y$ .
3. **Test understanding:** suppose  $X_1, X_2, \dots, X_n$  are independent and identically distributed, with finite absolute mean  $\mathbb{E}[|X_i|] < \infty$ . Use symmetry and linearity to show  $\mathbb{E}[X_1|X_1 + \dots + X_n] = \frac{1}{n}(X_1 + \dots + X_n)$ .

## Conditional Expectation (II): some other properties

Many facts about conditional expectation follow easily from this property-based approach. For example:

1. Linearity:  $\mathbb{E}[aX + bY|Z] = a\mathbb{E}[X|Z] + b\mathbb{E}[Y|Z]$ ;
  2. "Tower property":  $\mathbb{E}[\mathbb{E}[X|Y, Z]|Y] = \mathbb{E}[X|Y]$ ;
  3. "Taking out what is known":  $\mathbb{E}[f(Y)X|Y] = f(Y)\mathbb{E}[X|Y]$ ;
- and variations involving more than one or two conditioning random variables ...



Conditional Expectation (II): some other properties  
 Many facts about conditional expectation follow easily from this property-based approach. For example:  
 1. Linearity:  $\mathbb{E}[aX + bY|Z] = a\mathbb{E}[X|Z] + b\mathbb{E}[Y|Z]$   
 2. "Tower property":  $\mathbb{E}[\mathbb{E}[X|Y, Z]|Y] = \mathbb{E}[X|Y]$   
 3. "Taking out what is known":  $\mathbb{E}[f(Y)X|Y] = f(Y)\mathbb{E}[X|Y]$   
 and variations involving more than one or two conditioning random variables

**Test understanding:** explain how these follow from the property-based definition. Hints:

1. Use  $\mathbb{E}[aX + bY; A] = a\mathbb{E}[X; A] + b\mathbb{E}[Y; A]$ .
2. Take a deep breath and use property (c) of conditional expectation twice. Suppose  $A$  is defined in terms of  $Y$ . Then  $\mathbb{E}[\mathbb{E}[X|Y, Z]; A] = \mathbb{E}[\mathbb{E}[X|Y, Z]; A]$  and  $\mathbb{E}[\mathbb{E}[X|Y, Z]; A] = \mathbb{E}[X; A]$ .
3. Just consider when  $f$  has finite range, and use the (finite) sum  $\mathbb{E}[\mathbb{E}[f(Y)X|Y]; A] = \sum_t \mathbb{E}[\mathbb{E}[f(Y)X|Y]; A \cap \{f(Y) = t\}]$ . But then use  $\mathbb{E}[\mathbb{E}[f(Y)X|Y]; A \cap \{f(Y) = t\}] = \mathbb{E}[\mathbb{E}[tX|Y]; A \cap \{f(Y) = t\}] = \mathbb{E}[t\mathbb{E}[X|Y]; A \cap \{f(Y) = t\}] = \mathbb{E}[f(Y)\mathbb{E}[X|Y]; A \cap \{f(Y) = t\}]$ . General case now follows by approximation arguments.

### Conditional Expectation (III): Jensen's inequality

This is powerful and yet rather easy to prove.

**Theorem**

Let  $\phi$  be a convex function ("curves upwards",  $\phi'' \geq 0$  if smooth). Suppose the random variable  $X$  is such that  $\mathbb{E}[|X|] < \infty$  and  $\mathbb{E}[\phi(X)] < \infty$ . Then

$$\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)],$$

and the same is true for conditional expectations:

$$\phi(\mathbb{E}[X|\mathcal{G}]) \leq \mathbb{E}[\phi(X)|\mathcal{G}]$$

for some conditioning information  $\mathcal{G}$ .

Clue to proof: any convex function can be represented as supremum of all affine functions  $ax + b$  lying below it.



**Conditional Expectation (III): Jensen's inequality**  
 This is powerful and yet rather easy to prove.  
**Theorem**  
 Let  $\phi$  be a convex function ("curves upwards",  $\phi'' \geq 0$  if smooth). Suppose the random variable  $X$  is such that  $\mathbb{E}[|X|] < \infty$  and  $\mathbb{E}[\phi(X)] < \infty$ . Then  

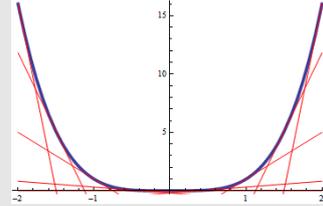
$$\phi(\mathbb{E}[X]) \leq \mathbb{E}[\phi(X)]$$
  
 and the same is true for conditional expectations:  

$$\phi(\mathbb{E}[X|\mathcal{G}]) \leq \mathbb{E}[\phi(X)|\mathcal{G}]$$
  
 for some conditioning information  $\mathcal{G}$ .  
 Clue to proof: any convex function can be represented as supremum of all affine functions  $ax + b$  lying below it.

Consider the simple convex function  $\phi(x) = x^2$ . We deduce, if  $X$  has finite second moment then

$$(\mathbb{E}[X|\mathcal{G}])^2 \leq \mathbb{E}[X^2|\mathcal{G}].$$

Here's a picture to illustrate the clue to the proof of Jensen's inequality in case  $\phi(x) = x^2$ :



### Limits versus expectations

- Often the crux of a piece of mathematics is whether one can exchange limiting operations such as  $\lim \sum \leftrightarrow \sum \lim$ . Here are a few very useful results on this, expressed in the language of expectations.
- Monotone Convergence Theorem:** If  $\mathbb{P}(X_n \uparrow Y) = 1$  and  $\mathbb{E}[X_1] > -\infty$  then  $\lim_n \mathbb{E}[X_n] = \mathbb{E}[\lim_n X_n] = \mathbb{E}[Y]$ .
- Dominated Convergence Theorem:** If  $\mathbb{P}(X_n \rightarrow Y) = 1$  and  $|X_n| \leq Z$  where  $\mathbb{E}[Z] < \infty$  then  $\lim_n \mathbb{E}[X_n] = \mathbb{E}[\lim_n X_n] = \mathbb{E}[Y]$ .
- Fubini's Theorem:** If  $\mathbb{E}[|f(X, Y)|] < \infty$ ,  $X, Y$  are independent,  $g(y) = \mathbb{E}[f(X, y)]$ ,  $h(x) = \mathbb{E}[f(x, Y)]$  then  $\mathbb{E}[g(Y)] = \mathbb{E}[f(X, Y)] = \mathbb{E}[h(X)]$ .
- Fatou's lemma:** If  $\mathbb{P}(X_n \rightarrow Y) = 1$  and  $X_n \geq 0$  for all  $n$  then  $\mathbb{E}[Y] \leq \lim_n \inf_{m \geq n} \mathbb{E}[X_m]$ .



**Limits versus expectations**  
 1. Often the crux of a piece of mathematics is whether one can exchange limiting operations such as  $\lim \sum \leftrightarrow \sum \lim$ . Here are a few very useful results on this, expressed in the language of expectations.  
 2. **Monotone Convergence Theorem:** If  $\mathbb{P}(X_n \uparrow Y) = 1$  and  $\mathbb{E}[X_1] > -\infty$  then  $\lim_n \mathbb{E}[X_n] = \mathbb{E}[\lim_n X_n] = \mathbb{E}[Y]$ .  
 3. **Dominated Convergence Theorem:** If  $\mathbb{P}(X_n \rightarrow Y) = 1$  and  $|X_n| \leq Z$  where  $\mathbb{E}[Z] < \infty$  then  $\lim_n \mathbb{E}[X_n] = \mathbb{E}[\lim_n X_n] = \mathbb{E}[Y]$ .  
 4. **Fubini's Theorem:** If  $\mathbb{E}[|f(X, Y)|] < \infty$ ,  $X, Y$  are independent,  $g(y) = \mathbb{E}[f(X, y)]$ ,  $h(x) = \mathbb{E}[f(x, Y)]$  then  $\mathbb{E}[g(Y)] = \mathbb{E}[f(X, Y)] = \mathbb{E}[h(X)]$ .  
 5. **Fatou's lemma:** If  $\mathbb{P}(X_n \rightarrow Y) = 1$  and  $X_n \geq 0$  for all  $n$  then  $\mathbb{E}[Y] \leq \lim_n \inf_{m \geq n} \mathbb{E}[X_m]$ .

- As we formulate this in expectation language, our results apply equally to sums and integrals.
- Note that the  $X_n$  must form an *increasing* sequence. We need  $\mathbb{E}[X_1] > -\infty$ . **Test understanding:** consider case of  $X_n = -1/(nU)$  for a fixed Uniform(0, 1) random variable.
- Note that convergence need not be monotonic here or in following. **Test understanding:** explain why it would be enough to have finite upper and lower bounds  $\alpha \leq X_n \leq \beta$ .
- Fubini exchanges expectations rather than an expectation and a limit.
- Try Fatou if all else fails. Note that something like  $X_n \geq 0$  is essential (a constant lower bound would suffice, though).

### Preliminary material

#### Markov chains

- Discrete-time countable-state-space basics:
  - Markov property, transition matrices;
  - irreducibility and aperiodicity;
  - transience and recurrence;
  - equilibrium equations and convergence to equilibrium.
- Discrete-time *countable*-state-space: why 'limit of sum need not always equal sum of limit'.
- Continuous-time countable-state-space: rates and  $Q$ -matrices.
- Definition and basic properties of Poisson counting process.



**Preliminary material**  
 below are:  
 • Discrete-time countable-state-space basics:  
 • Markov property, transition matrices  
 • irreducibility and aperiodicity  
 • transience and recurrence  
 • equilibrium equations and convergence to equilibrium  
 • Discrete-time countable-state-space  
 why 'limit of sum need not always equal sum of limit'  
 • Continuous-time countable-state-space: rates and  $Q$ -matrices.  
 • Definition and basic properties of Poisson counting process.

If some of this material is not well-known to you, then invest some time in looking over (for example) chapter 6 of [Grimmett and Stirzaker \(2001\)](#).

Instead of "countable-state-space" Markov chains, we'll use the shorter phrase "discrete Markov chains".

### Basic properties for discrete time and space case

1. Markov chain  $X = \{X_0, X_1, X_2, \dots\}$ :  $X$  at time  $t$  is in state  $X_t = x$ . View states  $x$  as integers.
2.  $X$  must have the **Markov property**:  $p_{xy} = p(x, y) = \mathbb{P}(X_{t+1} = y | X_t = x, X_{t-1}, \dots)$  must depend only on  $x, y$ , not on rest of past. (Our chains will be *time-homogeneous*, meaning no  $t$  dependence either.)
3. Chain behaviour is specified by (a) initial state  $X_0$  (could be random) and (b) table of **transition probabilities**  $p_{xy}$ .
4. Important **matrix** structure: if  $p_{xy}$  are arranged in matrix  $P$  then  $(i, j)$ <sup>th</sup> entry of  $P^n = P \cdot \dots \cdot P$  ( $n$  times) is  $p_{ij}^{(n)} = \mathbb{P}(X_n = j | X_0 = i)$ .  
 Equivalent: **Chapman-Kolmogorov equations**

$$p_{ij}^{(n+m)} = \sum_k p_{ik}^{(n)} p_{kj}^{(m)}$$



**Basic properties for discrete time and space case**  
 1. Markov chain  $X = \{X_0, X_1, X_2, \dots\}$ :  $X$  at time  $t$  is in state  $X_t = x$ . View states  $x$  as integers.  
 2.  $X$  must have the **Markov property**:  $p_{xy} = p(x, y) = \mathbb{P}(X_{t+1} = y | X_t = x, X_{t-1}, \dots)$  must depend only on  $x, y$ , not on rest of past. (Our chains will be time-homogeneous, meaning no  $t$  dependence either.)  
 3. Chain behaviour is specified by (a) initial state  $X_0$  (could be random) and (b) table of transition probabilities  $p_{xy}$ .  
 4. Important matrix structure: if  $p_{xy}$  are arranged in matrix  $P$  then  $(i, j)$ <sup>th</sup> entry of  $P^n = P \cdot \dots \cdot P$  ( $n$  times) is  $p_{ij}^{(n)} = \mathbb{P}(X_n = j | X_0 = i)$ .  
 Equivalent: Chapman-Kolmogorov equations  

$$p_{ij}^{(n+m)} = \sum_k p_{ik}^{(n)} p_{kj}^{(m)}$$

1. More general countable discrete state-spaces can always be indexed by integers
2. The example of "Markov's other chain" below shows we need to **insist** on the possibility of conditioning by further past  $X_{t-1}, \dots$  in this definition.  
 Note  $\sum_y p_{xy} = 1$  by "law of total probability".
3. Example: some word transition probabilities arising in the "random English" example given immediately below:  

$$\begin{aligned} P(\text{"round"}|\text{"all"}) &= 0.50 & P(\text{"contact"}|\text{"all"}) &= 0.50 & P(\text{"hearing"}|\text{"ocean,"}) &= 1.00 \\ P(\text{"first","go"}) &= 1.00 & P(\text{"As"}|\text{"up,"}) &= 1.00 & P(\text{"Every"}|\text{"day,"}) &= 1.00 \\ P(\text{"woman"}|\text{"young"}) &= 0.33 & & & P(\text{"prince","young"}) &= 0.33 \\ P(\text{"man"}|\text{"young"}) &= 0.33 & & & P(\text{"on"}|\text{"enjoyed"}) &= 1.00 \quad \dots \end{aligned}$$
4. **Test understanding**: show how the Chapman-Kolmogorov equations follow from considerations of conditional probability and the Markov property.

### Example: Models for language following Markov

How to generate "random English" as a Markov chain:

1. Take a large book in electronic form, for example Tolstoy's "War and Peace" (English translation).
2. Use it to build a table of digram frequencies (digram = pair of consecutive letters).
3. Convert frequencies into conditional probabilities of one letter following another, and use these to form a Markov chain to generate "random English".

It is an amusing if substantial exercise to use this as a prior for Bayesian decoding of simple substitution codes.



**Example: Models for language following Markov**  
 How to generate "random English" as a Markov chain  
 1. Take a large book in electronic form, for example Tolstoy's "War and Peace" (English translation).  
 2. Use it to build a table of digram frequencies (digram = pair of consecutive letters).  
 3. Convert frequencies into conditional probabilities of one letter following another, and use these to form a Markov chain to generate "random English".  
 It is an amusing if substantial exercise to use this as a prior for Bayesian decoding of simple substitution codes.

1. The World-Web Web has made this part much easier: try Project Gutenberg ([www.gutenberg.org/etext/2600](http://www.gutenberg.org/etext/2600)).
2. Skill is required in deciding *which* letters to use: should one use all, or some, punctuation? Certainly need to use spaces.
3. Trigrams would be more impressive. Indeed, one needs to work at the level of words to simulate something like English.  
 Here is example output based on a children's fable:  
*It was able to the end of great daring but which when Rapunzel was a guardian has enjoyed on a time, after a faked morning departure more directly; over its days in a stratagem, which supported her hair into the risk of endless figures on supplanted sorrow. The prince's directive, to clamber down would come up easily, and perceived a grudge against humans for a convincing simulation of a nearby robotic despot. But then a computer typing in a convincing simulation of the traditional manner. However they settled in quality, and the prince thought for Rapunzel made its ward's face, that as she then a mere girl.*

### (Counter)example: Markov's other chain

Conditional probability can be subtle. Consider:

1. Independent Bernoulli  $X_0, X_2, X_4, \dots$  such that  $\mathbb{P}(X_{2n} = \pm 1) = \frac{1}{2}$ ;
2. Define  $X_{2n+1} = X_{2n}X_{2n+2}$  for  $n = 0, 1, \dots$ ; these also form an independent identically distributed sequence.
3.  $\mathbb{P}(X_{n+1} = \pm 1 | X_n) = \frac{1}{2}$  for any  $n \geq 1$ .
4. Chapman-Kolmogorov equations hold for any  $0 \leq k \leq n + k$ :  

$$\mathbb{P}(X_{n+k} = \pm 1 | X_0) = \sum_{y=\pm 1} \mathbb{P}(X_{n+k} = \pm 1 | X_k = y) \mathbb{P}(X_k = y | X_0)$$
5. Nevertheless,  $\mathbb{P}(X_2 = \pm 1 | X_1 = 1, X_0 = u)$  depends on  $u = \pm 1$ , so Markov property **fails** for  $X$ .



**(Counter)example: Markov's other chain**  
 Conditional probability can be subtle. Consider:  
 1. Independent Bernoulli  $X_0, X_2, X_4, \dots$  such that  $\mathbb{P}(X_{2n} = \pm 1) = \frac{1}{2}$ ;  
 2. Define  $X_{2n+1} = X_{2n}X_{2n+2}$  for  $n = 0, 1, \dots$ ; these also form an independent identically distributed sequence.  
 3.  $\mathbb{P}(X_{n+1} = \pm 1 | X_n) = \frac{1}{2}$  for any  $n \geq 1$ .  
 4. Chapman-Kolmogorov equations hold for any  $0 \leq k \leq n + k$ :  

$$\mathbb{P}(X_{n+k} = \pm 1 | X_0) = \sum_{y=\pm 1} \mathbb{P}(X_{n+k} = \pm 1 | X_k = y) \mathbb{P}(X_k = y | X_0)$$
  
 5. Nevertheless,  $\mathbb{P}(X_2 = \pm 1 | X_1 = 1, X_0 = u)$  depends on  $u = \pm 1$ , so Markov property fails for  $X$ .

Example taken from **Grimmett and Stirzaker (2001)**.

Note that the entirety of random variables  $X_0, X_1, X_2, \dots$  are most certainly *not* independent!  
**Test understanding** by checking these calculations.

It is usual in stochastic modelling to *start* by specifying that a given random process  $X = \{X_0, X_1, X_2, \dots\}$  is Markov, so this kind of issue is not often encountered in practice. However it is as well to be aware of it: conditioning is a subtle concept and should be treated with respect!

## Irreducibility and aperiodicity

1. A discrete Markov chain is *irreducible* if for all  $i$  and  $j$  it has a positive chance of visiting  $j$  at some positive time, if it is started at  $i$ .
2. It is *aperiodic* if one cannot divide state-space into non-empty subsets such that the chain progresses through the subsets in a periodic way. Simple symmetric walk (jumps  $\pm 1$ ) is *not* aperiodic.
3. If the chain is not irreducible, then we can compute the chance of it getting from one state to another using *first passage equations*: if

$$f_{ij} = \mathbb{P}(X_n = j \text{ for some positive } n | X_0 = i)$$

then solve linear equations for the  $f_{ij}$ .



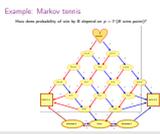
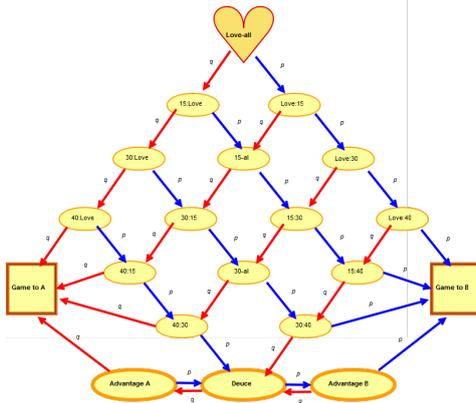
**Irreducibility and aperiodicity**

1. A discrete Markov chain is *irreducible* if for all  $i$  and  $j$  it has a positive chance of visiting  $j$  at some positive time, if it is started at  $i$ .
2. It is *aperiodic* if one cannot divide state space into non-empty subsets such that the chain progresses through the subsets in a periodic way. Simple symmetric walk (jumps  $\pm 1$ ) is not aperiodic.
3. If the chain is not irreducible, then we can compute the chance of it getting from one state to another using the passage equations: if
 
$$f_{ij} = \mathbb{P}(X_n = j \text{ for some positive } n | X_0 = i)$$
 then solve linear equations for the  $f_{ij}$ .

1. Consider the word game: change “good” to “evil” through other English words by altering just one letter at a time. Illustrative question (compare Gardner 1996): does your vocabulary of 4-letter English words form an irreducible Markov chain under moves which attempt random changes of letters? You can find an algorithmic approach to this question in Knuth (1993).
2. Equivalent definition: an irreducible chain  $X$  is aperiodic if its “independent double”  $\{(X_0, Y_0), (X_1, Y_1), \dots\}$  (for  $Y$  an independent copy of  $X$ ) is irreducible.
3. Because of the connection with matrices noted above, this can be cast in terms of rather basic linear algebra. First passage equations are still helpful in analyzing irreducible chains: for example the chance of visiting  $j$  before  $k$  is the same as computing  $f_{ij}$  for the modified chain which stops on hitting  $k$ .

## Example: Markov tennis

How does probability of win by B depend on  $p = \mathbb{P}(B \text{ wins point})$ ?



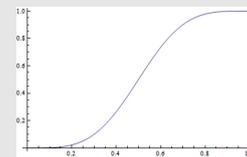
Use first passage equations, then solve linear equations for the  $f_{ij}$ , noting in particular

$$f_{\text{Game to A, Game to B}} = 0, \quad f_{\text{Game to B, Game to B}} = 1.$$

I obtain

$$f_{\text{Love-All, Game to B}} = \frac{p^4(15 - 34p + 28p^2 - 8p^3)}{1 - 2p + 2p^2},$$

graphed against  $p$  below:



## Transience and recurrence

1. Is it possible for a Markov chain  $X$  never to return to a starting state  $i$ ? If so then that state is said to be *transient*.
2. Otherwise the state is said to be *recurrent*.
3. Moreover if the return time  $T$  has finite mean then the state is said to be *positive-recurrent*.
4. Recurrent states which are not positive-recurrent are called *null-recurrent*.
5. States of an irreducible Markov chain are all recurrent if one is, all positive-recurrent if one is.



**Transience and recurrence**

1. Is it possible for a Markov chain  $X$  never to return to a starting state  $i$ ? If so then that state is said to be *transient*.
2. Otherwise the state is said to be *recurrent*.
3. Moreover if the return time  $T$  has finite mean then the state is said to be *positive-recurrent*.
4. Recurrent states which are not positive-recurrent are called *null-recurrent*.
5. States of an irreducible Markov chain are all recurrent if one is, all positive-recurrent if one is.

1. Example: asymmetric simple random walk (jumps  $\pm 1$ ): see Cox and Miller (1965) for a pretty explanation using strong law of large numbers.
2. Example: symmetric simple random walk (jumps  $\pm 1$ ).
3. As we will see, there exist infinite positive-recurrent chains (eg, “discrete AR(1)”).
4. Why “null”, “positive”? Terminology is motivated by the limiting behaviour of probability of being found in that state at large time. (Asymptotically zero if null-recurrent or transient: tends to  $1/\mathbb{E}[T]$  if aperiodic positive-recurrent.)
5. This is based on the criterion for recurrence of state  $i$ :  $\sum_n p_{ii}^{(n)} = \infty$ , which in turn arises from an application of generating functions. The criterion amounts to asserting, the chain is sure to return to a state  $i$  exactly when the *mean number* of returns is infinite.

### Recurrence/transience for random walks on $\mathbb{Z}$

Let  $X$  be a random walk on  $\mathbb{Z}$  which takes steps of size 1 with prob  $p$  and minus one with prob  $q = 1 - p$ . Define  $T_{0,1}$  to be the first time at which  $X$  hits 1, if it starts at 0. The probability generating function for this RV satisfies

$$G(z) = \mathbb{E} [z^{T_{0,1}}] = zp + zqG(z)^2$$

Solving this (and noting that we need to take the negative root!) we see that

$$G(z) = \frac{1 - \sqrt{1 - 4pqz^2}}{2qz},$$

and so  $\mathbb{P}(T_{0,1} < \infty) = \lim_{z \rightarrow 1} G(z) = \min\{p/q, 1\}$ . Thus if  $p < 1/2$  there is a positive chance that  $X$  never reaches state 1; by symmetry,  $X$  is recurrent iff  $p = 1/2$ .



Recurrence/transience for random walks on  $\mathbb{Z}$   
 Let  $X$  be a random walk on  $\mathbb{Z}$ , which takes steps of size 1 with prob  $p$  and minus one with prob  $q = 1 - p$ . Define  $T_{0,1}$  to be the first time at which  $X$  hits 1, if it starts at 0. The probability generating function for this RV satisfies  
 $G(z) = \mathbb{E}[z^{T_{0,1}}] = zp + zqG(z)^2$   
 Solving this (and noting that we need to take the negative root!) we see that  
 $G(z) = \frac{1 - \sqrt{1 - 4pqz^2}}{2qz}$   
 and so  $\mathbb{P}(T_{0,1} < \infty) = \lim_{z \rightarrow 1} G(z) = \min\{p/q, 1\}$ . Thus if  $p < 1/2$  there is a positive chance that  $X$  never reaches state 1; by symmetry,  $X$  is recurrent iff  $p = 1/2$ .

- Note that it's certainly possible to have  $\mathbb{P}(T_{0,1} < \infty) < 1$ , that is, for the random variable  $T_{0,1}$  to take the value  $\infty$ !
- Test understanding:** Show that the quadratic formula for  $G(z)$  holds by considering what can happen at time 1: argue that if  $X_1 = -1$  the time taken to get from -1 to 1 has the same distribution as the time taken to get from -1 to 0 plus the time to get from 0 to 1; these random variables are independent, and so the pgf of the sum is easy to work with...
- If we take the positive root then  $G(z) \rightarrow \infty$  as  $z \rightarrow 0$ , rather than to 0!
- Here we are using the fact that, since our state space is irreducible, state  $i$  is recurrent iff  $\mathbb{P}(T_{i,j} < \infty) = 1$  for all states  $j$ , where  $T_{i,j}$  is the first time that  $X$  hits  $j$  when started from  $i$ .

### Equilibrium of Markov chains

- If  $X$  is irreducible and positive-recurrent then it has a unique *equilibrium distribution*  $\pi$ : if  $X_0$  is random with distribution given by  $\mathbb{P}(X_0 = i) = \pi_i$  then  $\mathbb{P}(X_n = i) = \pi_i$  for any  $n$ .
- Moreover the equilibrium distribution viewed as a row vector solves the *equilibrium equations*:

$$\pi P = \pi, \quad \text{or} \quad \pi_j = \sum_i \pi_i p_{ij}.$$

- If in addition  $X$  is aperiodic then the equilibrium distribution is also the limiting distribution:

$$\mathbb{P}(X_n = i) \rightarrow \pi_i \quad \text{as } n \rightarrow \infty.$$



Equilibrium of Markov chains  
 1. If  $X$  is irreducible and positive-recurrent then it has a unique equilibrium distribution  $\pi$ . If  $X_0$  is random with distribution given by  $\mathbb{P}(X_0 = i) = \pi_i$  then  $\mathbb{P}(X_n = i) = \pi_i$  for any  $n$ .  
 2. Moreover the equilibrium distribution viewed as a row vector solves the equilibrium equations  
 $\pi P = \pi, \quad \text{or} \quad \pi_j = \sum_i \pi_i p_{ij}.$   
 3. If in addition  $X$  is aperiodic then the equilibrium distribution is also the limiting distribution  
 $\mathbb{P}(X_n = i) \rightarrow \pi_i \quad \text{as } n \rightarrow \infty.$

- In general the chain continues moving, but the marginal probabilities at time  $n$  do not change.
- Test understanding:** Show that the 2-state Markov chain with transition probability matrix  $\begin{pmatrix} 0.1 & 0.9 \\ 0.8 & 0.2 \end{pmatrix}$  has equilibrium distribution  $\pi = (0.470588 \dots, 0.529412 \dots)$ . Note that you need to use the fact that  $\pi_1 + \pi_2 = 1$ : this is *always* an important extra fact to use in determining a Markov chain's equilibrium distribution!
- This limiting result is of great importance in MCMC. If aperiodicity fails then it is always possible to sub-sample to convert to the aperiodic case on a subset of state-space. Note 4 of previous segment shows possibility of computing mean recurrence time using matrix arithmetic. NB:  $\pi_i$  can also be interpreted as "mean time in state  $i$ ".

### Sums of limits and limits of sums

- Finite state-space discrete Markov chains have a useful simplifying property: they are always positive-recurrent if they are irreducible.
- This can be proved by using a result, that for null-recurrent or transient states  $j$  we find  $p_{ij}^{(n)} \rightarrow 0$  as  $n \rightarrow \infty$ , for all other states  $i$ . Hence a contradiction:

$$\sum_j \lim_{n \rightarrow \infty} p_{ij}^{(n)} = \lim_{n \rightarrow \infty} \sum_j p_{ij}^{(n)}$$

and the right-hand sum equals 1 from "law of total probability", while left-hand sum equals  $\sum 0 = 0$  by null-recurrence.

- This argument fails for infinite state-space as it is incorrect arbitrarily to exchange infinite limiting operations:  
 $\lim \sum \neq \sum \lim$  in general.



Sums of limits and limits of sums  
 1. Finite state-space discrete Markov chains have a useful simplifying property: they are always positive-recurrent if they are irreducible.  
 2. This can be proved by using a result, that for null-recurrent or transient states  $j$  we find  $p_{ij}^{(n)} \rightarrow 0$  as  $n \rightarrow \infty$ , for all other states  $i$ . Hence a contradiction  
 $\sum_j \lim_{n \rightarrow \infty} p_{ij}^{(n)} = \lim_{n \rightarrow \infty} \sum_j p_{ij}^{(n)}$   
 and the right-hand sum equals 1 from "law of total probability", while left-hand sum equals  $\sum 0 = 0$  by null-recurrence.  
 3. This argument fails for infinite state-space as it is incorrect arbitrarily to exchange infinite limiting operations.  
 $\lim \sum \neq \sum \lim$  in general.

- Some argue that *all* Markov chains met in practice are finite, since we work on finite computers with finite floating point arithmetic. Do you find this argument convincing or not?
- The result used here puts the "null" in *null-recurrence*.
- We have earlier summarized the principal theorems which deliver checkable conditions as to when one can make this exchange.

Note that the simple random walk (irreducible but null-recurrent or transient) is the simplest practical example of why one must not carelessly exchange infinite limiting operations!

Continuous-time countable state-space Markov chains (a rough guide)

1. Definition of continuous-time (countable) discrete state-space (time-homogeneous) Markov chain  $X = \{X_t : t \geq 0\}$ : for  $s, t > 0$

$$p_t(x, y) = \mathbb{P}(X_{s+t} = y | X_s = x, X_u \text{ for various } u \leq s)$$

depends only on  $x, y, t$ , not on rest of past.

2. Organize  $p_t(x, y)$  into matrices  $P(t) = \{p_t(x, y) : \text{states } x, y\}$ ; as in discrete case  $P(t) \cdot P(s) = P(t + s)$  and  $P(0)$  is identity matrix.
3. (Try to) compute time derivative:  $Q = (d/dt)P(t)|_{t=0}$  is matrix of transition rates  $q(x, y)$ .



Continuous-time countable state-space Markov chains (rough guide)  
 1. Definition of continuous-time (countable) discrete state-space (time-homogeneous) Markov chain  $X = \{X_t : t \geq 0\}$  for  $s, t > 0$   
 $p_t(x, y) = \mathbb{P}(X_{s+t} = y | X_s = x, X_u \text{ for various } u \leq s)$   
 depends only on  $x, y, t$ , not on rest of past.  
 2. Organize  $p_t(x, y)$  into matrices  $P(t) = \{p_t(x, y) : \text{states } x, y\}$ , as in discrete case  $P(t) \cdot P(s) = P(t + s)$  and  $P(0)$  is identity matrix.  
 3. (Try to) compute time derivative:  $Q = (d/dt)P(t)|_{t=0}$  is matrix of transition rates  $q(x, y)$ .

This is a very rough guide: I pondered for a while whether to add this to prerequisites, since most of what I want to talk about will be in discrete time. I decided to add it in the end because sometimes the easiest examples in Markov chains are in continuous-time. The important point to grasp is that if we know the transition rates  $q(x, y)$  then we can write down differential equations to define the transition probabilities and so the chain. We don't necessarily try to solve the equations ...

1. For short, write  $p_t(x, y) = \mathbb{P}(X_{s+t} = y | X_s = x, \mathcal{F}_s)$  where  $\mathcal{F}_s$  represents all possible information about the past at time  $s$ .
2. From here on I omit many "under sufficient regularity" statements. Norris (1998) gives a careful treatment.
3. The row-sums of  $P(t)$  all equal 1 ("law of total probability"). Hence the row sums of  $Q$  ought to be 0 with non-positive diagonal entries  $q(x, x) = -q(x)$  measuring rate of leaving  $x$ .

Continuous-time countable state-space Markov chains (a rough guide continued)

For suitably regular continuous-time countable state-space Markov chains, we can use the  $Q$ -matrix  $Q$  to simulate the chain as follows:

1. rate of leaving state  $x$  is  $q(x) = \sum_{y \neq x} q(x, y)$  (since row sums of  $Q$  should be zero). Time till departure is  $\text{Exponential}(q(x))$ ;
2. on departure from  $x$ , go straight to state  $y \neq x$  with probability  $q(x, y)/q(x)$ .



Continuous-time countable state-space Markov chains (rough guide continued)  
 For suitably regular continuous-time countable state-space Markov chains, we can use the  $Q$ -matrix  $Q$  to simulate the chain as follows:  
 1. rate of leaving state  $x$  is  $q(x) = \sum_{y \neq x} q(x, y)$  (since row sums of  $Q$  should be zero). Time till departure is  $\text{Exponential}(q(x))$ ;  
 2. on departure from  $x$ , go straight to state  $y \neq x$  with probability  $q(x, y)/q(x)$ .

1. Why an exponential distribution? Because an effect of the Markov property is to require the holding time until the first transition to have a memory-less property—which characterizes Exponential distributions. Here it is relevant to note that "minimum of independent Exponential random variables is Exponential".
2. This also follows rather directly from the Markov property. Note that this shows two strong limitations of continuous-time Markov chains as stochastic models: the Exponential distribution of holding times may be unrealistic; and the state to which a transition is made does not depend on actual length of holding time. Of course, people have worked on generalizations (keyword: semi-Markov processes).

Continuous-time countable state-space Markov chains (a rough guide continued)

1. Compute the  $s$ -derivative of  $P(s) \cdot P(t) = P(s + t)$ . This yields the famous "Kolmogorov backwards equations":

$$Q \cdot P(t) = P(t)'$$

The other way round yields the "Kolmogorov forwards equations":

$$P(t) \cdot Q = P(t)'$$

2. If statistical equilibrium holds then the transition probabilities should converge to limiting values as  $t \rightarrow \infty$ : applying this to the forwards equation we expect the equilibrium distribution  $\pi$  to solve

$$\pi \cdot Q = 0.$$



Continuous-time countable state-space Markov chains (rough guide continued)  
 1. Compute the  $s$ -derivative of  $P(s) \cdot P(t) = P(s + t)$ . This yields the famous "Kolmogorov backwards equations":  
 $Q \cdot P(t) = P(t)'$   
 The other way round yields the "Kolmogorov forwards equations":  
 $P(t) \cdot Q = P(t)'$   
 2. If statistical equilibrium holds then the transition probabilities should converge to limiting values as  $t \rightarrow \infty$ : applying this to the forwards equation we expect the equilibrium distribution  $\pi$  to solve  
 $\pi \cdot Q = 0.$

1. **Test understanding:** use calculus to derive  

$$\sum_z p_s(x, z) p_t(z, y) = p_{s+t}(x, y) \text{ gives } \sum_z q(x, z) p_t(z, y) = \frac{\partial}{\partial t} p_t(x, y),$$

$$\sum_z p_t(x, z) p_s(z, y) = p_{t+s}(x, y) \text{ gives } \sum_z p_t(x, z) q(z, y) = \frac{\partial}{\partial t} p_t(x, y).$$

Note the shameless exchange of differentiation and summation over potentially infinite state-space ...
2. **Test understanding:** applying this idea to the backwards equation gets us nothing, as a consequence of the vanishing of row sums of  $Q$ . In extended form  $\pi \cdot Q = 0$  yields the important equilibrium equations  

$$\sum_z \pi(z) q(z, y) = 0.$$

## Example: the Poisson process

We use the above theory to *define* chains by specifying the non-zero rates. Consider the case when  $X$  counts the number of people arriving at random at constant rate:

1. Stipulate that the number  $X_t$  of people in system at time  $t$  forms a Markov chain.
2. Transition rates: people arrive one-at-a-time at constant rate, so  $q(x, x + 1) = \lambda$ .

One can solve the Kolmogorov differential equations in this case:

$$\mathbb{P}(X_t = n | X_0 = 0) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}.$$



Example: the Poisson process  
 We use the above theory to define chains by specifying the non-zero rates. Consider the case when  $X$  counts the number of people arriving at random at constant rate.  
 1. Stipulate that the number  $X_t$  of people in system at time  $t$  forms a Markov chain.  
 2. Transition rates: people arrive one-at-a-time at constant rate, so  $q(x, x + 1) = \lambda$ .  
 One can solve the Kolmogorov differential equations in this case:  

$$\mathbb{P}(X_t = n | X_0 = 0) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}.$$

For most Markov chains one makes progress *without* solving the differential equations.  
 The interplay between the simulation method above and the distributional information here is exactly the interplay between viewing the Poisson process as a counting process ("Poisson counts") and a sequence of inter-arrival times ("Exponential gaps"). The classic relationships between Exponential, Poisson, Gamma and Geometric distributions are all embedded in this one process.

Two significant extra facts are  
**superposition**: independent sum of Poisson processes is Poisson:  
**thinning**: if arrivals are censored i.i.d. at random then result is Poisson.

## Example: the M/M/1 queue

Consider a queue in which people arrive and are served (in order) at constant rates by a single server.

1. Stipulate that the number  $X_t$  of people in system at time  $t$  forms a Markov chain.
2. Transition rates (I): people arrive one-at-a-time at constant rate, so  $q(x, x + 1) = \lambda$ .
3. Transition rates (II): people are served in order at constant rate, so  $q(x, x - 1) = \mu$  if  $x > 0$ .

One can solve the equilibrium equations to deduce: the equilibrium distribution of  $X$  exists and is Geometric if and only if  $\lambda < \mu$ .



Example: the M/M/1 queue  
 Consider a queue in which people arrive and are served (in order) at constant rates by a single server.  
 1. Stipulate that the number  $X_t$  of people in system at time  $t$  forms a Markov chain.  
 2. Transition rates (I): people arrive one-at-a-time at constant rate, so  $q(x, x + 1) = \lambda$ .  
 3. Transition rates (II): people are served in order at constant rate, so  $q(x, x - 1) = \mu$  if  $x > 0$ .  
 One can solve the equilibrium equations to deduce: the equilibrium distribution of  $X$  exists and is Geometric if and only if  $\lambda < \mu$ .

**Don't** try to solve the equilibrium equations at home (unless you enjoy that sort of thing). In this case it is do-able, but during the module we'll discuss a much quicker way to find the equilibrium distribution in favourable cases.

Here is the equilibrium distribution in more explicit form: in equilibrium

$$\mathbb{P}(X = x) = (1 - \rho)\rho^x \quad \text{for } x = 0, 1, \dots,$$

where  $\rho = \lambda/\mu \in (0, 1)$  (the traffic intensity).

## Some useful texts (I)

At increasing levels of mathematical sophistication:

1. Haggström (2002) "Finite Markov chains and algorithmic applications".
2. Grimmett and Stirzaker (2001) "Probability and random processes".
3. Norris (1998) "Markov chains".
4. Williams (1991) "Probability with martingales".



Some useful texts (I)  
 At increasing levels of mathematical sophistication:  
 1. Haggström (2002) "Finite Markov chains and algorithmic applications".  
 2. Grimmett and Stirzaker (2001) "Probability and random processes".  
 3. Norris (1998) "Markov chains".  
 4. Williams (1991) "Probability with martingales".

1. Delightful introduction to finite state-space discrete-time Markov chains, from point of view of computer algorithms.
2. Standard undergraduate text on mathematical probability. This is the book I advise my students to buy, because it contains so much material.
3. Markov chains at a more graduate level of sophistication, revealing what I have concealed, namely the full gory story about  $Q$ -matrices.
4. Excellent graduate text for theory of martingales: mathematically demanding.

## Some useful texts (II): free on the web

1. Doyle and Snell (1984) "Random walks and electric networks" available on web at <http://arxiv.org/abs/math/0001057>.
2. Kindermann and Snell (1980) "Markov random fields and their applications" available on web at [http://www.ams.org/online\\_bks/conm1/](http://www.ams.org/online_bks/conm1/).
3. Meyn and Tweedie (1993) "Markov chains and stochastic stability" available on web at <http://probability.ca/MT/>.
4. Aldous and Fill (2001) "Reversible Markov Chains and Random Walks on Graphs" *only* available on web at <http://www.stat.berkeley.edu/~aldous/RWG/book.html>.



## ↳ Some useful texts (II): free on the web

1. Lays out (in simple and accessible terms) an important approach to Markov chains using relationship to resistance in electrical networks.
2. Sublimely accessible treatment of Markov random fields (Markov property, but in space not time).
3. The place to go if you need to get informed about theoretical results on rates of convergence for Markov chains (eg, because you are doing MCMC).
4. The best unfinished book on Markov chains known to me.

Some useful texts (II): free on the web

1. Doyle and Snell (1984) "Random walks and electric networks" available on web at <http://arxiv.org/abs/math/0001057>.
2. Kindermann and Snell (1980) "Markov random fields and their applications" available on web at [http://www.ams.org/online\\_bks/conm1/](http://www.ams.org/online_bks/conm1/).
3. Meyn and Tweedie (1993) "Markov chains and stochastic stability" available on web at <http://probability.ca/MT/>.
4. Aldous and Fill (2001) "Reversible Markov Chains and Random Walks on Graphs" only available on web at <http://www.stat.berkeley.edu/~aldous/RWG/book.html>.

## Some useful texts (III): going deeper

1. Kingman (1993) "Poisson processes".
2. Kelly (1979) "Reversibility and stochastic networks".
3. Lindvall (1992) "Lectures on the coupling method".
4. Steele (2004) "The Cauchy-Schwarz master class".
5. Aldous (1989) "Probability approximations via the Poisson clumping heuristic" see [www.stat.berkeley.edu/~aldous/Research/research80.html](http://www.stat.berkeley.edu/~aldous/Research/research80.html).
6. Øksendal (2003) "Stochastic differential equations".
7. Stoyan, Kendall, and Mecke (1987) "Stochastic geometry and its applications".



## ↳ Some useful texts (III): going deeper

Here are a few of the many texts which go much further

1. Very good introduction to the wide circle of ideas surrounding the Poisson process.
2. We'll cover reversibility briefly in the lectures, but this shows just how powerful the technique is.
3. We'll also talk briefly about the beautiful concept of coupling for Markov chains; this book gives a very nice introduction.
4. The book to read if you decide you need to know more about (mathematical) inequality.
5. A book full of what *ought* to be true; hence good for stimulating research problems and also for ways of computing heuristic answers.
6. An accessible introduction to Brownian motion and stochastic calculus, which we do not cover at all.
7. Discusses a range of techniques used to handle probability in geometric contexts.

Some useful texts (III): going deeper

1. Kingman (1993) "Poisson processes".
2. Kelly (1979) "Reversibility and stochastic networks".
3. Lindvall (1992) "Lectures on the coupling method".
4. Steele (2004) "The Cauchy-Schwarz master class".
5. Aldous (1989) "Probability approximations via the Poisson clumping heuristic" see <http://www.stat.berkeley.edu/~aldous/Research/research80.html>.
6. Øksendal (2003) "Stochastic differential equations".
7. Stoyan, Kendall, and Mecke (1987) "Stochastic geometry and its applications".

Aldous, D. J. (1989). *Probability approximations via the Poisson clumping heuristic*, Volume 77 of *Applied Mathematical Sciences*. New York: Springer-Verlag.

Aldous, D. J. and J. A. Fill (2001). *Reversible Markov Chains and Random Walks on Graphs*. Unpublished.

Cox, D. R. and H. D. Miller (1965). *The theory of stochastic processes*. New York: John Wiley & Sons Inc.

Doyle, P. G. and J. L. Snell (1984). *Random walks and electric networks*, Volume 22 of *Carus Mathematical Monographs*. Washington, DC: Mathematical Association of America.

Gardner, M. (1996). Word ladders: Lewis Carroll's doublets. *The Mathematical Gazette* 80(487), 195–198.



Grimmett, G. R. and D. R. Stirzaker (2001). *Probability and random processes* (Third ed.). New York: Oxford University Press.

Haggström, O. (2002). *Finite Markov chains and algorithmic applications*, Volume 52 of *London Mathematical Society Student Texts*. Cambridge: Cambridge University Press.

Kelly, F. P. (1979). *Reversibility and stochastic networks*. Chichester: John Wiley & Sons Ltd. Wiley Series in Probability and Mathematical Statistics.

Kindermann, R. and J. L. Snell (1980). *Markov random fields and their applications*, Volume 1 of *Contemporary Mathematics*. Providence, R.I.: American Mathematical Society.

Kingman, J. F. C. (1993). *Poisson processes*, Volume 3 of *Oxford Studies in Probability*. New York: The Clarendon Press Oxford University Press. Oxford Science Publications.



- Knuth, D. E. (1993).  
*The Stanford GraphBase: a platform for combinatorial computing.*  
New York, NY, USA: ACM.
- Lindvall, T. (1992).  
*Lectures on the coupling method.*  
Wiley Series in Probability and Mathematical Statistics: Probability and  
Mathematical Statistics. New York: John Wiley & Sons Inc.
- Meyn, S. P. and R. L. Tweedie (1993).  
*Markov chains and stochastic stability.*  
Communications and Control Engineering Series. London: Springer-Verlag London  
Ltd.
- Norris, J. R. (1998).  
*Markov chains*, Volume 2 of *Cambridge Series in Statistical and Probabilistic  
Mathematics.*  
Cambridge: Cambridge University Press.  
Reprint of 1997 original.
- Øksendal, B. (2003).  
*Stochastic differential equations* (Sixth ed.).  
Universitext. Berlin: Springer-Verlag.  
An introduction with applications.



- Steele, J. M. (2004).  
*The Cauchy-Schwarz master class.*  
MAA Problem Books Series. Washington, DC: Mathematical Association of  
America.  
An introduction to the art of mathematical inequalities.
- Stoyan, D., W. S. Kendall, and J. Mecke (1987).  
*Stochastic geometry and its applications.*  
Wiley Series in Probability and Mathematical Statistics: Applied Probability and  
Statistics. Chichester: John Wiley & Sons Ltd.  
With a foreword by D. G. Kendall.
- Williams, D. (1991).  
*Probability with martingales.*  
Cambridge Mathematical Textbooks. Cambridge: Cambridge University Press.

