

APTS Applied Stochastic Processes

Wilfrid Kendall
w.s.kendall@warwick.ac.uk

Department of Statistics, University of Warwick

12th July 2008

Introduction

- 1: Markov chains and reversibility
- 2: Martingales
- 3: Stopping times
- 4: Counting and compensating
- 5: Central Limit Theorem
- 6: Recurrence
- 7: Foster-Lyapunov criteria
- 8: Cutoff

Introduction

"... you never learn anything unless you are willing to take a risk and tolerate a little randomness in your life."
- Heinz Pagels,
The Dreams of Reason, 1988.

This module is intended to 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 significant ways in which these notions are used in statistics.

These notes outline the content of the module; they represent work-in-progress and will grow, be corrected, and be modified as the module lectures progress.

Introduction

... you never learn anything unless you are willing to take a risk and tolerate a little randomness in your life."
- Heinz Pagels, 1988

This module is intended to 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 significant ways in which these notions are used in statistics.

These notes outline the content of the module; they represent work-in-progress and will grow, be corrected, and be modified as the module lectures progress.

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 and to build their own mental map of the subject. The *Applied Stochastic Processes* module aims to contribute towards this end.

Corrections and suggestions are of course welcome! Email w.s.kendall@warwick.ac.uk. All images in these notes either are constructed by the author or have been released into the public domain.

Learning Outcomes

After successfully completing this module an APTS student will be able to:

- ▶ describe and calculate with the notion of a reversible Markov chain, both in discrete and continuous time;
- ▶ describe the basic properties of discrete-parameter martingales and check whether the martingale property holds;
- ▶ recall and apply some significant concepts from martingale theory;
- ▶ explain how to use Foster-Lyapunov criteria to establish recurrence and speed of convergence to equilibrium for Markov chains.

Learning Outcomes

After successfully completing this module an APTS student will be able to:

- ▶ describe and calculate with the notion of a reversible Markov chain, both in discrete and continuous time;
- ▶ describe the basic properties of discrete-parameter martingales and check whether the martingale property holds;
- ▶ recall and apply some significant concepts from martingale theory;
- ▶ explain how to use Foster-Lyapunov criteria to establish recurrence and speed of convergence to equilibrium for Markov chains.

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.

First of all, read the preliminary notes . . .

They provide notes and examples concerning a basic framework covering:

- ▶ Probability and conditional probability;
- ▶ Expectation and conditional expectation;
- ▶ Discrete-time countable-state-space Markov chains;
- ▶ Continuous-time countable-state-space Markov chains;
- ▶ Poisson processes.

First of all, read the preliminary notes . . .

They provide notes and examples concerning a basic framework covering:

- ▶ Probability and conditional probability;
- ▶ Expectation and conditional expectation;
- ▶ Discrete-time countable-state-space Markov chains;
- ▶ Continuous-time countable-state-space Markov chains;
- ▶ Poisson processes.

The purpose of the preliminary notes is not to provide all the information you might require concerning probability, but to serve as a prompt about material you may need to revise, and to introduce and to establish some basic choices of notation.

Some useful texts (I)

"There is no such thing as a moral or an immoral book. Books are well written or badly written."
- Oscar Wilde (1854-1900),
The Picture of Dorian Gray, 1891, preface

The next three slides list various useful textbooks.
At increasing levels of mathematical sophistication:

1. Häggström (2002) "Finite Markov chains and algorithmic applications".
2. Grimmett and Stirzaker (2001) "Probability and random processes".
3. Breiman (1992) "Probability".
4. Norris (1998) "Markov chains".
5. Williams (1991) "Probability with martingales".

Some useful texts (I)
There is no such thing as a moral or an immoral book. Books are well written or badly written.
- Oscar Wilde (1854-1900),
The Picture of Dorian Gray, 1891, preface
The next three slides list various useful textbooks.
At increasing levels of mathematical sophistication:
1. Häggström (2002) "Finite Markov chains and algorithmic applications".
2. Grimmett and Stirzaker (2001) "Probability and random processes".
3. Breiman (1992) "Probability".
4. Norris (1998) "Markov chains".
5. Williams (1991) "Probability with martingales".

1. Häggström (2002) is a delightful introduction to finite state-space discrete-time Markov chains, from point of view of computer algorithms.
2. Grimmett and Stirzaker (2001) is the standard undergraduate text on mathematical probability. This is the book I advise my undergraduate students to buy, because it contains so much material.
3. Breiman (1992) is a first-rate graduate-level introduction to probability.
4. Norris (1998) presents the theory of Markov chains at a more graduate level of sophistication, revealing what I have concealed, namely the full gory story about Q -matrices.
5. Williams (1991) provides an excellent graduate treatment 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 www.arxiv.org/abs/math/0001057.
2. Kindermann and Snell (1980) "Markov random fields and their applications" available on web at www.ams.org/online_bks/conm1/.
3. Meyn and Tweedie (1993) "Markov chains and stochastic stability" available on web at www.probability.ca/MT/.
4. Aldous and Fill (2001) "Reversible Markov Chains and Random Walks on Graphs" *only* available on web at www.stat.berkeley.edu/~aldous/RWC/book.html.

Some useful texts (II): free on the web
1. Doyle and Snell (1984) "Random walks and electric networks" available on web at www.arxiv.org/abs/math/0001057.
2. Kindermann and Snell (1980) "Markov random fields and their applications" available on web at www.ams.org/online_bks/conm1/.
3. Meyn and Tweedie (1993) "Markov chains and stochastic stability" available on web at www.probability.ca/MT/.
4. Aldous and Fill (2001) "Reversible Markov Chains and Random Walks on Graphs" only available on web at www.stat.berkeley.edu/~aldous/RWC/book.html.

1. Doyle and Snell (1984) lays out (in simple and accessible terms) an important approach to Markov chains using relationship to resistance in electrical networks.
2. Kindermann and Snell (1980) is a sublimely accessible treatment of Markov random fields (Markov property, but in space not time).
3. Consult Meyn and Tweedie (1993) if you need to get informed about theoretical results on rates of convergence for Markov chains (eg, because you are doing MCMC).
4. Aldous and Fill (2001) is the best unfinished book on Markov chains known to me (at the time of writing these notes).

Some useful texts (III): going deeper

1. Kingman (1993) "Poisson processes".
2. Kelly (1979) "Reversibility and stochastic networks".
3. Steele (2004) "The Cauchy-Schwarz master class".
4. Aldous (1989) "Probability approximations via the Poisson clumping heuristic".
5. Øksendal (2003) "Stochastic differential equations".
6. Stoyan, Kendall, and Mecke (1995) "Stochastic geometry and its applications".

Some useful texts (III): going deeper
1. Kingman (1993) "Poisson processes".
2. Kelly (1979) "Reversibility and stochastic networks".
3. Steele (2004) "The Cauchy-Schwarz master class".
4. Aldous (1989) "Probability approximations via the Poisson clumping heuristic".
5. Øksendal (2003) "Stochastic differential equations".
6. Stoyan, Kendall, and Mecke (1995) "Stochastic geometry and its applications".

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

1. Kingman (1993) gives a very good introduction to the wide circle of ideas surrounding the Poisson process.
2. We'll cover reversibility briefly in the lectures, but Kelly (1979) shows just how powerful the technique can be.
3. Steele (2004) is the book to read if you decide you need to know more about (mathematical) inequality.
4. Aldous (1989) is a book full of what *ought* to be true; hence good for stimulating research problems and also for ways of computing heuristic answers. See www.stat.berkeley.edu/~aldous/Research/research80.html.
5. Øksendal (2003) is an accessible introduction to Brownian motion and stochastic calculus, which we do not cover at all.
6. Stoyan et al. (1995) discusses a range of techniques used to handle probability in geometric contexts.

Markov chains and reversibility

"People assume that time is a strict progression of cause to effect, but actually from a non-linear, non-subjective viewpoint, it's more like a big ball of wibbly-wobbly, timey-wimey ... stuff."
The Tenth Doctor,
Doctor Who, in the episode "Blink", 2007

Introduction and simplest non-trivial example
Birth, death and immigration
Detailed balance definition and theorem
M/M/1 queue
Random chess
Ising model



Markov chains and reversibility
"People assume that time is a strict progression of cause to effect, but actually from a non-linear, non-subjective viewpoint, it's more like a big ball of wibbly-wobbly, timey-wimey ... stuff."
The Tenth Doctor,
Doctor Who, in the episode "Blink", 2007
Introduction and simplest non-trivial example
Birth, death and immigration
Detailed balance definition and theorem
M/M/1 queue
Random chess
Ising model

We begin our module with the important, simple and subtle idea of a **reversible** Markov chain, and the associated notion of **detailed balance**; we will return to these ideas periodically through the module. This first major theme isolates a class of Markov chains for which computation of the equilibrium distribution is relatively straightforward. (Remember from the pre-requisites: if a chain is irreducible and positive-recurrent then it has an equilibrium distribution $\underline{\pi}$; and if it is aperiodic then $\underline{\pi}$ is also the limiting long-time empirical distribution. Moreover $\underline{\pi} \cdot P = \underline{\pi}$. However if there are k states then these matrix equation presents k equations each potentially involving all k unknowns ... a complexity issue if k is large!)

- 1: Markov chains and reversibility
- Introduction and simplest non-trivial example

Markov chains and reversibility

Here is detailed balance in a nutshell:
 Suppose we could solve for $\underline{\pi}$ in $\pi_x p_{xy} = \pi_y p_{yx}$ (discrete-time) or $\pi_x q_{xy} = \pi_y q_{yx}$ (continuous-time). In both cases simple algebra then shows $\underline{\pi}$ solves the equilibrium equations.
 So on a prosaic level it is always worth trying this easy route; if the detailed balance equations are insoluble then revert to the more complicated equilibrium equations $\underline{\pi} \cdot \underline{Q} = \underline{\pi}$, respectively $\underline{\pi} \cdot \underline{Q} = 0$.
 We will consider reversibility of Markov chains in both discrete and continuous time, the computation of equilibrium distributions for such chains, and application to some illustrative examples.

- 1: Markov chains and reversibility
- Introduction and simplest non-trivial example
- Markov chains and reversibility

We will consider:

- simple symmetric random walk;
- the birth-death-immigration process;
- the M/M/1 queue;
- a discrete-time chain on a 8×8 state space;
- Gibbs' samplers (briefly);
- and Metropolis-Hastings samplers (briefly).

Test understanding: show the detailed balance equations (discrete-case) lead to equilibrium equations by applying them and then $\sum_x p_{yx} = 1$ to $\sum_x \pi_x p_{xy}$.

Markov chains and reversibility
 Here is detailed balance in a nutshell:
 Suppose we could solve for $\underline{\pi}$ in $\pi_x p_{xy} = \pi_y p_{yx}$ (discrete-time) or $\pi_x q_{xy} = \pi_y q_{yx}$ (continuous-time). In both cases simple algebra then shows $\underline{\pi}$ solves the equilibrium equations.
 So on a prosaic level it is always worth trying this easy route; if the detailed balance equations are insoluble then revert to the more complicated equilibrium equations $\underline{\pi} \cdot \underline{Q} = \underline{\pi}$, respectively $\underline{\pi} \cdot \underline{Q} = 0$.
 We will consider reversibility of Markov chains in both discrete and continuous time, the computation of equilibrium distributions for such chains, and application to some illustrative examples.

- 1: Markov chains and reversibility
- Introduction and simplest non-trivial example

Simplest non-trivial example (I)

Consider **doubly-reflected simple symmetric random walk** X on $\{0, 1, \dots, k\}$, with reflection "by prohibition": moves $0 \rightarrow -1, k \rightarrow k+1$ are replaced by $0 \rightarrow 0, k \rightarrow k$. **ANIMATION**

1. X is **irreducible** and **aperiodic**, so there is a unique equilibrium distribution $\underline{\pi} = (\pi_0, \pi_1, \dots, \pi_k)$.
2. The **equilibrium equations** $\underline{\pi} \cdot \underline{P} = \underline{\pi}$ are solved by $\pi_i = \frac{1}{k+1}$ for all i .
3. Consider X in equilibrium and **run backwards in time**. Calculation then shows, $\mathbb{P}[X_{n-1} = x | X_n = y] = \pi_x \mathbb{P}[X_n = y | X_{n-1} = x] / \pi_y = \mathbb{P}[X_n = y | X_{n-1} = x]$ so in this case **by symmetry of the kernel** the equilibrium chain has the same transition kernel (so looks the same) whether run forwards or backwards in time.

- 1: Markov chains and reversibility
- Introduction and simplest non-trivial example
- Simplest non-trivial example (I)

1. **Test understanding:** explain why X is aperiodic when *non-reflected* simple symmetric random walk has period 2.
2. **Test understanding:** verify solution of equilibrium equations.
3. Develop Markov property to deduce X_0, X_1, \dots, X_{n-1} is conditionally independent of X_{n+1}, X_{n+2}, \dots given X_n . Hence reversed Markov chain is *still* Markov (though not necessarily time-homogeneous in more general circumstances). Suppose the reversed chain has kernel $\bar{p}_{y,x}$.
 - Use definition of conditional probability to compute $\bar{p}_{y,x} = \mathbb{P}[X_{n-1} = x, X_n = y] / \mathbb{P}[X_n = y]$,
 - then $\mathbb{P}[X_{n-1} = x, X_n = y] / \mathbb{P}[X_n = y] = \mathbb{P}[X_{n-1} = x | X_n = y]$.
 - now substitute, using $\mathbb{P}[X_n = i] = \frac{1}{k+1}$ for all i so $\bar{p}_{y,x} = p_{x,y}$.
 - Symmetry of kernel ($p_{x,y} = p_{y,x}$) then shows backwards kernel $\bar{p}_{y,x}$ is same as forwards kernel $p_{y,x} = p_{y,x}$.
 The construction generalizes ... so the link between reversibility and detailed balance holds generally.

Simplest non-trivial example (I)
 Consider doubly-reflected simple symmetric random walk X on $\{0, 1, \dots, k\}$, with reflection "by prohibition": moves $0 \rightarrow -1, k \rightarrow k+1$ are replaced by $0 \rightarrow 0, k \rightarrow k$.
 1. X is irreducible and aperiodic, so there is a unique equilibrium distribution $\underline{\pi} = (\pi_0, \pi_1, \dots, \pi_k)$.
 2. The equilibrium equations $\underline{\pi} \cdot \underline{P} = \underline{\pi}$ are solved by $\pi_i = \frac{1}{k+1}$ for all i .
 3. Consider X in equilibrium and run backwards in time. Calculation then shows, $\mathbb{P}[X_{n-1} = x | X_n = y] = \pi_x \mathbb{P}[X_n = y | X_{n-1} = x] / \pi_y = \mathbb{P}[X_n = y | X_{n-1} = x]$ so in this case by symmetry of the kernel the equilibrium chain has the same transition kernel (so looks the same) whether run forwards or backwards in time.

- 1: Markov chains and reversibility
- Introduction and simplest non-trivial example

Simplest non-trivial example (II)

There is a computational aspect to this.

1. Even in more general cases, if the π_i depend on i then above computations show reversibility holds if equilibrium distribution exists and **equations of detailed balance** hold: $\pi_x p_{x,y} = \pi_y p_{y,x}$.
2. Moreover if one can solve for π_i in $\pi_x p_{x,y} = \pi_y p_{y,x}$ then it is easy to show $\underline{\pi} \cdot \underline{P} = \underline{\pi}$.
3. Consequently if one can solve the equations of detailed balance, and if the solution can be normalized to have unit total probability, then the result also solves the equilibrium equations.

- 1: Markov chains and reversibility
- Introduction and simplest non-trivial example
- Simplest non-trivial example (II)

1. **Test understanding:** check this.
2. **Test understanding:** check this.
3. Even in this simple example there is an evident improvement in complexity. Detailed balance involves k equations each with two unknowns, easily "chained together". The equilibrium equations involve k equations of which $k-2$ involve three unknowns.
 In general the detailed balance equations can be solved unless "chaining together by different routes" delivers inconsistent results. Kelly (1979) goes into more detail about this.
Test understanding: show detailed balance doesn't work for 3-state chain with transition probabilities $\frac{1}{3}$ for $0 \rightarrow 1, 1 \rightarrow 2, 2 \rightarrow 0$ and $\frac{2}{3}$ for $2 \rightarrow 1, 1 \rightarrow 0, 0 \rightarrow 2$.
Test understanding: show detailed balance *does* work for doubly reflected *asymmetric* simple random walk.
 We will see there are still major computational issues for more general Markov chains, connected with determining the normalizing constant to ensure $\sum_i \pi_i = 1$.

Simplest non-trivial example (II)
 There is a computational aspect to this.
 1. Even in more general cases, if the π_i depend on i then above computations show reversibility holds if equilibrium distribution exists and equations of detailed balance hold: $\pi_x p_{x,y} = \pi_y p_{y,x}$.
 2. Moreover if one can solve for π_i in $\pi_x p_{x,y} = \pi_y p_{y,x}$ then it is easy to show $\underline{\pi} \cdot \underline{P} = \underline{\pi}$.
 3. Consequently if one can solve the equations of detailed balance, and if the solution can be normalized to have unit total probability, then the result also solves the equilibrium equations.

- 1: Markov chains and reversibility
- Birth, death and immigration

Birth-death-immigration process

The same idea works for continuous-time Markov chains: replace transition probabilities $p_{x,y}$ by rates $q_{x,y}$ and equilibrium equation $\underline{\pi} \cdot \underline{Q} = \underline{\pi}$ by differentiated variant using Q -matrix: $\underline{\pi} \cdot \underline{Q} = 0$.

Definition

The birth-death-immigration process has transitions:

- ▶ Birth ($X \rightarrow X+1$ at rate λX);
- ▶ Death ($X \rightarrow X-1$ at rate μX);
- ▶ plus an extra **Immigration** term ($X \rightarrow X+1$ at rate α).

Hence $q_{x,x+1} = \lambda x + \alpha; q_{x,x-1} = \mu x$. **ANIMATION**

Equilibrium is derived easily from detailed balance:

$$\pi_x = \frac{\lambda(x-1) + \alpha}{\mu x} \cdot \frac{\lambda(x-2) + \alpha}{\mu(x-1)} \cdot \dots \cdot \frac{\alpha}{\mu} \cdot \pi_0$$

- 1: Markov chains and reversibility
- Birth, death and immigration
- Birth-death-immigration process

Reversibility here is decidedly non-trivial ... We need $0 \leq \lambda < \mu$ and $\alpha > 0$. Note that for this population process the rates $q_{x,x\pm 1}$ make sense and are defined only for $x = 0, 1, 2, \dots$
 Detailed balance equations:

$$\pi_x \times \mu x = \pi_{x-1} \times (\lambda(x-1) + \alpha)$$

Test understanding: check the calculations!
 Normalizing constant can be computed exactly when $\lambda < \mu$ via

$$\pi_0^{-1} = \sum_{x=0}^{\infty} \frac{\lambda(x-1) + \alpha}{\mu x} \cdot \frac{\lambda(x-2) + \alpha}{\mu(x-1)} \cdot \dots \cdot \frac{\alpha}{\mu} = \left(\frac{\mu}{\mu - \lambda} \right)^{\frac{\alpha}{\lambda}}$$

If the condition $\lambda < \mu$ is not satisfied then the sum does not converge and therefore there can be **no** equilibrium!
 If $\alpha = 0$ then equilibrium = extinction ...

Poisson process: $\lambda = \mu = 0$.

Birth-death-immigration process
 The same idea works for continuous-time Markov chains: replace transition probabilities $p_{x,y}$ by rates $q_{x,y}$ and equilibrium equation $\underline{\pi} \cdot \underline{Q} = \underline{\pi}$ by differentiated variant using Q -matrix: $\underline{\pi} \cdot \underline{Q} = 0$.
Definition
 The birth-death-immigration process has transitions:
 - Birth ($X \rightarrow X+1$ at rate λX);
 - Death ($X \rightarrow X-1$ at rate μX);
 - plus an extra **Immigration** term ($X \rightarrow X+1$ at rate α).
 Hence $q_{x,x+1} = \lambda x + \alpha; q_{x,x-1} = \mu x$.
 Equilibrium is derived easily from detailed balance:

$$\pi_x \cdot \mu x = \pi_{x-1} \cdot (\lambda(x-1) + \alpha)$$

Detailed balance and reversibility

Definition

The Markov chain X satisfies **detailed balance** if

Discrete time: there is a non-trivial solution of

$$\pi_x p_{x,y} = \pi_y p_{y,x};$$

Continuous time: there is a non-trivial solution of

$$\pi_x q_{x,y} = \pi_y q_{y,x}.$$

Theorem

The irreducible Markov chain X satisfies **detailed balance** and the solution $\{\pi_x\}$ can be normalized by $\sum_x \pi_x = 1$ if and only if $\{\pi_x\}$ is an equilibrium distribution for X and X started in equilibrium is statistically the same whether run forwards or backwards in time.

Detailed balance and reversibility

Definition
 The Markov chain X satisfies detailed balance if there is a non-trivial solution of $\pi_x p_{x,y} = \pi_y p_{y,x}$.
Question: there is a non-trivial solution of $\pi_x p_{x,y} = \pi_y p_{y,x}$.

Theorem
 The irreducible Markov chain X satisfies detailed balance and the solution $\{\pi_x\}$ can be normalized by $\sum_x \pi_x = 1$ if and only if $\{\pi_x\}$ is an equilibrium distribution for X and X started in equilibrium is statistically the same whether run forwards or backwards in time.

1. Proof of theorem is routine: see example of random walk above.
2. The reversibility phenomenon has surprisingly deep ramifications! Consider birth-death-immigration example above and ask yourself whether it is apparent that the time-reversed process in equilibrium looks statistically the same as the original process. (Note: both immigrations and births convert to deaths, and vice versa ...)

M/M/1 queue

Here we have

- ▶ Arrivals: $X \rightarrow X + 1$ at rate λ ;
- ▶ Departures: $X \rightarrow X - 1$ at rate μ if $X > 0$.

Hence detailed balance: $\mu \pi_x = \lambda \pi_{x-1}$ and therefore when $\lambda < \mu$ (stability) the equilibrium distribution is $\pi_x = \rho^x (1 - \rho)$ for $x = 0, 1, \dots$, where $\rho = \frac{\lambda}{\mu}$ (the traffic intensity). **ANIMATION**

Reversibility/detailed balance is more than a computational device: consider **Burke's theorem**, if a stable M/M/1 queue is in equilibrium then people leave according to a Poisson process of rate λ . Hence if a stable M/M/1 queue feeds into another stable M/M/1 queue then in equilibrium the second queue on its own behaves as an M/M/1 queue in equilibrium.

M/M/1 queue

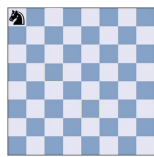
Here we have
 • Arrivals: $X \rightarrow X + 1$ at rate λ ;
 • Departures: $X \rightarrow X - 1$ at rate μ if $X > 0$.
 Hence detailed balance $\mu \pi_x = \lambda \pi_{x-1}$ and therefore when $\lambda < \mu$ (stability) the equilibrium distribution is $\pi_x = \rho^x (1 - \rho)$ for $x = 0, 1, \dots$ where $\rho = \frac{\lambda}{\mu}$ is the traffic intensity.
Burke's theorem: if a stable M/M/1 queue is in equilibrium then people leave according to a Poisson process of rate λ .
 Hence if a stable M/M/1 queue feeds into another stable M/M/1 queue then in equilibrium the second queue on its own behaves as an M/M/1 queue in equilibrium.

We recall the M/M/1 queue example discussed in the preliminary notes. Birth-death-immigration processes and queueing processes are examples of **generalized birth-death processes**; only $X \rightarrow X \pm 1$ transitions, hence detailed balance equations easily solved. Note: the M/M/1 queue is **non-linear**. Linearity allows solution of forwards equations: we do not discuss this here. Detailed balance is also a subtle and important tool for the study of Markovian queueing networks (e.g. Kelly 1979). The argument connecting reversibility to detailed balance runs both ways. If detailed balance equations can be solved to derive equilibrium then the process is reversible if run in equilibrium. Hence a one-line proof of Burke's theorem: if queue is run backwards in time then departures become arrivals. **Test understanding:** use Burke's theorem for a feed-forward M/M/1 queueing network (no loops) to show that in equilibrium each queue viewed in isolation is M/M/1. This uses the fact that independent thinnings and superpositions of Poisson processes are still Poisson ...

Random chess (Aldous and Fill 2001, Ch1, Ch3§2)

Example (A mean Knight's tour)

Place a chess Knight at the corner of a standard 8×8 chessboard. Move it randomly, at each move choosing uniformly from available legal chess moves independently of the past.



ANIMATION

1. What is the equilibrium distribution? (use detailed balance)
2. Is the resulting Markov chain periodic? (what if you sub-sample at even times?)
3. What is the mean time till the Knight returns to its starting point? (inverse of equilibrium probability)

Random chess (Aldous and Fill 2001, Ch1, Ch3§2)

Example (A mean Knight's tour)
 Place a chess Knight at the corner of a standard 8×8 chessboard. Move it randomly, at each move choosing uniformly from available legal chess moves independently of the past.

1. What is the equilibrium distribution? (use detailed balance)
2. Is the resulting Markov chain periodic? (what if you sub-sample at even times?)
3. What is the mean time till the Knight returns to its starting point? (inverse of equilibrium probability)

1. Use $\pi_v/d_v = \pi_u/d_u$ if $u \leftrightarrow v$, where d_u is the degree of u . Also use fact, there are $168 = (2 + 2 \times 3 + 5 \times 4 + 4 \times 6 + 6 \times 8) \times 4/2$ different edges. So total degree is 2×168 and equilibrium probability at corner is $2/(2 \times 168)$.
2. Period 2 (white versus black). Sub-sampling at even times makes chain aperiodic on squares of one colour.
3. Inverse of equilibrium probability shows that mean return time to corner is 168.

Gibbs' sampler for Ising model

(I) Ising model

- ▶ Pattern of spins $S_j = \pm 1$ on (finite fragment of) lattice (i is typical node of lattice).
- ▶ Probability mass function

$$\mathbb{P}[S_i = s_i \text{ all } i] \propto \begin{cases} \exp\left(J \sum_{i \sim j} s_i s_j\right), \\ \exp\left(J \sum_{i \sim j} s_i s_j + H \sum_i s_i \tilde{s}_i\right) \\ \text{if external field } \tilde{s}_i. \end{cases}$$

Gibbs' sampler for Ising model

• Pattern of spins $S_i = \pm 1$ on (finite fragment of) lattice (i is typical node of lattice).
 • Probability mass function

$$\mathbb{P}[S_i = s_i \text{ all } i] \propto \begin{cases} \exp\left(J \sum_{i \sim j} s_i s_j\right), \\ \exp\left(J \sum_{i \sim j} s_i s_j + H \sum_i s_i \tilde{s}_i\right) \\ \text{if external field } \tilde{s}_i. \end{cases}$$

1. Sample applications: idealized model for magnetism, simple binary image. Physics: interest in fragment expanding to fill whole lattice: cases of zero-interaction, sub-critical, critical ($\frac{KT}{J} = 2.269185$), super-critical. The Ising model is the nexus for a whole variety of scientific approaches, each bringing their own rather different questions.
2. $i \sim j$ if i and j are lattice neighbours. Note, physics treatments use a (physically meaningful) over-parametrization $J \rightarrow \frac{KT}{J}, H \rightarrow mH$. The $H \sum_i s_i \tilde{s}_i$ term can be interpreted physically as modelling an external magnetic field, or statistically as a noisy image conditioning the image. For a simulation physics view of the Ising model, see the expository article by David Landau in Kendall et al. (2005).
3. Actually computing the normalizing constant here is **hard** in the sense of complexity theory (see for example Jerrum 2003).

Gibbs' sampler for Ising model

(II) Gibbs' sampler (or heat-bath)

- Consider Markov chain with states which are Ising configurations on an $n \times n$ lattice, moving as follows:
 - Set \underline{s} to be a given configuration, with $\underline{s}^{(i)}$ obtained by flipping spin i ,
 - Choose a site i in the lattice at random;
 - Compute the conditional probability $\mathbb{P}[\underline{s}^{(i)} | \{\underline{s}^{(j)}, \underline{s}\}]$ of current configuration given configuration at other sites;
 - Flip the current value of S_i with probability $\mathbb{P}[\underline{s}^{(i)} | \{\underline{s}^{(j)}, \underline{s}\}]$, otherwise leave unchanged.
- Simple general calculations show,

$$\sum_i \frac{1}{n^2} \mathbb{P}[\underline{s}^{(i)}] \times \mathbb{P}[\underline{s} | \{\underline{s}^{(j)}, \underline{s}\}] = \mathbb{P}[\underline{s}]$$

so chain has Ising model as equilibrium distribution.

Gibbs' sampler for Ising model
 Consider Markov chain with states which are Ising configurations on an $n \times n$ lattice, moving as follows:
 - Set \underline{s} to be a given configuration, with $\underline{s}^{(i)}$ obtained by flipping spin i (at the lattice at random).
 - Choose a site i in the lattice at random.
 - Compute the conditional probability $\mathbb{P}[\underline{s}^{(i)} | \{\underline{s}^{(j)}, \underline{s}\}]$ of current configuration given configuration at other sites.
 - Flip the current value of S_i with probability $\mathbb{P}[\underline{s}^{(i)} | \{\underline{s}^{(j)}, \underline{s}\}]$, otherwise leave unchanged.
 - Simple general calculations show:

$$\sum_i \frac{1}{n^2} \mathbb{P}[\underline{s}^{(i)}] \times \mathbb{P}[\underline{s} | \{\underline{s}^{(j)}, \underline{s}\}] = \mathbb{P}[\underline{s}]$$

 so chain has Ising model as equilibrium distribution.

This is a *particular* example of the Gibbs' sampler in the special context of Ising models.

- Note that the configurations can be viewed as vectors of ± 1 's listing the various spins at different sites.
- In case of the Ising model, noting that $s_j^{(i)} = -s_j$,

$$\mathbb{P}[\underline{s} | \{\underline{s}^{(j)}, \underline{s}\}] \propto \frac{\exp\left(J \sum_{j \sim i} s_j s_i\right)}{\exp\left(J \sum_{j \sim i} s_j s_j\right) + \exp\left(-J \sum_{j \sim i} s_j s_j\right)}$$

- Obvious changes if external field.
- This is really a completely general computation!
 Note how complicated the equilibrium equations are: n^2 equations, each with n^2 terms on left-hand side.
- General pattern for Gibbs sampler: update individual random variables according to their conditional distributions given all other random variables.
- Conditional* distributions, so ratios, so normalizing constants cancel out.

Gibbs' sampler for Ising model

(III) Detailed balance

- Detailed balance calculations provide a much easier justification: merely check

$$\frac{1}{n^2} \mathbb{P}[\underline{s}^{(i)}] \times \mathbb{P}[\underline{s} | \{\underline{s}^{(j)}, \underline{s}\}] = \frac{1}{n^2} \mathbb{P}[\underline{s}] \times \mathbb{P}[\underline{s}^{(i)} | \{\underline{s}^{(j)}, \underline{s}\}]$$

- Here is an animation of a Gibbs' sampler producing an Ising model conditioned by a noisy image, produced by systematic scans: 128×128 , with 8 neighbours. Noisy image to left, draw from Ising model to right.



ANIMATION

Gibbs' sampler for Ising model
 Detailed balance calculations provide a much easier justification: merely check:

$$\frac{1}{n^2} \mathbb{P}[\underline{s}^{(i)}] \times \mathbb{P}[\underline{s} | \{\underline{s}^{(j)}, \underline{s}\}] = \frac{1}{n^2} \mathbb{P}[\underline{s}] \times \mathbb{P}[\underline{s}^{(i)} | \{\underline{s}^{(j)}, \underline{s}\}]$$

 Here is an animation of a Gibbs' sampler producing an Ising model conditioned by a noisy image, produced by systematic scans: 128×128 , with 8 neighbours. Noisy image to left, draw from Ising model to right.

- Test understanding:** check the detailed balance calculations. This also works for processes obtained from:
 - systematic scans
 - coding ("simultaneous updates on alternate colours of a chessboard") but *not* for wholly simultaneous updates.
- The example is taken from a discussion of "perfect simulation", but that is another story! See

www.warwick.ac.uk/go/wsk/ising-animations

for more on perfect sampling for the Ising model.

Metropolis-Hastings

- An important alternative to the Gibbs' sampler, even more closely connected to detailed balance:
 - Suppose $X_n = x$;
 - Pick y using a transition probability kernel $\alpha(x, y)$ (the *proposal kernel*);
 - accept* the proposed transition $x \rightarrow y$ with probability

$$q(x, y) = \min\left\{1, \frac{\pi(y)\alpha(y, x)}{\pi(x)\alpha(x, y)}\right\}$$

- if transition accepted, set $X_{n+1} = y$;
- otherwise set $X_{n+1} = x$.

- If π satisfies detailed balance then π is an equilibrium distribution.

ANIMATION

Metropolis-Hastings
 An important alternative to the Gibbs' sampler, even more closely connected to detailed balance:
 - Suppose $X_n = x$ using a transition probability kernel $\alpha(x, y)$ (the proposal kernel).
 - Pick y using $\alpha(x, y)$ with probability $\alpha(x, y)$.
 - If transition accepted, set $X_{n+1} = y$; otherwise set $X_{n+1} = x$.
 - If π satisfies detailed balance then π is an equilibrium distribution.

- Actually the Gibbs' sampler is a special case of the Metropolis-Hastings sampler.
- Test understanding:** write down the transition probability kernel for X .
Test understanding: check that π solves the detailed balance equations.
- Common variations on choice of proposal kernel:
 - independence sampler:* $\alpha(x, y) = f(y)$;
 - random-walk sampler:* $\alpha(x, y) = f(y - x)$;
 - Langevin sampler:* replace random-walk shift by shift depending on grad log π .
- Slice sampler: example of Metropolis-Hastings sampler used to deliver a uniform draw from region under the graph of a probability density function.
- Ratio* $\pi(x)/\pi(y)$, so normalizing constants cancel out.

Martingales

"One of these days... a guy is going to come up to you and show you a nice brand-new deck of cards on which the seal is not yet broken, and this guy is going to offer to bet you that he can make the Jack of Spades jump out of the deck and squirt cider in your ear. But, son, do not bet this man, for as sure as you are standing there, you are going to end up with an earful of cider."

Frank Loesser,
 Guys and Dolls musical, 1950, script



- Simplest possible example
- Thackeray's martingale
- Populations
- Definitions
- Finance example
- Martingales and likelihood
- Chicken Little

Martingales
 This is the second major theme of these notes: *martingales* are a class of random processes which are closely linked to ideas of conditional expectation. Briefly, martingales model your fortune if you are playing a fair game. (There are associated notions of "supermartingale", for a game unfair to you, and "submartingale", for a game fair to you.) But martingales can do so much more!

This is the second major theme of these notes: *martingales* are a class of random processes which are closely linked to ideas of conditional expectation. Briefly, martingales model your fortune if you are playing a fair game. (There are associated notions of "supermartingale", for a game unfair to you, and "submartingale", for a game fair to you.) But martingales can do so much more!

Martingales pervade modern probability

- We say the random process X is a martingale if it satisfies the **martingale property**:

$$\mathbb{E}[X_{n+1} | X_n, X_{n-1}, \dots] = \mathbb{E}[X_n \text{ plus jump at time } n | X_n, X_{n-1}, \dots] = X_n.$$

- Simplest possible example: simple symmetric random walk $X_0 = 0, X_1, X_2, \dots$. The martingale property follows from independence and distributional symmetry of jumps.
- For convenience and brevity, we often replace $\mathbb{E}[\dots | X_n, X_{n-1}, \dots]$ by $\mathbb{E}[\dots | \mathcal{F}_n]$ and think of "conditioning on \mathcal{F}_n " as "conditioning on all events which can be determined to have happened by time n ".

Martingales pervade modern probability

- We say the random process X is a martingale if it satisfies the martingale property: $\mathbb{E}[X_{n+1} | X_n, X_{n-1}, \dots] = X_n$. (It also jumps at time n , $X_n, X_{n-1}, \dots = X_n$.)
- Simplest possible example: simple symmetric random walk $X_0 = 0, X_1, X_2, \dots$. The martingale property follows from independence and distributional symmetry of jumps.
- For convenience and brevity, we often replace $\mathbb{E}[\dots | X_n, X_{n-1}, \dots]$ by $\mathbb{E}[\dots | \mathcal{F}_n]$ and think of "conditioning on \mathcal{F}_n " as "conditioning on all events which can be determined to have happened by time n ".

We use X as a convenient abbreviation for the stochastic process $\{X_n : n \geq 0\}$, et cetera.

- For a conversation with the inventor, see www.dartmouth.edu/~chance/Doob/conversation.html.
- Expected future level of X is current level.
- We use \mathcal{F}_n notation without comment in future, usually representing conditioning by X_0, X_1, \dots, X_n (if X is martingale in question). *Sometimes* further conditioning will be added; but \mathcal{F}_{n+1} always represents **at least as much** conditioning as \mathcal{F}_n . Crucially, the "Tower property" of conditional expectation then applies:

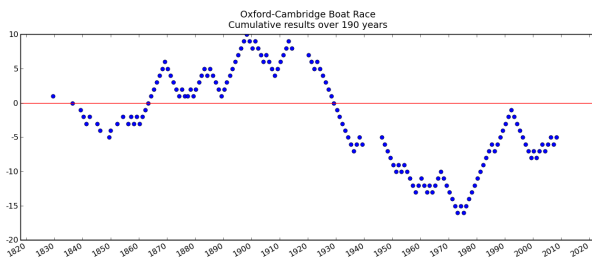
$$\mathbb{E}[\mathbb{E}[Z | \mathcal{F}_{n+1}] | \mathcal{F}_n] = \mathbb{E}[Z | \mathcal{F}_n].$$

Test understanding: deduce

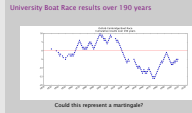
$$\mathbb{E}[X_{n+k} | \mathcal{F}_n] = X_n.$$

- There is an extensive theory about the notion of a *filtration of σ -algebras* or *σ -fields*, $\{\mathcal{F}_n : n \geq 0\}$. We avoid going into details ...

University Boat Race results over 190 years



Could this represent a martingale?



I first became aware of the boat race in about 1970, at which time the martingale property would have seemed not to apply.

There is now a much more satisfactory balance, but I still have my doubts as to the validity of the martingale property here ...



Thackeray's martingale

- MARTINGALE:**
 - spar under the bowsprit of a sailboat;
 - a harness strap that connects the nose piece to the girth; prevents the horse from throwing back its head.
- MARTINGALE in gambling:**
 The original sense is given in the OED: "a system in gambling which consists in doubling the stake when losing in the hope of eventually recouping oneself." The oldest quotation is from 1815 but the nicest is from 1854: Thackeray in *The Newcomes* I. 266 "You have not played as yet? Do not do so; above all avoid a martingale if you do."
 3. Result of playing Thackeray's martingale system and stopping on first win:
 set fortune at time n to be M_n .
 If $X_1 = -1, \dots, X_n = -n$ then
 $M_n = -1 - 2 - \dots - 2^{n-1} = 1 - 2^n$, otherwise $M_n = 1$.

Thackeray's martingale

- MARTINGALE:**
 - was under the bowsprit of a sailboat.
 - a harness strap that connects the nose piece to the girth; prevents the horse from throwing back its head.
- MARTINGALE in gambling:**
 The original sense is given in the OED: "a system in gambling which consists in doubling the stake when losing in the hope of eventually recouping oneself." The oldest quotation is from 1815 but the nicest is from 1854: Thackeray in *The Newcomes* I. 266 "You have not played as yet? Do not do so; above all avoid a martingale if you do."
 3. Result of playing Thackeray's martingale system and stopping on first win:
 set fortune at time n to be M_n .
 If $X_1 = -1, \dots, X_n = -n$ then
 $M_n = -1 - 2 - \dots - 2^{n-1} = 1 - 2^n$, otherwise $M_n = 1$.

- This is the "doubling" strategy. The equestrian meaning resembles the probabilistic definition to some extent.
- Notice how Thackeray's martingale is really based on a simple symmetric random walk.
 Test understanding: compute the expected value of M_n from first principles.
- Test understanding: what should be the value of $\mathbb{E}[\tilde{M}_n = 1]$ if \tilde{M} is computed as for M but stopping play if M hits level $1 - 2^N$? (Think about this, but note that a satisfactory answer has to await discussion of optional stopping theorem in next section.)

Martingales and populations

- Consider a **branching process** Y : population at time n is Y_n , where $Y_0 = 1$ (say) and Y_{n+1} is the sum $Z_{n,1} + \dots + Z_{n,Y_n}$ of Y_n independent copies of a non-negative integer-valued **family-size r.v.** Z .
- Suppose $\mathbb{E}[Z] = \mu < \infty$. Then $X_n = Y_n / \mu^n$ defines a martingale.
- Suppose $\mathbb{E}[s^Z] = G(s)$. Let $H_n = Y_0 + \dots + Y_n$ be total of all populations up to time n . Then $s^{H_n} / (G(s)^{H_{n-1}})$ defines a martingale.
- In both these examples we can use $\mathbb{E}[\dots | \mathcal{F}_n]$, representing conditioning by all $Z_{m,i}$ for $m \leq n$.

Martingales and populations

- Consider a **branching process** Y : population at time n is Y_n , where $Y_0 = 1$ (say) and Y_{n+1} is the sum $Z_{n,1} + \dots + Z_{n,Y_n}$ of Y_n independent copies of a non-negative integer-valued family size r.v. Z .
- Suppose $\mathbb{E}[Z] = \mu < \infty$. Then $X_n = Y_n / \mu^n$ defines a martingale.
- Suppose $\mathbb{E}[s^Z] = G(s)$. Let $H_n = Y_0 + \dots + Y_n$ be total of all populations up to time n . Then $s^{H_n} / (G(s)^{H_{n-1}})$ defines a martingale.
- In both these examples we can use $\mathbb{E}[\dots | \mathcal{F}_n]$, representing conditioning by all $Z_{m,i}$ for $m \leq n$.

- New Yorker's definition of branching process (to be read out aloud in strong New York accent): "You are born. You live a while. You have a random number of kids. You die. Your children are completely independent of you, but behave in exactly the same way." The formal definition requires the $Z_{n,i}$ to be independent of Y_0, \dots, Y_n .
- Test understanding: check this example.
- Test understanding: check this example.
- Indeed, we can also generalize to general Y_0 .

Definition of a martingale

Formally:

Definition

X is a **martingale** if $\mathbb{E}[|X_n|] < \infty$ (for all n) and

$$X_n = \mathbb{E}[X_{n+1} | \mathcal{F}_n].$$

1. It is important that the X_n are integrable.
2. It is a consequence that X_n is part of the conditioning expressed by \mathcal{F}_n .
3. Sometimes we expand the reference to \mathcal{F}_n :

$$X_n = \mathbb{E}[X_{n+1} | X_n, X_{n-1}, \dots, X_1, X_0].$$

Supermartingales and submartingales

Two associated definitions

Definition

$\{X_n\}$ is a **supermartingale** if $\mathbb{E}[|X_n|] < \infty$ (for all n) and

$$X_n \geq \mathbb{E}[X_{n+1} | \mathcal{F}_n],$$

(and X_n forms part of conditioning expressed by \mathcal{F}_n).

Definition

$\{X_n\}$ is a **submartingale** if $\mathbb{E}[|X_n|] < \infty$ (for all n) and

$$X_n \leq \mathbb{E}[X_{n+1} | \mathcal{F}_n],$$

(and X_n forms part of conditioning expressed by \mathcal{F}_n).

1. It is important that the X_n are integrable. It is now *not* automatic that X_n forms part of the conditioning expressed by \mathcal{F}_n , and it is important that this is part of the definition.
2. It is important that the X_n are integrable. Again it is important that X_n forms part of the conditioning expressed by \mathcal{F}_n . How to remember the difference between "sub-" and "super-?" Suppose $\{X_n\}$ measures your fortune in a casino gambling game. Then "sub-" is bad and "super-" is good for the casino! Wikipedia: life is a supermartingale, as one's expectations are always no greater than one's present state.

Examples of supermartingales and submartingales

1. Consider asymmetric simple random walk: supermartingale if jumps have negative expectation, submartingale if jumps have positive expectation.
2. This holds even if the walk is stopped on first return to 0.
3. Consider Thackeray's martingale based on asymmetric random walk. This is a supermartingale or a submartingale depending on whether jumps have negative or positive expectation.
4. Consider branching process $\{Y_n\}$ and consider Y_n on its own instead of Y_n/μ^n . This is a supermartingale if $\mu < 1$ (sub-critical case), a submartingale if $\mu > 1$ (super-critical case), a martingale if $\mu = 1$ (critical case).

Test understanding: check all these examples. In each case the general procedure is as follows: compare $\mathbb{E}[X_{n+1} | \mathcal{F}_n]$ to X_n .

An example of importance in finance

1. Suppose N_1, N_2, \dots are independent identically distributed normal random variables of mean 0 and variance σ^2 , and put $S_n = N_1 + \dots + N_n$.
2. Then the following is a martingale:

$$Y_n = \exp\left(S_n - \frac{n}{2}\sigma^2\right).$$

▶ ANIMATION

3. A modification exists for which the N_i have non-zero mean μ .
Hint: $S_n \rightarrow S_n - n\mu$.

Here (modifications of) Y_n provides the simplest model for market price fluctuations appropriately discounted.

1. In fact $\{S_n\}$ is a martingale, though this is not the point here.
2. **Test understanding:** Prove this!
Hint: $\mathbb{E}[\exp(N_1)] = e^{\sigma^2/2}$.
3. **Test understanding:** figure out the modification!
4. A continuous-time variation on this (using Brownian motion) is an important baseline model in mathematical finance. Note that the martingale can be expressed as

$$Y_{n+1} = Y_n \exp\left(N_{n+1} - \frac{\sigma^2}{2}\right).$$

Martingales and likelihood

- ▶ Suppose X_1, X_2, \dots are observed at times 1, 2, ... Write down likelihood at time n :

$$L(\theta; X_1, \dots, X_n) = p(X_1, \dots, X_n | \theta).$$

- ▶ If θ_0 is "true" value then (computing expectation with $\theta = \theta_0$)

$$\mathbb{E} \left[\frac{L(\theta_1; X_1, \dots, X_{n+1})}{L(\theta_0; X_1, \dots, X_{n+1})} \middle| \mathcal{F}_n \right] = \frac{L(\theta_1; X_1, \dots, X_n)}{L(\theta_0; X_1, \dots, X_n)}$$

Martingales and likelihood

- ▶ Suppose X_1, X_2, \dots are observed at times 1, 2, ... Write down likelihood at time n :
 $L(\theta; X_1, \dots, X_n) = p(X_1, \dots, X_n | \theta)$.
- ▶ If θ_0 is "true" value then computing expectation with $\theta = \theta_0$:
 $\mathbb{E} \left[\frac{L(\theta_1; X_1, \dots, X_{n+1})}{L(\theta_0; X_1, \dots, X_{n+1})} \middle| \mathcal{F}_n \right] = \frac{L(\theta_1; X_1, \dots, X_n)}{L(\theta_0; X_1, \dots, X_n)}$

1. Simple case of normal data with unknown mean θ :

$$L(\theta; X_1, \dots, X_n) \propto \exp \left(\sum_1^n X_i - \frac{n}{2} \theta^2 \right).$$

2. Hence likelihood ratios are really the same thing as martingales.
3. The martingale in the finance example can also arise in this way, as the likelihood ratio between two different values of θ if the model is that the X_i are independent identically distributed $N(\theta, \sigma^2)$.

The "Chicken Little" example

1. An comet may or may not collide with Earth in n days time. Chaotic dynamics: model by supposing comet may follow one of n possible paths, of which just one leads to collision at day n .
2. Each day, new observations eliminate exactly one of possible paths: path to be eliminated on day r is chosen from $n - r + 1$ surviving paths uniformly at random and independently of the past.
3. Compute conditional collision probability at day r , supposing collision path is not yet been eliminated. Deduce that conditional collision probabilities at days $r = 0, 1, \dots, n$ form a martingale.

The "Chicken Little" example

1. An comet may or may not collide with Earth in n days time. Chaotic dynamics: model by supposing comet may follow one of n possible paths, of which just one leads to collision at day n .
2. Each day, new observations eliminate exactly one of possible paths: path to be eliminated on day r is chosen from $n - r + 1$ surviving paths uniformly at random and independently of the past.
3. Compute conditional collision probability at day r , supposing collision path is not yet been eliminated. Deduce that conditional collision probabilities at days $r = 0, 1, \dots, n$ form a martingale.

1. This is a considerable simplification of chaotic dynamics, but not unreasonable.
2. This models the fact that observations are hard to come by, and do not provide much information.
3. Let D be indicator random variable indicating event that collision occurs, and compute $\mathbb{E}[D | \mathcal{F}_r]$ where \mathcal{F}_r captures information of whether or not collision occurs by day r . Probability of collision grows more and more rapidly ($1/(n-r)$ on day r) till either it suddenly falls to zero (if collision path eliminated before n) or collision actually occurs (if collision path not eliminated before day n). Therefore collision probability increases day by day (engendering increasing despair), hopefully until it falls to zero (engendering mass relief).

Stopping times

"Hurry please it's time."
 T. S. Eliot,
 The Waste Land, 1922

- "No-look-ahead" condition
- Random walk example
- Branching process example
- Events revealed by stopping time
- Optional Stopping Theorem
- Application to gambling
- Hitting times
- Martingale convergence
- Harmonic functions



Stopping times

"Hurry please it's time."
 T. S. Eliot,
 The Waste Land, 1922

"No-look-ahead" condition
 Random walk example
 Branching process example
 Events revealed by stopping time
 Optional Stopping Theorem
 Application to gambling
 Hitting times
 Martingale convergence
 Harmonic functions

Playing a fair game, what happens if you adopt a strategy of leaving the game at a random time? For "reasonable" random times, this should offer you no advantage. Here we seek to make sense of the term "reasonable". Note that the gambling motivation is less frivolous than it might appear. Mathematical finance is about developing trading strategies (complex gambles!) aimed at controlling uncertainty.

Stopping times

The big idea

Martingales M stopped at "nice" times are still martingales. In particular, for "nice" random T ,

$$\mathbb{E}[M_T] = \mathbb{E}[M_0]$$

For a random time T to be "nice", two things are required:

1. T must not "look ahead";
2. T must not be "too big".
3. Note that in general a useful random time T can have positive chance of being infinite.

▶ ANIMATION

Stopping times
 The big idea

Martingales M stopped at "nice" times are still martingales. In particular, for "nice" random T ,

$$\mathbb{E}[M_T] = \mathbb{E}[M_0]$$

For a random time T to be "nice", two things are required:

1. T must not "look ahead"
2. T must not be "too big"
3. Note that in general a useful random time T can have positive chance of being infinite.

How can T fail to be "nice"? Consider simple symmetric random walk X begun at 0.

1. Example of "looking ahead": Set $S = \sup\{X_n : 0 \leq n \leq 10\}$ and set $T_2 = \inf\{n : X_n = S\}$. Then $\mathbb{E}[X_{T_2}] \geq \mathbb{P}[S > 0] > 0 = \mathbb{E}[X_0]$
2. Example of being "too big": $T_1 = \inf\{n : X_n = 1\}$ so (assuming T is almost surely finite) $\mathbb{E}[X_{T_1}] = 1 > 0 = \mathbb{E}[X_0]$
3. Example of possibly being infinite: asymmetric simple random walk X begun at 0, $\mathbb{E}[X_1] < 0$, $T_1 = \inf\{n : X_n = 1\}$ as above.

Non-obvious "no-look-ahead" condition

Definition

A non-negative integer-valued random variable T is said to be a **stopping time** if (equivalently) for all n

- ▶ $[T \leq n]$ is determined by information at time n ;
- ▶ or $[T \leq n] \in \mathcal{F}_n$
- ▶ or we can write **rules** (Bernoulli random variables) ζ_0, ζ_1, \dots with ζ_n in \mathcal{F}_n , such that

$$[\zeta_n = 1] = [T \leq n].$$

Non-obvious "no-look-ahead" condition
 Definition
 A non-negative integer-valued random variable T is said to be a **stopping time** if (equivalently) for all n
 • $[T \leq n]$ is determined by information at time n
 • $\equiv [T \leq n] \in \mathcal{F}_n$
 • or we can write rules (Bernoulli random variables) ζ_0, ζ_1, \dots with ζ_n in \mathcal{F}_n , such that
 $[\zeta_n = 1] = [T \leq n]$.

Note that we need to have a clear notion of exactly what might be \mathcal{F}_n , information revealed by time n .

Here is a poetical illustration of a **non-stopping time**, due to my father David Kendall:

*There is a rule for timing toast,
 You never need to guess;
 Just wait until it starts to smoke,
 And then ten seconds less.*

(Adapted from a "grook" by Piet Hein, *Grooks II* MIT Press, 1968.)

Recall the example on previous slide of T being the time to hit 1 for a negatively-biased simple random walk begun at 0: stopping times can have positive chance of being infinite.

Example using random walks

Let X be a random walk begun at 0.

- ▶ The random time $T = \inf\{n > 0 : X_n \geq 10\}$ is a stopping time.
- ▶ Indeed $[T \leq n]$ is clearly determined by information at time n :

$$[T \leq n] = [X_1 \geq 10] \cup \dots \cup [X_n \geq 10].$$

- ▶ Finally, T is typically "too big": so long as it is almost surely finite, we find that $0 = \mathbb{E}[X_0] < \mathbb{E}[X_T]$.
 Finiteness is the case if $\mathbb{E}[X_1] > 0$ or if $\mathbb{E}[X_1] = 0$ and $\mathbb{P}[X_1 > 0] > 0$.

Example using random walks
 Let X be a random walk begun at 0.
 • The random time $T = \inf\{n > 0 : X_n \geq 10\}$ is a stopping time.
 • Indeed $[T \leq n]$ is clearly determined by information at time n :
 $[T \leq n] = [X_1 \geq 10] \cup \dots \cup [X_n \geq 10]$.
 • Finally, T is typically "too big": so long as it is almost surely finite, we find that $0 = \mathbb{E}[X_0] < \mathbb{E}[X_T]$.
 Finiteness is the case if $\mathbb{E}[X_1] > 0$ or if $\mathbb{E}[X_1] = 0$ and $\mathbb{P}[X_1 > 0] > 0$.

1. X need not be symmetric, need not be simple. Indeed a Markov chain or even a general random process would do.
2. We could replace $n > 0$ by $n \geq 0, X \geq 10$ by $X \in A$ for some subset A of state-space, ...:
 thus we could have $T_A = \inf\{n > 0 : X_n \in A\}$ (the "hitting time on A ").
3. In case of hitting time on A ,

$$[T_A \leq n] = [X_1 \in A] \cup \dots \cup [X_n \in A]$$

so $[T_A \leq n]$ is determined by information at time n , so T_A is a stopping time.

4. General hitting times T_A need not be "too big": example if X is simple symmetric random walk begun at 0 and $A = \{\pm 10\}$.

Example using branching processes

Let Y be a branching process of mean-family-size μ (so Y_n/μ^n determines a martingale), with $Y_0 = 1$.

- ▶ The random time $T = \inf\{n : Y_n = 0\} = \inf\{n : X_n = 0\}$ is a stopping time.
- ▶ Indeed $[T \leq n]$ is clearly determined by information at time n :

$$[T \leq n] = [Y_n = 0]$$

since $Y_{n-1} = 0$ implies $Y_n = 0$ et cetera.

- ▶ Again T here is "too big": so long as it is almost surely finite then $1 = \mathbb{E}[X_0] > \mathbb{E}[X_T]$.
 Finiteness occurs if $\mu < 1$, or if $\mu = 1$ and there is positive chance of zero family size.

Example using branching processes
 Let Y be a branching process of mean-family-size μ (so Y_n/μ^n determines a martingale), with $Y_0 = 1$.
 • The random time $T = \inf\{n : Y_n = 0\} = \inf\{n : X_n = 0\}$ is a stopping time.
 • Indeed $[T \leq n]$ is clearly determined by information at time n :
 $[T \leq n] = [Y_n = 0]$
 since $Y_{n-1} = 0$ implies $Y_n = 0$ et cetera.
 • Again T here is "too big": so long as it is almost surely finite then $1 = \mathbb{E}[X_0] > \mathbb{E}[X_T]$.
 Finiteness occurs if $\mu < 1$, or if $\mu = 1$ and there is positive chance of zero family size.

1. So $Y_n = Z_{n-1,1} + \dots + Z_{n-1,Y_{n-1}}$ for independent family sizes $Z_{m,j}$.
2. For a more interesting example, consider

$$S = \inf\{n : \text{at least one family of size 0 before } n\}$$

3. In case of S , consider

$$[S \leq n] = A_0 \cup A_1 \cup \dots \cup A_{n-1}$$

where $A_i = [Z_{i,j} = 0 \text{ for some } j \leq Y_i]$. Thus $[S \leq n]$ is determined by information at time n , so S is a stopping time.

4. It is important to be clear about what *is* information provided at time n . Here we suppose it to be made up only of the sizes of families produced by individuals in generations $0, 1, \dots, n-1$. Other choices are possible, of course.

Events revealed by the time of a stopping time T

Suppose T is a stopping time.

Definition

The "pre- T σ -algebra" \mathcal{F}_T is composed of events which, if T does not occur later than time n , are themselves determined at time n . Thus:

$$A \in \mathcal{F}_T \text{ if } A \cap [T \leq n] \in \mathcal{F}_n \text{ for all } n.$$

Definition

Random variables Z are said to be " \mathcal{F}_T -measurable" if events made up from them ($[Z \leq z], \dots$) are in pre- T σ -algebra \mathcal{F}_T .

Events revealed by the time of a stopping time T
 Suppose T is a stopping time.
 Definition
 The "pre- T σ -algebra" \mathcal{F}_T is composed of events which, if T does not occur later than time n , are themselves determined at time n . Thus:
 $A \in \mathcal{F}_T$ if $A \cap [T \leq n] \in \mathcal{F}_n$ for all n .
 Definition
 Random variables Z are said to be " \mathcal{F}_T -measurable" if events made up from them ($[Z \leq z], \dots$) are in pre- T σ -algebra \mathcal{F}_T .

1. Consider random walk X begun at 0 and the stopping time $T = \inf\{n : X_n \geq 10\}$. Then the event $[X_{15} < 5 \text{ and } T > 15]$ is in the pre- T σ -algebra \mathcal{F}_T .
2. The random variable $X_{\min\{15, T\}}$ is \mathcal{F}_T -measurable.
3. Consider the branching process example with S being the time at which a zero-size family is first encountered. Then

$$Y_0 + Y_1 + \dots + Y_S \in \mathcal{F}_S.$$

Optional stopping theorem

Theorem

Suppose M is a martingale and $S \leq T$ are two **bounded** stopping times. Then

$$\mathbb{E}[M_T | \mathcal{F}_S] = M_S.$$

We can generalize to general stopping times $S \leq T$ if either M is bounded or M is "uniformly integrable".

Uniform integrability: note we can take expectation of a single random variable X exactly when $\mathbb{E}[|X|; |X| > n] \rightarrow 0$ as $n \rightarrow \infty$. (This fails when $\mathbb{E}[|X|; |X| > n] = \infty$!). Uniform integrability requires this to hold **uniformly** for a whole collection of random variables X_i :

$$\lim_{n \rightarrow \infty} \sup_i \mathbb{E}[|X_i|; |X_i| > n] = 0.$$

Examples: if the X_i are bounded, if there is a single non-negative random variable Z with $\mathbb{E}[Z] < \infty$ and $|X_i| \leq Z$ for all i , if the p -moments $\mathbb{E}[X_i^p]$ are bounded for some $p > 1$.

Theorem
 Suppose M is a martingale and S, T are two bounded stopping times. Then
 $\mathbb{E}[M_T | \mathcal{F}_S] = M_S$
 We can generalize to general stopping times $S \leq T$ if either M is bounded or M is uniformly integrable.

Gambling: you shouldn't expect to win

Suppose your fortune in a gambling game is X , a martingale begun at 0 (for example, a simple symmetric random walk. If N is the maximum time you can spend playing the game, and if $T \leq N$ is a bounded stopping time, then

$$\mathbb{E}[X_T] = 0.$$

Contrast Fleming (1953):

"Then the Englishman, Mister Bond, increased his winnings to exactly three million over the two days. He was playing a progressive system on red at table five. . . . It seems that he is persevering and plays in maximums. He has luck."

There are exceptions, for example Blackjack (using card-counting: en.wikipedia.org/wiki/Card_counting). I find proposed strategies in other games less convincing, for example the Labouchère system favoured by Ian Fleming (en.wikipedia.org/wiki/Labouch%C3%A8re_system):

The Labouchère system, also called the cancellation system, is a gambling strategy used in roulette. The user of such a strategy decides before playing how much money they want to win, and writes down a list of positive numbers that sum to the predetermined amount. With each bet, the player stakes an amount equal to the sum of the first and last numbers on the list. If only one number remains, that number is the amount of the stake. If bet is successful, the two amounts are removed from the list. If the bet is unsuccessful, the amount lost is appended to the end of the list. This process continues until either the list is completely crossed out, at which point the desired amount of money has been won, or until the player runs out of money to wager.

Gambling: you shouldn't expect to win
 Suppose your fortune in a gambling game is X , a martingale begun at 0 (for example, a simple symmetric random walk. If N is the maximum time you can spend playing the game, and if $T \leq N$ is a bounded stopping time, then
 $\mathbb{E}[X_T] = 0$.
Contrast Fleming (1953)
 Then the Englishman, Mister Bond, increased his winnings to exactly three million over the two days. He was playing a progressive system on red at table five. . . . It seems that he is persevering and plays in maximums. He has luck."

Martingales and hitting times

Suppose X_1, X_2, \dots are independent Gaussian random variables of mean $-\mu < 0$ and variance 1. Let $Y_n = X_1 + \dots + X_n$ and let T be the time when Y first exceeds level $\ell > 0$. Then $\exp\left(\alpha(Y_n + \mu n) - \frac{\alpha^2}{2}n\right)$ determines a martingale, and the optional stopping theorem can be applied to show

$$\mathbb{E}[\exp(-pT)] \sim e^{-(\mu + \sqrt{\mu^2 + 2p})\ell}.$$

This improves to an equality, at the expense of using more advanced theory, if we replace the Gaussian random walk Y by Brownian motion.

So $T = \inf\{n : Y_n \geq \ell\}$. Use the optional stopping theorem on the bounded stopping time $\min\{T, n\}$:

$$\mathbb{E}\left[\exp\left(\alpha Y_{\min\{T, n\}} + \alpha\left(\mu - \frac{\alpha}{2}\right)\min\{T, n\}\right)\right] = 1.$$

Use careful analysis of the left-hand side, letting $n \rightarrow \infty$, large ℓ ,

$$\mathbb{E}\left[\exp\left(\alpha\ell + \alpha\left(\mu - \frac{\alpha}{2}\right)T\right)\right] \sim 1.$$

Now set $\alpha = \mu + \sqrt{\mu^2 + 2p} > 0$, so $\alpha(\mu - \frac{\alpha}{2}) = -p$:

$$\mathbb{E}[\exp(-pT)] \sim \exp\left(-(\mu + \sqrt{\mu^2 + 2p})\ell\right).$$

Improvement: Brownian motion is continuous in time and so cannot jump over the level ℓ without hitting it.

Martingales and hitting times
 Suppose X_1, X_2, \dots are independent Gaussian random variables of mean $-\mu < 0$ and variance 1. Let $Y_n = X_1 + \dots + X_n$ and let T be the time when Y first exceeds level $\ell > 0$. Then $\exp\left(\alpha(Y_n + \mu n) - \frac{\alpha^2}{2}n\right)$ determines a martingale, and the optional stopping theorem can be applied to show
 $\mathbb{E}[\exp(-pT)] \sim e^{-(\mu + \sqrt{\mu^2 + 2p})\ell}$.
 This improves to an equality, at the expense of using more advanced theory, if we replace the Gaussian random walk Y by Brownian motion.

Martingale convergence

Theorem

Suppose X is a non-negative supermartingale. Then $Z = \lim X_n$ exists, moreover $\mathbb{E}[Z | \mathcal{F}_n] \leq X_n$.

Theorem

Suppose X is a bounded martingale (or, more generally, uniformly integrable). Then $Z = \lim X_n$ exists, moreover $\mathbb{E}[Z | \mathcal{F}_n] = X_n$.

Theorem

Suppose X is a martingale and $\mathbb{E}[X_n^2] \leq K$ for some fixed constant K . Then one can prove directly that $Z = \lim X_n$ exists, moreover $\mathbb{E}[Z | \mathcal{F}_n] = X_n$.

1. Consider symmetric simple random walk begun at 1 and stopped at 0: $X_n = Y_{\min\{n, T\}}$ if $T = \inf\{n : Y_n = 0\}$ and Y is symmetric simple random walk. Clearly $X_n = Y_{\min\{n, T\}}$ is non-negative; clearly $X_n = Y_{\min\{n, T\}} \rightarrow Z = 0$, since Y will eventually hit 0; clearly $0 = \mathbb{E}[Z | \mathcal{F}_n] \leq X_n$ since $Y_n \geq 0$.
2. Thus symmetric simple random walk Y begin at 0 and stopped at ± 10 must converge to a limiting value Z . Evidently $Z = \pm 10$. Moreover since $\mathbb{E}[Z | \mathcal{F}_n] = Y_n$ we deduce $\mathbb{P}[Z = 10 | \mathcal{F}_n] = \frac{Y_n + 10}{20}$.
3. Sketch argument: from martingale property

$$\mathbb{E}\left[(X_{m+n} - X_n)^2 | \mathcal{F}_n\right] = \mathbb{E}\left[X_{m+n}^2 | \mathcal{F}_n\right] - X_n^2;$$

hence $\mathbb{E}[X_n^2]$ is non-decreasing; hence it converges to a limiting value;

hence $\mathbb{E}[(X_{m+n} - X_n)^2]$ tends to 0.

Martingale convergence
Theorem
 Suppose X is a non-negative supermartingale. Then $Z = \lim X_n$ exists, moreover $\mathbb{E}[Z | \mathcal{F}_n] \leq X_n$.
Theorem
 Suppose X is a bounded martingale (or, more generally, uniformly integrable). Then $Z = \lim X_n$ exists, moreover $\mathbb{E}[Z | \mathcal{F}_n] = X_n$.
Theorem
 Suppose X is a martingale and $\mathbb{E}[X_n^2] \leq K$ for some fixed constant K . Then one can prove directly that $Z = \lim X_n$ exists, moreover $\mathbb{E}[Z | \mathcal{F}_n] = X_n$.

- ↳ 3: Stopping times
- ↳ Harmonic functions

Martingales and bounded harmonic functions

- ▶ Consider a discrete state-space Markov chain X with transition kernel p_{ij} . Suppose $f(i)$ is a bounded **harmonic function**: a function for which $f(i) = \sum_j f(j)p_{ij}$. Then $f(X)$ is a bounded martingale, hence must converge as time increases to infinity.
- ▶ The simplest example: consider simple random walk X absorbed at boundaries $a < b$. Then $f(x) = \frac{x-a}{b-a}$ is a bounded harmonic function, and can be shown to satisfy

$$f(x) = \mathbb{P}[X \text{ hits } b \text{ before } a | X_0 = x].$$

- ▶ Another example: given branching process Y and family size generating function $G(s)$, suppose ζ is smallest non-negative root of $\zeta = G(\zeta)$. Set $f(y) = \zeta^y$. Check this is a non-negative martingale (and therefore harmonic).

- ↳ 3: Stopping times
- ↳ Harmonic functions
- ↳ Martingales and bounded harmonic functions

Martingales and bounded harmonic functions

- Consider a discrete state-space Markov chain X with transition kernel p_{ij} . Suppose f is a bounded harmonic function for which $f(i) = \sum_j p_{ij} f(j)$. Then $f(X)$ is a bounded martingale, hence must converge as time increases to infinity.
- The simplest example: consider simple random walk X absorbed at boundaries $a < b$. Then $f(x) = \frac{x-a}{b-a}$ is a bounded harmonic function, and can be shown to satisfy $f(x) = \mathbb{P}[X \text{ hits } b \text{ before } a | X_0 = x]$.
- Another example: given branching process Y and family size generating function $G(s)$, suppose ζ is smallest non-negative root of $\zeta = G(\zeta)$. Set $f(y) = \zeta^y$. Check this is a non-negative martingale (and therefore harmonic).

1. The terminology supermartingale/submartingale was actually chosen to mirror the potential-theoretic terminology superharmonic/subharmonic.
2. Use martingale convergence theorem and optional stopping theorem.
3. We'd like to say, therefore $f(y) = \mathbb{P}[Y \text{ becomes extinct} | Y_0 = y]$. Since $\zeta \leq 1$, it follows f is bounded, so this follows as before.
4. Further significant examples come from, for example, multidimensional random walk absorbed at boundary of a geometric region.

- ↳ 4: Counting and compensating

Counting and compensating

"It is a law of nature we overlook, that intellectual versatility is the compensation for change, danger, and trouble."
 H. G. Wells,
The Time Machine, 1896

- Simplest example: Poisson process
- Compensators
- Examples
- Variance of compensated counting process
- Counting processes and Poisson processes
- Compensation of population processes



- ↳ 4: Counting and compensating
- ↳ Counting and compensating

Counting and compensating

There are two directions to go with the Poisson process:

- view as a pattern of points
- Slivnyak's theorem: condition on t being a transition / incident
- PASTA principle: if a Markov chain has "arrivals" following a Poisson distribution, then in statistical equilibrium Poisson Arrivals See Time Averages.
- How to make points "interact"?
- Generalize to Poisson patterns of geometric objects.

view as counting process and generalize:

- varying "hazard rate";
- relate to martingales?

Here we follow the second direction.

We can now make a connection between martingales and Markov chains. We start with the Poisson process, viewed as a process used for counting incidents, and show how martingales can be used to describe much more general counting processes.

- ↳ 4: Counting and compensating
- ↳ Simplest example: Poisson process

Simplest example: Poisson process

Consider birth-death-immigration process from above, with birth and death rates set to zero: $\lambda = \mu = 0$. The result is a **Poisson process** of rate α as described before:

Definition

A continuous-time Markov chain N is a Poisson process of rate $\alpha > 0$ if the only transitions are $N \rightarrow N + 1$ of rate α .

Theorem

If N is Poisson process of rate α then

$$\mathbb{P}[N_t = k] = \mathbb{P}[Poisson(\alpha t) = k] = \frac{(\alpha t)^k}{k!} e^{-\alpha t}.$$

The times of transitions are often referred to as **incidents**.

- ↳ 4: Counting and compensating
- ↳ Simplest example: Poisson process
- ↳ Simplest example: Poisson process

Simplest example: Poisson process

Consider birth-death-immigration process from above, with birth and death rates set to zero: $\lambda = \mu = 0$. The result is a Poisson process of rate α as described before:

Definition

A continuous-time Markov chain N is a Poisson process of rate $\alpha > 0$ if the only transitions are $N \rightarrow N + 1$ of rate α .

Theorem

If N is Poisson process of rate α then

$$\mathbb{P}[N_t = k] = \mathbb{P}[Poisson(\alpha t) = k] = \frac{(\alpha t)^k}{k!} e^{-\alpha t}.$$

The times of transitions are often referred to as incidents.

1. This has a claim to be the simplest possible continuous-time Markov chain. Its state-space is *very* reducible, so it does not supply good examples for questions of equilibrium!
2. In one approach to stochastic processes this serves as a fundamental building block for more complicated processes.
3. Times between consecutive incidents are independent Exponential(α). Thence a whole wealth of distributional relationships between Exponential, Poisson, and indeed Gamma, Geometric, Hypergeometric, . . .
4. A more general result is suggestive about how to generalize to Poisson point patterns: if $A \subset [0, \infty)$ has length measure a then

$$\mathbb{P}[k \text{ incidents in } A] = \mathbb{P}[Poisson(\alpha a) = k].$$

5. A significant converse: given a random point pattern such that
- $$\mathbb{P}[\text{No incidents in } A] = \exp(-\alpha a)$$

for any A of length measure a , the point pattern marks the incidents of a Poisson counting process of rate α .

- ↳ 4: Counting and compensating
- ↳ Simplest example: Poisson process

Poisson process directions

There are two directions to go with the Poisson process:

- ▶ view as a pattern of points:
 - ▶ Slivnyak's theorem: condition on t being a transition / incident. Then remaining incidents form transitions of Poisson process of same rate.
 - ▶ PASTA principle: if a Markov chain has "arrivals" following a Poisson distribution, then in statistical equilibrium Poisson Arrivals See Time Averages.
 - ▶ How to make points "interact"?
 - ▶ Generalize to Poisson patterns of geometric objects.
- ▶ view as counting process and generalize:
 - ▶ varying "hazard rate";
 - ▶ relate to martingales?

Here we follow the second direction.

- ↳ 4: Counting and compensating
- ↳ Simplest example: Poisson process
- ↳ Poisson process directions

Poisson process directions

There are two directions to go with the Poisson process:

- view as a pattern of points
- Slivnyak's theorem: condition on t being a transition / incident
- PASTA principle: if a Markov chain has "arrivals" following a Poisson distribution, then in statistical equilibrium Poisson Arrivals See Time Averages.
- How to make points "interact"?
- Generalize to Poisson patterns of geometric objects.

view as counting process and generalize:

- varying "hazard rate";
- relate to martingales?

Here we follow the second direction.

1. Slivnyak's theorem generalizes directly to Poisson point patterns. The trick is, of course, to make sense of conditioning on an event of probability 0.
2. That is to say, at "just before" the arrival time, the probability that the system is in state k is π_k the equilibrium probability. Easy consequence of Slivnyak's theorem.
3. Crucial for calculations: the chance of seeing no object of given kind in given region is $\exp(-\mu)$ where μ is mean number of such objects.
4. What is the hazard rate? does it suggest generalizations?

Hazard rate and compensators

Starting point: if N is Poisson process of rate α then

- ▶ (“mean”) $N_t - \alpha t$ determines a martingale;
- ▶ (“variance”) $(N_t - \alpha t)^2 - \alpha t$ determines a martingale;

Consider processes which “count” incidents:

Definition

A counting process is a continuous-time process—not necessarily Markov—changing by single jumps of +1. Try to subtract something to turn it into a martingale.

Definition

We say $\int_0^t \ell(s) ds$ compensates a counting process N if

- ▶ the (possibly random) $\ell(s)$ is in \mathcal{F}_s ;
- ▶ $N_t - \int_0^t \ell(s) ds$ determines a martingale.

1. Calculation based on $\mathbb{E}[N_{t+s} - N_s | \mathcal{F}_s] = \alpha t$.
2. Calculation based on $\text{Var}[N_{t+s} - N_s | \mathcal{F}_s] = \alpha t$.
3. Later we will also consider population processes counting births +1 and deaths -1.
4. It is possible to make a more general definition which replaces $\int_0^t \ell(s) ds$ by a non-decreasing process Λ_t —but then we have to require “ $\Lambda_t \in \mathcal{F}_{t-}$ ”.
5. It can then be shown that
 - compensators always exist
 - and are essentially unique.

Hazard rate and compensators
 Starting point: if N is Poisson process of rate α then
 - “mean”: $N_t - \alpha t$ determines a martingale;
 - “variance”: $(N_t - \alpha t)^2 - \alpha t$ determines a martingale.
 Consider processes which “count” incidents.
Definition
 A counting process is a continuous-time process—not necessarily Markov—changing by single jumps of +1. Try to subtract something to turn it into a martingale.
Definition
 We say $\int_0^t \ell(s) ds$ compensates a counting process N if
 - the (possibly random) $\ell(s)$ is in \mathcal{F}_s ;
 - $N_t - \int_0^t \ell(s) ds$ determines a martingale.

Example: random sample of lifetimes

Suppose X_1, \dots, X_n are independent and identically distributed non-negative random variables (lifetimes) with common density f .

- ▶ Set $\mathbb{P}[X_i > t] = 1 - \int_0^t f(s) ds = \exp(-\int_0^t h(s) ds)$.
- ▶ Counting process $N_t = \#\{i : X_i \leq t\}$ increases by +1 jumps in continuous time.
- ▶ Observe:
 - ▶ $N_t - \int_0^t h(s) N_s ds$ is a martingale.
 - ▶ $(N_t - \int_0^t h(s) N_s ds)^2 - \int_0^t h(s) N_s ds$ is a martingale.

1. Resolves to showing the following is a martingale:

$$\mathbb{P}[X_i \leq t] - \int_0^{\min\{t, X_i\}} h(u) du.$$

Key calculation: the expectation of the above is

$$\mathbb{P}[X_i \leq t] - \int_0^t h(u) \mathbb{P}[X_i > u] du,$$

which vanishes if we substitute in $\mathbb{P}[X_i > u] = \exp(-\int_0^u h(s) ds)$. This of course is computation of an absolute probability: **Test understanding:** make changes to get the relevant conditional probability calculation.

2. This follows most directly by noting independence of the $\mathbb{P}[X_i \leq t] - \int_0^{\min\{t, X_i\}} h(s) ds$. **However** it is actually true for a more general reason

Example: random sample of lifetimes
 Suppose X_1, \dots, X_n are independent and identically distributed non-negative random variables (lifetimes) with common density f .
 - Set $\mathbb{P}[X_i > t] = 1 - \int_0^t f(s) ds = \exp(-\int_0^t h(s) ds)$.
 - Counting process $N_t = \#\{i : X_i \leq t\}$ increases by +1 jumps in continuous time.
 - Observe:
 - $N_t - \int_0^t h(s) N_s ds$ is a martingale.
 - $(N_t - \int_0^t h(s) N_s ds)^2 - \int_0^t h(s) N_s ds$ is a martingale.

Example: pure birth process

Example (Pure birth process)

If the pure birth process N makes transitions $N \rightarrow N + 1$ at rate λN then

$$N_t - \int_0^t \lambda N_s ds \quad \text{is a martingale.}$$

Here again one can check that the expression of variance type $(N_t - \int_0^t \lambda N_s ds)^2 - \int_0^t \lambda N_s ds$ also determines a martingale.

1. A direct proof can be obtained by computing the distribution of N_t given N_0 . Alternatively here is a plausibility argument: in a small period of time $[t, t + \Delta t)$ it is most likely no transition will occur; the chance of one transition is about $\lambda N_t \Delta t$, and the chance of more is infinitesimal. So the conditional mean increment is $\lambda N_t \Delta t$ which is exactly matched by the compensator.
2. Direct computations would permit a direct proof; but a similar plausibility argument also applies. The conditional variance of the increment is about $\lambda N_t \Delta t (1 - \lambda N_t \Delta t) \approx \lambda N_t \Delta t$, again matching the compensator.

Example: pure birth process
 Example (Pure birth process)
 If the pure birth process N makes transitions $N \rightarrow N + 1$ at rate λN then
 $N_t - \int_0^t \lambda N_s ds$ is a martingale.
 Here again one can check that the expression of variance type $(N_t - \int_0^t \lambda N_s ds)^2 - \int_0^t \lambda N_s ds$ also determines a martingale.

Variance of compensated counting process

The above expression of variance type holds more generally:

Theorem

Suppose N is a counting process compensated by $\int \ell(s) ds$. Then

$$\left(N_t - \int_0^t \ell(s) ds\right)^2 - \int_0^t \ell(s) ds \quad \text{is a martingale.}$$

Rigorous proof, or heuristic limiting argument

1. The key point of the rigorous proof, which we omit, is that “ $\Lambda_t = \int_0^t \ell(s) ds \in \mathcal{F}_{t-}$ ”.
2. But again one can argue plausibly, starting with the comment that the increment over $(t, t + \Delta t)$ has conditional expectation $\int_t^{t+\Delta t} \ell(s) ds$ and takes values 0 or 1. Hence we can deduce the conditional probability of a +1-jump as being $\int_t^{t+\Delta t} \ell(s) ds$, and so argue as above.

Variance of compensated counting process
 The above expression of variance type holds more generally:
 Suppose N is a counting process compensated by $\int \ell(s) ds$. Then
 $(N_t - \int_0^t \ell(s) ds)^2 - \int_0^t \ell(s) ds$ is a martingale.
 Rigorous proof, or heuristic limiting argument

Counting processes and Poisson processes

The compensator of a counting process can be used to tell whether the counting process is Poisson:

Theorem

Suppose N is a counting process which has compensator αt . Then N is a Poisson process of rate α .

Better still, counting processes with compensators approximating αt are approximately Poisson of rate α . Here is a nice way to see this:

Theorem

Suppose N is a counting process with compensator $\Lambda = \int \ell(s) ds$. Consider the random time change $\tau(t) = \inf\{s : \Lambda_s = \alpha t\}$. Then the time-changed counting process $N_{\tau(t)}$ is Poisson of rate α .

The above gives a good pay-off for this theory.

Counting processes and Poisson processes
 The compensator of a counting process can be used to tell whether the counting process is Poisson.
Theorem
 Suppose N is a counting process which has compensator αt . Then N is a Poisson process of rate α .
 Better still, counting processes with compensators approximating αt are approximately Poisson of rate α . Here is a nice way to see this.
Theorem
 Suppose N is a counting process with compensator $\Lambda = \int \ell(s) ds$. Consider the random time change $\tau(t) = \inf\{s : \Lambda_s = \alpha t\}$. Then the time-changed counting process $N_{\tau(t)}$ is Poisson of rate α .
 The above gives a good pay-off for this theory.

- Again there is a plausibility argument: the increment over $(t, t + \Delta t)$ has conditional probability $\alpha \Delta t$, hence is approximately independent of past; hence N_t is approximately the sum of many Bernoulli random variables each of the same small mean, hence is approximately approximately Poisson
- Begs the question, is $N_{\tau(t)}$ a counting process? (Yes, but needs proof.)
- There is an amazing multivariate generalization of this time-change result, related to Cox's proportional hazards model.
- If the compensator approximates αt then it is immediate that $\tau(t)$ approximates t , and hence good approximation results can be derived!

Compensation of population processes

The notion of compensation works for much more general processes, such as population processes:

Example (Birth-death-immigration process)

If the birth-death-immigration process X makes transitions $X \rightarrow X + 1$ at rate $\lambda X + \alpha$ and $X \rightarrow X - 1$ at rate μX then

$$X_t - \int_0^t ((\lambda - \mu)X_s + \alpha) ds \text{ is a martingale.}$$

But the compensator no longer converts $(X_t - \int_0^t ((\lambda - \mu)X_s + \alpha) ds)^2$ into a martingale. More generally a continuous-time Markov chain X relates to martingales obtained from $f(X)$ (for given functions f) by compensation using the rates of X .

Compensation of population processes
 The notion of compensation works for much more general processes, such as population processes.
Example (Birth-death-immigration process)
 If the birth-death-immigration process X makes transitions $X \rightarrow X + 1$ at rate $\lambda X + \alpha$ and $X \rightarrow X - 1$ at rate μX then $X_t - \int_0^t ((\lambda - \mu)X_s + \alpha) ds$ is a martingale.
 But the compensator no longer converts $(X_t - \int_0^t ((\lambda - \mu)X_s + \alpha) ds)^2$ into a martingale. More generally a continuous-time Markov chain X relates to martingales obtained from $f(X)$ (for given functions f) by compensation using the rates of X .

- Plausibility argument much as before.
- The plausibility argument fails for the variance case! However it is possible to use a slightly different integral here. In fact

$$(X_t - \int_0^t ((\lambda - \mu)X_s + \alpha) ds)^2 - \int_0^t ((\lambda + \mu)X_s + \alpha) ds \text{ is a martingale.}$$

This is best understood using ideas of *stochastic integrals* (of rather simple form), which we will not explore here.

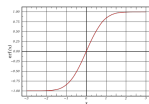
- This is the heart of the famous "Stroock-Varadhan martingale formulation", which allows one to use martingales to study and to define very general Markov chains.

Central Limit Theorem

"Everybody believes in the exponential law of errors: the experimenters, because they think it can be proved by mathematics; and the mathematicians, because they believe it has been established by observation"

Lippmann, quoted in E. T. Whittaker and G. Robinson, *Normal Frequency Distribution*. Ch. 8 in *The Calculus of Observations: A Treatise on Numerical Mathematics*, 1967.

Classical Central Limit Theorem
 Lindeberg's Central Limit Theorem
 Rates of convergence
 Martingale case



Central Limit Theorem
 The Central Limit Theorem is one of the jewels of classical probability theory, with a huge literature developing such questions as, how may the assumptions be relaxed? and how at what speed does the convergence actually occur?

The Central Limit Theorem is one of the jewels of classical probability theory, with a huge literature developing such questions as, how may the assumptions be relaxed? and how at what speed does the convergence actually occur?

The classical Central Limit Theorem

Definition

Random variables Y_n are said to converge in distribution to a random variable Z (or its distribution) if

$$\mathbb{P}[Y_n \leq y] \rightarrow \mathbb{P}[Z \leq y] \text{ whenever } \mathbb{P}[Z \leq y] \text{ is continuous at } y.$$

Theorem

Suppose X_1, \dots, X_n are independent and identically distributed, with finite mean μ and finite variance σ^2 . Then

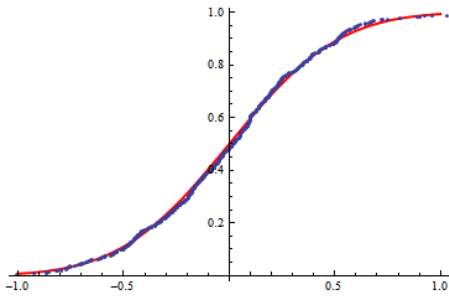
$$Y_n = \frac{X_1 + \dots + X_n - n\mu}{\sqrt{n}\sigma} \rightarrow N(0, 1),$$

where convergence is in distribution.

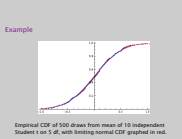
The classical Central Limit Theorem
Definition
 Random variables Y_n are said to converge in distribution to a random variable Z (or its distribution) if $\mathbb{P}[Y_n \leq y] \rightarrow \mathbb{P}[Z \leq y]$ whenever $\mathbb{P}[Z \leq y]$ is continuous at y .
Theorem
 Suppose X_1, \dots, X_n are independent and identically distributed, with finite mean μ and finite variance σ^2 . Then $Y_n = \frac{X_1 + \dots + X_n - n\mu}{\sqrt{n}\sigma} \rightarrow N(0, 1)$, where convergence is in distribution.

- $N(0, 1)$ denotes a random variable with standard normal distribution.
- Common notations: $Y_n \xrightarrow{d} Z$ or $Y_n \xrightarrow{D} Z$ or $Y_n \Rightarrow Z$.
- Cleanest proof involves *characteristic functions* $\mathbb{E}[\exp(iuY_n)]$, $\mathbb{E}[\exp(iuZ)] = e^{-u^2/2}$ and hence complex numbers. A Taylor series expansion shows $\mathbb{E}[\exp(iuX_n)] \approx 1 + iu\mu - \frac{u^2}{2}\sigma^2$; hence $\mathbb{E}[\exp(iuY_n)] \approx \left(1 - \frac{u^2}{2}\right)^n \rightarrow e^{-u^2/2}$. Result follows from theory of characteristic function transform.

Example



Empirical CDF of 500 draws from mean of 10 independent Student t on 5 df, with limiting normal CDF graphed in red.



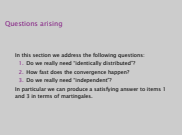
It is appropriate to use the CDF (cumulative distribution function) here, because that is the approximation which the CLT describes. Note there is good agreement!

Questions arising

In this section we address the following questions:

1. Do we really need "identically distributed"?
2. How fast does the convergence happen?
3. Do we really need "independent"?

In particular we can produce a satisfying answer to items 1 and 3 in terms of martingales.



1. No we don't need exactly "identically distributed", and we can produce a useful answer.
2. Something really rather definite can be said about rate of convergence.
3. No we do not need exactly "independent", and we can produce a useful answer.
4. (Our answer to items 1 and 3 is satisfying though not necessarily as good as possible!)

Lindeberg's Central Limit Theorem

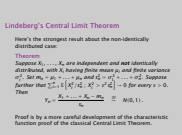
Here's the strongest result about the non-identically distributed case:

Theorem

Suppose X_1, \dots, X_n are independent and **not** identically distributed, with X_i having finite mean μ_i and finite variance σ_i^2 . Set $m_n = \mu_1 + \dots + \mu_n$ and $s_n^2 = \sigma_1^2 + \dots + \sigma_n^2$. Suppose further that $\sum_{i=1}^n \mathbb{E} \left[\frac{X_i^2}{s_n^2}; X_i^2 > \varepsilon^2 s_n^2 \right] \rightarrow 0$ for every $\varepsilon > 0$. Then

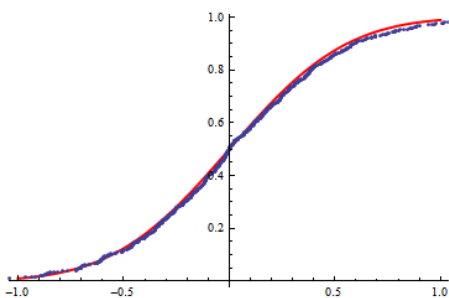
$$Y_n = \frac{X_1 + \dots + X_n - m_n}{s_n} \xrightarrow{D} N(0, 1).$$

Proof is by a more careful development of the characteristic function proof of the classical Central Limit Theorem.

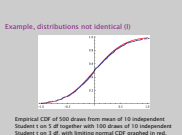


1. The beauty of the Lindeberg condition is that it simply requires none of the individual components to contribute too much to the total variance relative to the intended limit. Put this way, it is rather easy to remember the final result!
2. However the Lindeberg condition can be tricky to check. The Lyapunov condition is easier, and implies the Lindeberg condition: the sum of the third central moments $r_n^3 = \sum_{i=1}^n \mathbb{E} \left[|X_i - \mu_i|^3 \right]$ is finite and satisfies $(r_1 + \dots + r_n) / s_n \rightarrow 0$.

Example, distributions not identical (I)

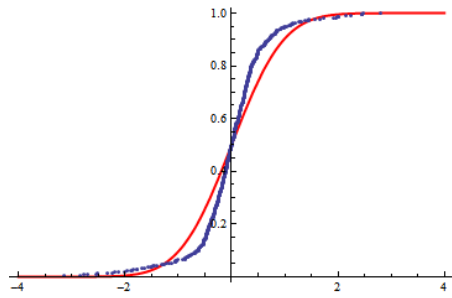


Empirical CDF of 500 draws from mean of 10 independent Student t on 5 df together with 100 draws of 10 independent Student t on 3 df, with limiting normal CDF graphed in red.

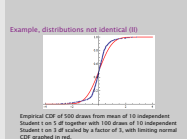


Here the distributions are not all the same. There is still reasonably good agreement!

Example, distributions not identical (II)



Empirical CDF of 500 draws from mean of 10 independent Student t on 5 df together with 100 draws of 10 independent Student t on 3 df scaled by a factor of 3, with limiting normal CDF graphed in red.



Now agreement is rather poorer.

Rates of convergence

Remarkably, we can capture how fast convergence occurs if we are given some extra information about the X_i . Reverting to the classical conditions (identically distributed, finite mean and variance), using above notation, suppose $\rho^{(3)} = \mathbb{E}[|X_i - \mu|^3] < \infty$. Let $F_n(x)$ be the distribution function of $\frac{(X_1 + \dots + X_n) - n\mu}{\sqrt{n\sigma}}$, and let $\Phi(x)$ be the standard normal distribution function. Then there is a universal constant $C > 0$ such that

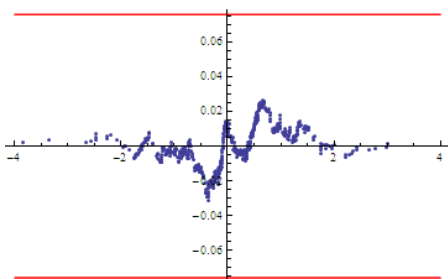
$$|F_n(x) - \Phi(x)| \leq \frac{C\rho^{(3)}}{\sigma^3\sqrt{n}}$$

Rates of convergence
 Remarkably, we can capture how fast convergence occurs if we are given some extra information about the X_i . Reverting to the classical conditions (identically distributed, finite mean and variance), using above notation, suppose $\rho^{(3)} = \mathbb{E}[|X_i - \mu|^3] < \infty$. Let $F_n(x)$ be the distribution function of $\frac{(X_1 + \dots + X_n) - n\mu}{\sqrt{n\sigma}}$, and let $\Phi(x)$ be the standard normal distribution function. Then there is a universal constant $C > 0$ such that

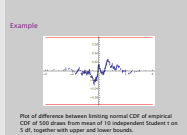
$$|F_n(x) - \Phi(x)| \leq \frac{C\rho^{(3)}}{\sigma^3\sqrt{n}}$$

There are many variants and many improvements on this result, whose proof requires much detailed mathematical analysis. For example, what is C ? (Latest: we can take $C = 0.7655$.) And so forth ...

Example



Plot of difference between limiting normal CDF of empirical CDF of 500 draws from mean of 10 independent Student t on 5 df, together with upper and lower bounds.



It is apparent that the bound on CLT discrepancy is not top bad ...

Martingale case

Theorem
 Suppose $X_0 = 0, X_1, \dots$ is a martingale for which $\mathbb{E}[X_n^2]$ is finite for each n . Set $s_n^2 = \mathbb{E}[X_n^2]$ and suppose $s_n^2 \rightarrow \infty$. The following two conditions taken together imply that X_n/s_n converges to a standard normal distribution:

$$\frac{1}{s_n^2} \sum_{m=0}^{n-1} \mathbb{E}[|X_{m+1} - X_m|^2 | \mathcal{F}_m] \rightarrow 1,$$

$$\frac{1}{s_n^2} \sum_{m=0}^{n-1} \mathbb{E}[|X_{m+1} - X_m|^2; |X_{m+1} - X_m|^2 \geq \varepsilon^2 s_n^2] \rightarrow 0 \text{ for each } \varepsilon > 0.$$

Martingale case
 Theorem
 Suppose $X_0 = 0, X_1, \dots$ is a martingale for which $\mathbb{E}[X_n^2]$ is finite for each n . Set $s_n^2 = \mathbb{E}[X_n^2]$ and suppose $s_n^2 \rightarrow \infty$. The following two conditions taken together imply that X_n/s_n converges to a standard normal distribution:

$$\frac{1}{s_n^2} \sum_{m=0}^{n-1} \mathbb{E}[|X_{m+1} - X_m|^2 | \mathcal{F}_m] \rightarrow 1,$$

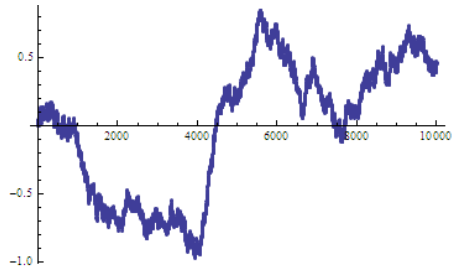
$$\frac{1}{s_n^2} \sum_{m=0}^{n-1} \mathbb{E}[|X_{m+1} - X_m|^2; |X_{m+1} - X_m|^2 \geq \varepsilon^2 s_n^2] \rightarrow 0 \text{ for each } \varepsilon > 0.$$

1. There are central limit theorems for martingales, typically close in spirit to the Lindeberg theorem. Namely: the total variance needs to be nearly constant, and there must be no relatively large contributions to the variance.
2. In fact $s_n^2 \rightarrow \infty$ is forced by the second (Lindeberg-type) condition.
3. Even more is true! the linear interpolation of the X_n , suitably rescaled, then converges to a Brownian motion.
4. There are many references, and many variations and generalizations. See for example Brown (1971).
 (Practical remarks about contrast between theory and practice ...)

- └ 5: Central Limit Theorem
- └ Martingale case

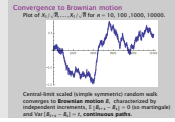
Convergence to Brownian motion

Plot of $X_1/\sqrt{n}, \dots, X_n/\sqrt{n}$ for $n = 10, 100, 1000, 10000$.



Central-limit scaled (simple symmetric) random walk converges to **Brownian motion** B , characterized by independent increments, $\mathbb{E}[B_{t+s} - B_s] = 0$ (so martingale) and $\text{Var}[B_{t+s} - B_s] = t$, **continuous paths**.

- └ 5: Central Limit Theorem
- └ Martingale case
- └ Convergence to Brownian motion



If paths weren't continuous, then compensated Poisson process would produce another example of this!
 In fact any random walk with jumps of zero mean and finite variance also converges to Brownian motion under central-limit scaling.
 There are also similar theorems for martingales . . .

- └ 6: Recurrence

Recurrence

"A bad penny always turns up"
 Old English proverb.

- Speed of convergence
- Irreducibility for general chains
- Regeneration and small sets
- Harris-recurrence
- Examples



- └ 6: Recurrence
- └ Recurrence



We have a theory of recurrence for discrete state space Markov chains. But what if the state space is not discrete? and how can we describe speed of convergence?

- └ 6: Recurrence

Motivation from MCMC

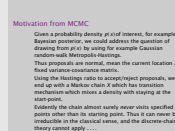
Given a probability density $p(x)$ of interest, for example a Bayesian posterior, we could address the question of drawing from $p(x)$ by using for example Gaussian random-walk Metropolis-Hastings.

Thus proposals are normal, mean the current location x , fixed variance-covariance matrix.

Using the Hastings ratio to accept/reject proposals, we end up with a Markov chain X which has transition mechanism which mixes a density with staying at the start-point.

Evidently the chain almost surely *never* visits specified points other than its starting point. Thus it can never be irreducible in the classical sense, and the discrete-chain theory cannot apply . . .

- └ 6: Recurrence
- └ Motivation from MCMC



1. Clearly the discrete-chain theory needs major rehabilitation if it is to be helpful in the continuous state space case!

- └ 6: Recurrence

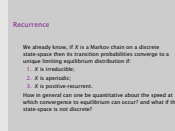
Recurrence

We already know, if X is a Markov chain on a discrete state-space then its transition probabilities converge to a unique limiting equilibrium distribution if:

1. X is irreducible;
2. X is aperiodic;
3. X is positive-recurrent.

How in general can one be quantitative about the speed at which convergence to equilibrium can occur? and what if the state-space is not discrete?

- └ 6: Recurrence
- └ Recurrence



Recurrence and rates of convergence for Markov chains in discrete case (uniform and geometric ergodicity). Making sense of continuous state-space, ϕ -irreducibility, Harris-recurrence. Small sets. Application to important examples.

1. the state space of X cannot be divided into regions some of which are inaccessible from others;
2. the state space of X cannot be broken into periodic cycles;
3. the mean time for X to return to its starting point is finite.

Measuring speed of convergence to equilibrium (I)

Total variation distance

- ▶ Speed of convergence of a Markov chain X to equilibrium can be measured as discrepancy between two probability measures: $\mathcal{L}(X_t|X_0 = x)$ (distribution of X_t) and π (equilibrium measure).
- ▶ Simple possibility: **total variation distance**. Let \mathcal{X} be state-space, for $A \subseteq \mathcal{X}$ maximize discrepancy between $\mathcal{L}(X_t|X_0 = x)(A) = \mathbb{P}[X_t \in A|X_0 = x]$ and $\pi(A)$:

$$\text{dist}_{TV}(\mathcal{L}(X_t|X_0 = x), \pi) = \sup_{A \subseteq \mathcal{X}} \{ \mathbb{P}[X_t \in A|X_0 = x] - \pi(A) \}.$$

- ▶ Alternative expression in case of discrete state-space:

$$\text{dist}_{TV}(\mathcal{L}(X_t|X_0 = x), \pi) = \frac{1}{2} \sum_{y \in \mathcal{X}} | \mathbb{P}[X_t = y|X_0 = x] - \pi_y |.$$

(Many other possible measures of distance . . .)

Measuring speed of convergence to equilibrium (I)
 Definition
 • Speed of convergence of a Markov chain X to equilibrium can be measured as discrepancy between two probability measures: $\mathcal{L}(X_t|X_0 = x)$ (distribution of X_t) and π (equilibrium measure).
 • Simple possibility: **total variation distance**. Let \mathcal{X} be state-space. For $A \subseteq \mathcal{X}$ maximize discrepancy between $\mathcal{L}(X_t|X_0 = x)(A) = \mathbb{P}[X_t \in A|X_0 = x]$ and $\pi(A)$.

$$\text{dist}_{TV}(\mathcal{L}(X_t|X_0 = x), \pi) = \sup_{A \subseteq \mathcal{X}} \{ \mathbb{P}[X_t \in A|X_0 = x] - \pi(A) \}.$$

 • Alternative expression in case of discrete state-space:

$$\text{dist}_{TV}(\mathcal{L}(X_t|X_0 = x), \pi) = \frac{1}{2} \sum_{y \in \mathcal{X}} | \mathbb{P}[X_t = y|X_0 = x] - \pi_y |.$$

 (Many other possible measures of distance . . .)

1. $\mathcal{L}(X_t|X_0 = x)(A)$ is probability that X_t belongs to A .
2. **Test understanding:** why is it not necessary to consider $|\mathbb{P}[X_t \in A|X_0 = x] - \pi(A)|$?
 (Hint: consider $\mathbb{P}[X_t \in A^c|X_0 = x] - \pi(A^c)$.)
3. **Test understanding:** prove this by considering $A = \{y : \mathbb{P}[X_t = y|X_0 = x] > \pi_y\}$.
4. It is not even clear that total variation is best notion: in the case of MCMC one might consider a spectral approach (which we will pick up again when we come to consider cutoff):

$$\sup_{f: \int f(x)^2 \pi(dx) < \infty} \left(\mathbb{E}[f(X_t)|X_0 = x] - \int f(x)\pi(dx) \right)^2.$$

5. Nevertheless the concept of total variation isolates a desirable kind of rapid convergence.

Measuring speed of convergence to equilibrium (II)

Uniform ergodicity

Definition

The Markov chain X is **uniformly ergodic** if its distribution converges to equilibrium in total variation *uniformly in the starting point* $X_0 = x$: for some fixed $C > 0$ and for fixed $\gamma \in (0, 1)$,

$$\sup_{x \in \mathcal{X}} \text{dist}_{TV}(\mathcal{L}(X_n|X_0 = x), \pi) \leq C\gamma^n.$$

In theoretical terms, for example when carrying out MCMC, this is a very satisfactory property. No account need be taken of the starting point, and accuracy improves in proportion to the length of the simulation.

Measuring speed of convergence to equilibrium (II)
 Definition
 The Markov chain X is **uniformly ergodic** if its distribution converges to equilibrium in total variation *uniformly in the starting point* $X_0 = x$: for some fixed $C > 0$ and for fixed $\gamma \in (0, 1)$.

$$\sup_{x \in \mathcal{X}} \text{dist}_{TV}(\mathcal{L}(X_n|X_0 = x), \pi) \leq C\gamma^n.$$

 In theoretical terms, for example when carrying out MCMC, this is a very satisfactory property. No account need be taken of the starting point, and accuracy improves in proportion to the length of the simulation.

1. In fact this is a consequence of the apparently weaker assertion, as $n \rightarrow \infty$ so

$$\sup_{x \in \mathcal{X}} \text{dist}_{TV}(\mathcal{L}(X_t|X_0 = x), \pi) \rightarrow 0.$$

2. Much depends on size of C and on how small is γ .
3. Typically theoretical estimates of C and γ are very conservative.
4. Other things being equal(!), given a choice, consider choosing a uniformly ergodic Markov chain for your MCMC algorithm.

Measuring speed of convergence to equilibrium (III)

Geometric ergodicity

Definition

The Markov chain X is **geometrically ergodic** if its distribution converges to equilibrium in total variation for some $C(x) > 0$ *depending on the starting point* x and for fixed $\gamma \in (0, 1)$,

$$\text{dist}_{TV}(\mathcal{L}(X_t|X_0 = x), \pi) \leq C(x)\gamma^n.$$

Here account does need to be taken of the starting point, but still accuracy improves in proportion to the length of the simulation.

Measuring speed of convergence to equilibrium (III)
 Definition
 The Markov chain X is **geometrically ergodic** if its distribution converges to equilibrium in total variation for some $C(x) > 0$ depending on the starting point x and for fixed $\gamma \in (0, 1)$.

$$\text{dist}_{TV}(\mathcal{L}(X_n|X_0 = x), \pi) \leq C(x)\gamma^n.$$

 Here account does need to be taken of the starting point, but still accuracy improves in proportion to the length of the simulation.

A significant question is, how might one get a sense of whether a specified chain *is* indeed geometrically ergodic (because at least that indicates the rate at which the distribution of X_t gets closer to equilibrium) and how one might obtain upper bounds on γ . We shall see later on that even given good information about γ and C , and even if total variation is of primary interest, geometric ergodicity still leaves important phenomena untouched!

ϕ -irreducibility (I)

We make two observations about Markov chain irreducibility:

1. The discrete theory fails to apply directly even to well-behaved chains on non-discrete state-space.
2. Suppose ϕ is a measure on the state-space: then we could ask for the chain to be irreducible *on sets of positive ϕ measure*.

Definition

The Markov chain X is **ϕ -irreducible** if for any state x and for any subset B of state-space of positive ϕ -measure $\phi(B) > 0$ we find that X has positive chance of reaching B if begun at x .

ϕ -irreducibility (I)
 We make two observations about Markov chain irreducibility:
 1. The discrete theory fails to apply directly even to well-behaved chains on non-discrete state-space.
 2. Suppose ϕ is a measure on the state-space: then we could ask for the chain to be irreducible *on sets of positive ϕ measure*.
 Definition
 The Markov chain X is **ϕ -irreducible** if for any state x and for any subset B of state-space of positive ϕ -measure $\phi(B) > 0$ we find that X has positive chance of reaching B if begun at x .

1. Consider the Gaussian random walk X (jumps have standard normal distribution): if $X_0 = 0$ then we can assert that with probability one X *never* returns to its starting point.
2. "measure": like a probability measure, but not necessarily of finite total mass. Think of length, area, or volume as examples. Also, counting measure.
3. The Gaussian random walk is Lebesgue-measure-irreducible! (Here Lebesgue measure is just length measure.)

φ-irreducibility (II)

1. We call ϕ an *irreducibility measure*. It is possible to modify ϕ to construct a *maximal irreducibility measure* ψ ; one such that any set B of positive measure under some irreducibility measure for X is of positive measure for ψ .
2. Irreducible chains on countable state-space are c -irreducible where c is counting measure.
3. If a chain has unique equilibrium measure π then π will serve as a maximal irreducibility measure.

φ-irreducibility (II)

1. We call ϕ an irreducibility measure. It is possible to modify ϕ to construct a maximal irreducibility measure ψ ; one such that any set B of positive measure under some irreducibility measure for X is of positive measure for ψ .
2. Irreducible chains on countable state-space are c -irreducible where c is counting measure.
3. If a chain has unique equilibrium measure π then π will serve as a maximal irreducibility measure.

1. Lebesgue measure is a maximal irreducibility measure for the Gaussian random walk.
2. So ϕ -irreducibility simply generalizes the original notion of irreducibility.
3. Note that ϕ can be replaced by any other measure which is "measure-equivalent" (has the same null-sets). So while π will serve as a maximal irreducibility measure, we can use any alternative measure which has the same sets of measure zero.

Regeneration and small sets (I)

The discrete-state-space theory works because (a) the Markov chain *regenerates* each time it visits individual states, and (b) it has a positive chance of visiting specified individual states. So it is natural to consider regeneration when visiting sets.

Definition

A set E of ϕ -positive measure is a *small set of lag k* for X if there is $\alpha \in (0, 1)$ and a probability measure ν such that for all $x \in E$ the following *minorization condition* is satisfied

$$\mathbb{P}[X_k \in A | X_0 = x] \geq \alpha \nu(A) \quad \text{for all } A.$$

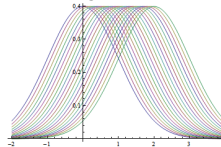
Regeneration and small sets (I)

The discrete-state-space theory works because (a) the Markov chain regenerates each time it visits individual states, and (b) it has a positive chance of visiting specified individual states. So it is natural to consider regeneration when visiting sets.

Definition
 A set E of ϕ -positive measure is a small set of lag k for X if there is $\alpha \in (0, 1)$ and a probability measure ν such that for all $x \in E$ the following minorization condition is satisfied

$$\mathbb{P}[X_k \in A | X_0 = x] \geq \alpha \nu(A) \quad \text{for all } A.$$

1. In effect this reduces the theory of convergence to equilibrium to a chapter in the theory of renewal processes, with renewals occurring each time the chain visits a specified state.
2. In effect, if we *sub-sample* X every k time-steps then, every time it visits E , there is a chance α that X forgets its entire past and starts again, using probability measure ν . Consider the Gaussian random walk described above. **Any bounded set is small of lag 1.**



3. In general α can be very small—reducing practical impact, but still helping theoretically.

Regeneration and small sets (II)

Small sets would not be interesting except that:

1. all ϕ -irreducible Markov chains X possess small sets;
2. consider chains X with continuous transition density kernels. They possess *many* small sets of lag 1;
3. consider chains X with measurable transition density kernels. They need possess *no* small sets of lag 1, but will possess many sets of lag 2;
4. given just one small set, X can be represented using a chain which has a single recurrent atom.

In a word, small sets *discretize* Markov chains.



▶ ANIMATION

Regeneration and small sets (II)

Small sets would not be interesting except that:

1. all ϕ -irreducible Markov chains X possess small sets;
2. consider chains X with continuous transition density kernels. They possess many small sets of lag 1;
3. consider chains X with measurable transition density kernels. They need possess *no* small sets of lag 1, but will possess many sets of lag 2;
4. given just one small set, X can be represented using a chain which has a single recurrent atom.

In a word, small sets *discretize* Markov chains.

1. This is a very old result: see Nummelin (1984) for a recent treatment.
2. Exercise: try seeing why this is obviously true!
3. Kendall and Montana (2002): so measurable transition density kernels lead to chains which possess latent discretizations.
4. "Split-chain construction" (Athreya and Ney 1978; Nummelin 1978).

Harris-recurrence

Now it is evident what we should mean by recurrence for non-discrete state spaces. Suppose X is ϕ -irreducible and ϕ is a maximal irreducibility measure.

Definition

X is (ϕ) -*recurrent* if, for ϕ -almost all starting points x and any subset B with $\phi(B) > 0$, when started at x the chain X is almost sure eventually to hit B .

Definition

X is *Harris-recurrent* if we can drop " ϕ -almost" in the above.

Harris-recurrence

Now it is evident what we should mean by recurrence for non-discrete state spaces. Suppose X is ϕ -irreducible and ϕ is a maximal irreducibility measure.

Definition
 X is (ϕ) -recurrent if, for ϕ -almost all starting points x and any subset B with $\phi(B) > 0$, when started at x the chain X is almost sure eventually to hit B .

Definition
 X is Harris-recurrent if we can drop " ϕ -almost" in the above.

1. So the irreducibility measure is used to focus attention on sets rather than points.
2. And in fact we don't even then need ϕ to be maximal.

Examples of ϕ -irreducibility

- ▶ Random walks with continuous jump densities. And in fact measurable jump densities suffice.
- ▶ Chains with continuous or even measurable transition densities with exception that chain may stay put.
- ▶ Vervaat perpetuities:

$$X_{n+1} = U_{n+1}^\alpha (X_n + 1)$$

where U_1, U_2, \dots are independent Uniform(0, 1).

- ▶ Volatility models:

$$\begin{aligned} X_{n+1} &= X_n + \sigma_n Z_{n+1} \\ \sigma_{n+1} &= f(\sigma_n, U_{n+1}) \end{aligned}$$

for suitable f , and independent Gaussian Z_{n+1}, U_{n+1} .

Examples of ϕ -irreducibility

- Random walks with continuous jump densities. And in fact measurable jump densities suffice.
- Chains with continuous or even measurable transition densities with exception that chain may stay put.
- Vervaat perpetuities:

$$X_{n+1} = U_{n+1}^\alpha (X_n + 1)$$
 where U_1, U_2, \dots are independent Uniform(0, 1).
- Volatility models:

$$\begin{aligned} X_{n+1} &= X_n + \sigma_n Z_{n+1} \\ \sigma_{n+1} &= f(\sigma_n, U_{n+1}) \end{aligned}$$
 for suitable f , and independent Gaussian Z_{n+1}, U_{n+1} .

1. Convolutions of measurable densities are continuous!
2. Many examples of Metropolis-Hastings samplers.
3. **Test understanding:** find a small set for the Vervaat perpetuity example!

Foster-Lyapunov criteria

"Even for the physicist the description in plain language will be the criterion of the degree of understanding that has been reached."
 Werner Heisenberg,
Physics and philosophy: The revolution in modern science, 1958

- Renewal and regeneration
- Positive recurrence
- Geometric ergodicity
- Examples



Foster-Lyapunov criteria

Foster-Lyapunov criteria

The Foster-Lyapunov criterion for positive recurrence of a ϕ -irreducible Markov chain X on a state space X .

Positive recurrence

Given $b, c > 0$, positive constants a, b, c , and a small set $C = \{x : \Lambda(x) \leq c\} \subseteq X$, with

then $\mathbb{E}[T_A | X_0 = x] < \infty$ for any A with $\phi(A) > 0$, where $T_A = \inf\{n \geq 0 : X_n \in A\}$ is the time when X first hits A , and moreover X has an equilibrium distribution.

Geometric and uniform ergodicity make sense for general Markov chains: how to find out whether they hold? and how to find out whether equilibrium distributions exist?
 We want simple criteria, and we can capture these using the language of martingales.

Renewal and regeneration

Suppose C is a small set for ϕ -recurrent X , with lag 1:

$$\mathbb{P}[X_1 \in A | X_0 = x \in C] \geq \alpha \nu(A).$$

Identify **regeneration events**: X regenerates at $x \in C$ with probability α and then makes transition with distribution ν ; otherwise it makes transition with distribution $\frac{p(x, \cdot) - \alpha \nu(\cdot)}{1 - \alpha}$. The regeneration events occur as a **renewal sequence**. Set

$$p_k = \mathbb{P}[\text{next regeneration at time } k | \text{regeneration at time } 0].$$

If the renewal sequence is **non-defective** if $\sum_k p_k = 1$ and **positive-recurrent** if $\sum_k k p_k < \infty$ then there exists a stationary version. This is the key to equilibrium theory whether for discrete or continuous state-space.

Renewal and regeneration

Suppose C is a small set for ϕ -recurrent X , with lag 1:

$$\mathbb{P}[X_1 \in A | X_0 = x \in C] \geq \alpha \nu(A).$$

Identify regeneration events: X regenerates at $x \in C$ with probability α and then makes transition with distribution ν ; otherwise it makes transition with distribution $\frac{p(x, \cdot) - \alpha \nu(\cdot)}{1 - \alpha}$. The regeneration events occur as a renewal sequence. Set

$$p_k = \mathbb{P}[\text{next regeneration at time } k | \text{regeneration at time } 0].$$

If the renewal sequence is non-defective ($\sum_k p_k = 1$) and positive-recurrent ($\sum_k k p_k < \infty$), then there exists a stationary version. This is the key to equilibrium theory whether for discrete or continuous state space.

1. If lag is $k > 1$ then sub-sample every k steps!
2. This is a coupling construction, linked to the split-chain construction (Athreya and Ney 1978; Nummelin 1978) and the Murdoch and Green (1998) approach to CFTP.
3. This is just the appropriate compensating distribution

$$\frac{p(x, \cdot) - \alpha \nu(\cdot)}{p(x, X) - \alpha \nu(X)} = \frac{p(x, \cdot) - \alpha \nu(\cdot)}{1 - \alpha}.$$

4. So there will always be a next regeneration.
5. So mean time to next regeneration is finite.
6. Richard Tweedie at WRASS 1998: "continuous is no harder than discrete!"

Positive recurrence

The Foster-Lyapunov criterion for positive recurrence of a ϕ -irreducible Markov chain X on a state-space X :

Theorem (Foster-Lyapunov criterion for positive recurrence)

Given $\Lambda : X \rightarrow [0, \infty)$, positive constants a, b, c , and a small set $C = \{x : \Lambda(x) \leq c\} \subseteq X$ with

$$\mathbb{E}[\Lambda(X_{n+1}) | \mathcal{F}_n] \leq \Lambda(X_n) - a + b \mathbb{1}_{[X_n \in C]};$$

then $\mathbb{E}[T_A | X_0 = x] < \infty$ for any A with $\phi(A) > 0$, where $T_A = \inf\{n \geq 0 : X_n \in A\}$ is the time when X first hits A , and moreover X has an equilibrium distribution.

Positive recurrence

The Foster-Lyapunov criterion for positive recurrence of a ϕ -irreducible Markov chain X on a state space X .

Positive recurrence

Given $b, c > 0$, positive constants a, b, c , and a small set $C = \{x : \Lambda(x) \leq c\} \subseteq X$, with

$$\mathbb{E}[\Lambda(X_{n+1}) | \mathcal{F}_n] \leq \Lambda(X_n) - a + b \mathbb{1}_{[X_n \in C]}$$

then $\mathbb{E}[T_A | X_0 = x] < \infty$ for any A with $\phi(A) > 0$, where $T_A = \inf\{n \geq 0 : X_n \in A\}$ is the time when X first hits A , and moreover X has an equilibrium distribution.

1. In words, we can find a non-negative $\Lambda(X)$ such that $\Lambda(X_n) - an$ determines a supermartingale until $\Lambda(X)$ becomes small enough for X to belong to a small set!
2. We can re-scale Λ so that $a = 1$.
3. In fact if the criterion holds it can then be shown, *any* sub-level set of Λ is small.
4. It is evident from the verbal description that reflected simple asymmetric random walk (negatively biased) is an example for which the criterion applies.

Sketch of proof

Proof.

- $Y_n = \Lambda(X_n) + an$ is non-negative supermartingale up to time $T = \inf\{m \geq 0 : X_m \in C\} > n$:

$$\mathbb{E}[Y_{\min\{n+1, T\}} | \mathcal{F}_n, T > n] \leq (\Lambda(X_n) - a) + a(n+1) = Y_n.$$

Hence $Y_{\min\{n, T\}}$ converges.

- So $\mathbb{P}[T < \infty] = 1$ (for otherwise $\Lambda(X) > c$ and $Y_n > c + an$). Moreover $\mathbb{E}[Y_T | X_0] \leq \Lambda(X_0)$.
- Now use finiteness of b to show $\mathbb{E}[T^* | X_0] < \infty$, where T^* first regeneration in C .
- ϕ -irreducibility: positive chance of hitting A before first regeneration in C . Hence $\mathbb{E}[T_A | X_0] < \infty$.

- There is a stationary version of the renewal process of successive regenerations on C .
- One can construct a "bridge" of X conditioned to regenerate on C at time 0 , and then to regenerate again on C at time n .
- Hence one can sew these together to form a stationary version of X , which therefore has the property that X_t has the equilibrium distribution for all time t .

A converse ...

Suppose on the other hand that $\mathbb{E}[T | X_0] < \infty$ for all starting points X_0 , where C is some small set and T is the first time for X to return to C . The Foster-Lyapunov criterion for positive recurrence follows for $\Lambda(x) = \mathbb{E}[T | X_0 = x]$ if $\mathbb{E}[T | X_0]$ is bounded on C .

- ϕ -irreducibility then follows automatically.
- Indeed, (supposing lag 1 for simplicity)

$$\mathbb{E}[\Lambda(X_{n+1}) | \mathcal{F}_n] \leq \Lambda(X_n) - 1 + b \mathbb{1}_{\{X_n \in C\}},$$

- where b is the mean value of $\mathbb{E}[Y_T | x]$ if x is chosen using the regeneration probability measure for C .
- Moreover if the renewal process of successive regenerations on C is aperiodic then a coupling argument shows general X will converge to equilibrium.
 - If the renewal process of successive regenerations on C is not aperiodic then one can sub-sample ...
 - Showing that X has an equilibrium is then a matter of probabilistic constructions using the renewal process of successive regenerations on C .

Geometric ergodicity

The Foster-Lyapunov criterion for geometric ergodicity of a ϕ -irreducible Markov chain X on a state-space \mathcal{X} :

Theorem (Foster-Lyapunov criterion for geometric ergodicity)

Given $\Lambda : \mathcal{X} \rightarrow [1, \infty)$, positive constants $\gamma \in (0, 1)$, $b, c \geq 1$, and a small set $C = \{x : \Lambda(x) \leq c\} \subseteq \mathcal{X}$ with

$$\mathbb{E}[\Lambda(X_{n+1}) | \mathcal{F}_n] \leq \gamma \Lambda(X_n) + b \mathbb{1}_{\{X_n \in C\}};$$

then $\mathbb{E}[\gamma^{-T_A} | X_0 = x] < \infty$ for any A with $\phi(A) > 0$, where $T_A = \inf\{n \geq 0 : X_n \in A\}$ is the time when X first hits A , and moreover (under suitable periodicity conditions) X is geometrically ergodic.

- In words, we can find a $\Lambda(X) \geq 1$ such that $\Lambda(X_n)/\gamma^n$ determines a supermartingale until $\Lambda(X)$ becomes small enough for X to belong to a small set!
- We can rescale Λ so that $b = 1$.
- The criterion for positive-recurrence is implied by this criterion.
- We can enlarge C and alter b so that the criterion holds simultaneously for all $\mathbb{E}[\Lambda(X_{n+m}) | \mathcal{F}_n]$.

Sketch of proof

Proof.

- $Y_n = \Lambda(X_n)/\gamma^n$ defines non-negative supermartingale up to time T when X first hits C :

$$\mathbb{E}[Y_{\min\{n+1, T\}} | \mathcal{F}_n, T > n] \leq \gamma \times \Lambda(X_n)/\gamma^{n+1} = Y_n.$$

Hence $Y_{\min\{n, T\}}$ converges.

- $\mathbb{P}[T < \infty] = 1$, for otherwise $\Lambda(X) > c$ and so $Y_n > c/\gamma^n$ does not converge. Moreover $\mathbb{E}[\gamma^{-T} | X_0] \leq \Lambda(X_0)$.
- Finiteness of b shows $\mathbb{E}[\gamma^{-T^*} | X_0] < \infty$, where T^* is time of regeneration in C .
- From ϕ -irreducibility there is positive chance of hitting A before regeneration in C . Hence $\mathbb{E}[\gamma^{-T_A} | X_0] < \infty$.

- Geometric ergodicity follows by a coupling argument which I do not specify here.
- The constant γ here provides an upper bound on the constant γ used in the definition of geometric ergodicity. However it is not necessarily a very good bound!

Two converses

- Suppose on the other hand that $\mathbb{E}[y^{-T}|X_0] < \infty$ for all starting points X_0 (and fixed $y \in (0, 1)$), where C is some small set and T is the first time for X to return to C . The Foster-Lyapunov criterion for geometric ergodicity then follows for $\Lambda(x) = \mathbb{E}[y^{-T}|X_0 = x]$ if $\mathbb{E}[y^{-T}|X_0]$ is bounded on C .

Uniform ergodicity follows if the Λ function is bounded above.

But more is true. Strikingly,

- For Harris-recurrent Markov chains the existence of a geometric Foster-Lyapunov condition is **equivalent** to the property of geometric ergodicity.

Examples

- General reflected random walk:**
 $X_{n+1} = \max\{X_n + Z_{n+1}, 0\}$ with independent Z_{n+1} of continuous density, $\mathbb{E}[Z_{n+1}] < 0$. Then
 - X is Lebesgue-irreducible on $[0, \infty)$;
 - Foster-Lyapunov criterion for positive recurrence applies.

Similar considerations often apply to Metropolis-Hastings Markov chains based on random walks.

- Reflected Simple Asymmetric Random Walk:**
 $X_{n+1} = \max\{X_n + Z_{n+1}, 0\}$ with independent Z_{n+1} such that $\mathbb{P}[Z_{n+1} = -1] = q = 1 - p = 1 - \mathbb{P}[Z_{n+1} = +1] > \frac{1}{2}$.
 - X is counting-measure-irreducible on non-negative integers;
 - Foster-Lyapunov criterion for geometric ergodicity applies.
 Aim for $\mathbb{E}[e^{aZ_{n+1}}] < 1$ for some positive a .

Reflected Simple asymmetric random walk (II)

- Positive recurrence criterion: check for $\Lambda(x) = x$, $C = \{0\}$:

$$\mathbb{E}[\Lambda(X_1)|X_0 = x_0] = \begin{cases} \Lambda(x_0) - (q - p) & \text{if } x_0 \notin C, \\ 0 + p & \text{if } x_0 \in C. \end{cases}$$

- Geometric ergodicity criterion: check for $\Lambda = e^{ax}$, $C = \{0\} = \Lambda^{-1}(\{1\})$:

$$\mathbb{E}[\Lambda(X_1)|X_0 = x_0] = \begin{cases} \Lambda(x_0) \times (pe^a + qe^{-a}) & \text{if } x_0 \notin C, \\ 1 \times (p + qe^{-a}) & \text{if } x_0 \in C. \end{cases}$$

This works when $pe^a + qe^{-a} < 1$; equivalently when $0 < a < \log(q/p)$ (solve the quadratic in e^a !).

Cutoff

"I have this theory of convergence, that good things always happen with bad things."
 Cameron Crowe, *Say Anything* film, 1989

The cutoff phenomenon
 Cutoff and eigenvalues
 Two metrics
 A special case



- This was used in [Kendall 2004](#) to provide perfect simulation *in principle*. The Markov inequality can be used to convert the condition on $\Lambda(X)$ into the existence of a Markov chain on $[0, \infty)$ whose exponential dominates $\Lambda(X)$. The chain in question turns out to be a kind of queue (in fact, $D/M/1$). For $y \geq e^{-1}$ the queue will not be recurrent; however one can sub-sample X to convert the situation into one in which the dominating queue will be positive-recurrent.

Two converses
 1. Suppose on the other hand that $\mathbb{E}[y^{-T}|X_0] < \infty$ for all starting points X_0 (and fixed $y \in (0, 1)$), where C is some small set and T is the first time for X to return to C . The Foster-Lyapunov criterion for geometric ergodicity then follows for $\Lambda(x) = \mathbb{E}[y^{-T}|X_0 = x]$ if $\mathbb{E}[y^{-T}|X_0]$ is bounded on C .
 Uniform ergodicity follows if the Λ function is bounded above.
 But more is true. Strikingly,
 2. For Harris-recurrent Markov chains the existence of a geometric Foster-Lyapunov condition is equivalent to the property of geometric ergodicity.

It is instructive to notice that the criteria continue to apply to a considerable variety of appropriately modified Markov chains.

- Test understanding:** Lebesgue-irreducibility follows from continuous jump density by writing down chains of transitions;
 - Test understanding:** Check Foster-Lyapunov criterion for positive recurrence for $\Lambda(x) = x$.
- Test understanding:** this is the same as ordinary irreducibility for discrete-state-space Markov chains!
 - Test understanding:** Check Foster-Lyapunov criterion for geometric ergodicity for $\Lambda(x) = e^{ax}$ for small positive a .

(Further practical remarks about contrast between theory and practice . . .)

Examples
 1. General reflected random walk:
 $X_{n+1} = \max\{X_n + Z_{n+1}, 0\}$ with independent Z_{n+1} of continuous density, $\mathbb{E}[Z_{n+1}] < 0$. Then
 (a) X is Lebesgue-irreducible on $[0, \infty)$;
 (b) Foster-Lyapunov criterion for positive recurrence applies. Similar considerations often apply to Metropolis-Hastings Markov chains based on random walks.
 2. Reflected simple asymmetric random walk:
 $X_{n+1} = \max\{X_n + Z_{n+1}, 0\}$ with independent Z_{n+1} such that $\mathbb{P}[Z_{n+1} = -1] = q = 1 - p = 1 - \mathbb{P}[Z_{n+1} = +1] > \frac{1}{2}$.
 (a) X is counting-measure-irreducible on non-negative integers.
 (b) Foster-Lyapunov criterion for geometric ergodicity applies.
 Aim for $\mathbb{E}[e^{aZ_{n+1}}] < 1$ for some positive a .

One may ask, does this kind of argument show that *all* positive-recurrent random walks can be shown to be geometrically ergodic simply by moving from $\Lambda(x) = x$ to $\Lambda(x) = e^{ax}$? The answer is no, essentially because there exist random walks whose jump distributions have negative mean but fail to have exponential moments

Reflected Simple asymmetric random walk (II)
 Positive recurrence criterion: check for $\Lambda(x) = x$, $C = \{0\}$:
 $\mathbb{E}[\Lambda(X_1)|X_0 = x_0] = \begin{cases} \Lambda(x_0) - (q - p) & \text{if } x_0 \notin C, \\ 0 + p & \text{if } x_0 \in C. \end{cases}$
 Geometric ergodicity criterion: check for $\Lambda = e^{ax}$, $C = \{0\} = \Lambda^{-1}(\{1\})$:
 $\mathbb{E}[\Lambda(X_1)|X_0 = x_0] = \begin{cases} \Lambda(x_0) \times (pe^a + qe^{-a}) & \text{if } x_0 \notin C, \\ 1 \times (p + qe^{-a}) & \text{if } x_0 \in C. \end{cases}$
 This works when $pe^a + qe^{-a} < 1$, equivalently when $0 < a < \log(q/p)$ (solve the quadratic in e^a !).

In what way does a Markov chain converge to equilibrium? Is it a gentle exponential process? Or might most of the convergence happen relatively quickly?
 Once again we focus on reversible Markov chains, as these make computations simpler.

Cutoff
 "I have this theory of convergence, that good things always happen with bad things."
 Cameron Crowe, *Say Anything* film, 1989
 The cutoff phenomenon
 Cutoff and eigenvalues
 Two metrics
 A special case

- └ 8: Cutoff
- └ The cutoff phenomenon

Convergence: cutoff or geometric decay?

What we have so far said about convergence to equilibrium will have left the misleading impression that the distance from equilibrium for a Markov chain is characterized by a gentle and rather geometric decay. It is true that this is typically the case after an extremely long time, and it can be the case over all time. However it is entirely possible for "most" of the convergence to happen quite suddenly at a specific threshold. The theory for this is developing fast, but many questions remain open. In this section we describe a specific easy example.

- └ 8: Cutoff
- └ The cutoff phenomenon
- └ Convergence: cutoff or geometric decay?

Random walk wrapped around a circle exhibits a gentle and rather geometric decay. Famously (Bayer and Diaconis 1992) the riffle shuffle does not! (For a pack of 52 cards, 7 shuffles suffices for essentially all practical purposes.)

What we have to be said about convergence to equilibrium will have left the misleading impression that the distance from equilibrium for a Markov chain is characterized by a gentle and rather geometric decay. It is true that this is typically the case after an extremely long time, and it can be the case over all time. However it is entirely possible for "most" of the convergence to happen quite suddenly at a specific threshold. The theory for this is developing fast, but many questions remain open. In this section we describe a specific easy example.

- └ 8: Cutoff
- └ Cutoff and eigenvalues

Cutoff (I): Markov chains and matrices

We need to understand something about eigenvalues for Markov chains. Fix attention on a finite state space X , with reversible aperiodic Markov chain of transition kernel $p_{x,y}$ and equilibrium distribution π . The vector space of functions on X can be give a weighted Euclidean norm:

$$\|f\|_{\pi}^2 = \sum_{x \in X} |f(x)|^2 \pi(x)$$

and hence an inner product $\langle f, g \rangle_{\pi}$. View transition kernel as linear operator $Pf(x) = \sum_y p_{x,y} f(y)$: by reversibility this is $\langle \cdot, \cdot \rangle_{\pi}$ symmetric.

- └ 8: Cutoff
- └ Cutoff and eigenvalues
- └ Cutoff (I): Markov chains and matrices

Cutoff (I): Markov chains and matrices
We need to understand something about eigenvalues for Markov chains. Fix attention on a finite state space X , with reversible aperiodic Markov chain of transition kernel $p_{x,y}$ and equilibrium distribution π . The vector space of functions on X can be give a weighted Euclidean norm:
$$\|f\|_{\pi}^2 = \sum_{x \in X} |f(x)|^2 \pi(x)$$
and hence an inner product $\langle f, g \rangle_{\pi}$. View transition kernel as linear operator $Pf(x) = \sum_y p_{x,y} f(y)$: by reversibility this is $\langle \cdot, \cdot \rangle_{\pi}$ symmetric.

Finite-state-space reversible Markov chains and (weighted) euclidean spaces.

- $\langle f, g \rangle_{\pi} = \sum_y f(y)g(y)\pi(y)$.
- Test understanding:** use detailed balance to show $\langle f, Pg \rangle_{\pi} = \sum_x f(x) \sum_y p_{x,y} g(y)\pi(x) = \langle Pf, g \rangle_{\pi}$
- Adam Willis (MMORSE student at Warwick, 2004-2008) recently wrote an excellent Integrated Masters project on this subject.
- The vector space of functions on a finite state space is finite-dimensional!

- └ 8: Cutoff
- └ Cutoff and eigenvalues

Cutoff (II): eigenvalues and eigenfunctions

So P can be viewed as a symmetric matrix and thus has a full set of eigenvalues $-1 \leq \lambda_k \leq \dots \leq \lambda_1 \leq 1$ (if X has k elements) and corresponding normalized eigenfunctions V_1, \dots, V_k . Because of symmetry of P we may take the V_i to be an orthonormal basis, so

$$\sum |f(y)|^2 \pi(y) = \sum_{i=1}^k \langle f, V_i \rangle_{\pi}^2$$

The law of total probability implies $\lambda_1 = 1$ and $V_1 \equiv 1$, and irreducibility implies $\lambda_2 < \lambda_1$. Aperiodicity implies $-1 < \lambda_k$.

- └ 8: Cutoff
- └ Cutoff and eigenvalues
- └ Cutoff (II): eigenvalues and eigenfunctions

Cutoff (II): eigenvalues and eigenfunctions
So P can be viewed as a symmetric matrix and thus has a full set of eigenvalues $-1 \leq \lambda_2 \leq \dots \leq \lambda_1 \leq 1$ if X has k elements and corresponding normalized eigenfunctions V_1, \dots, V_k . Because of symmetry of P we may take the V_i to be an orthonormal basis, so
$$\sum |f(y)|^2 \pi(y) = \sum_{i=1}^k \langle f, V_i \rangle_{\pi}^2$$
The law of total probability implies $\lambda_1 = 1$ and $V_1 \equiv 1$, and irreducibility implies $\lambda_2 < \lambda_1$. Aperiodicity implies $-1 < \lambda_k$.

- Normalized: $\|V_i\|_{\pi}^2 = 1$; eigen property: $PV_i = \lambda_i V_i$.
- In fact all eigenvalues cannot exceed 1 in absolute value, by an inequality argument. Two eigenvalues equal to 1 would allow us to split state space into 2 components which violated irreducibility.
- In passing, there is a useful analysis of rate of convergence of **expectations** of functions of Markov chains based on this spectral analysis. Good when you know *a priori* what you want to estimate ...

- └ 8: Cutoff
- └ Two metrics

Cutoff (III): metrics

We need to relate total variation distance to the weighted Euclidean distance. Recall

$$\text{dist}_{TV}(P_X^{(n)}, \pi) = \frac{1}{2} \sum_y |P_X^{(n)}(y) - \pi(y)| = \frac{1}{2} \sum_y \left| \frac{P_X^{(n)}(y)}{\pi(y)} - 1 \right| \pi(y)$$

But this relates to weighted Euclidean distance by using Cauchy-Schwartz inequality and $\sum_y \pi(y) = 1$:

$$2 \text{dist}_{TV}(P_X^{(n)}, \pi) \leq \sqrt{\left\| \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2} \sqrt{\sum_y \pi(y)} = \sqrt{\left\| \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2}$$

Now expand using orthonormal eigenfunctions and $V_1 \equiv 1$:

$$\left\| \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2 = \sum_{i=2}^k \left\langle \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)}, V_i \right\rangle_{\pi}^2 = \sum_{i=2}^k (P_X^n V_i)^2 = \sum_{i=2}^k \lambda_i^{2n} V_i(x)^2$$

- └ 8: Cutoff
- └ Two metrics
- └ Cutoff (III): metrics

Cutoff (III): metrics
We need to relate total variation distance to the weighted Euclidean distance. Recall
$$\text{dist}_{TV}(P_X^{(n)}, \pi) = \frac{1}{2} \sum_y |P_X^{(n)}(y) - \pi(y)| = \frac{1}{2} \sum_y \left| \frac{P_X^{(n)}(y)}{\pi(y)} - 1 \right| \pi(y)$$
But this relates to weighted Euclidean distance by using Cauchy-Schwartz inequality and $\sum_y \pi(y) = 1$:
$$2 \text{dist}_{TV}(P_X^{(n)}, \pi) \leq \sqrt{\left\| \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2} \sqrt{\sum_y \pi(y)} = \sqrt{\left\| \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2}$$
Now expand using orthonormal eigenfunctions and $V_1 \equiv 1$:
$$\left\| \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2 = \sum_{i=2}^k \left\langle \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)}, V_i \right\rangle_{\pi}^2 = \sum_{i=2}^k (P_X^n V_i)^2 = \sum_{i=2}^k \lambda_i^{2n} V_i(x)^2$$

- The key here is the Cauchy-Schwartz inequality, $(\mathbb{E}[XY])^2 \leq \mathbb{E}[X^2] \mathbb{E}[Y^2]$.

Applied probabilists and statisticians may be more comfortable with this if they recognize that it is proved in the same way as the statement that correlations are always bounded between ± 1 .

- Miss $i = 1$ since $V_1 \equiv 1$, so

$$\left\langle \frac{P_X^{(n)}(\cdot)}{\pi(\cdot)} - 1, V_1 \right\rangle_{\pi} = \sum_y P_X^{(n)}(y) - \langle V_1, V_1 \rangle_{\pi} = 1 - 1 = 0$$

Miss -1 in other terms by orthogonality, since for $i > 1$

$$\langle -1, V_i \rangle_{\pi} = -\langle V_1, V_i \rangle_{\pi} = 0$$

- Bear in mind that in this finite-state-space context eigenfunctions are the same as eigenvectors!

Cutoff (IV): upper bound in special case

Gibbs' sampler for zero-interaction Ising model

Model for Gibb's sampler. Consider $N \times N$ array of ± 1 . At each step choose entry at random, flip sign.

As above, identify $\binom{N^2}{r}$ eigenfunctions of eigenvalue $1 - \frac{2r}{N^2}$, for $0 \leq r \leq N^2$. Set $n = \frac{N^2}{4}(\log(N^2) + \theta)$.

$$\begin{aligned} \left\| \frac{P_x^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2 &= \sum_{r=1}^{N^2} \binom{N^2}{r} \left(1 - \frac{2r}{N^2}\right)^{2n} \\ &\leq \sum_{r=1}^{N^2} \binom{N^2}{r} \exp\left(-\frac{2r}{N^2} \left(\frac{N^2}{2}(\log(N^2) + \theta)\right)\right) \\ &= \sum_{r=1}^{N^2} \binom{N^2}{r} (N^2)^{-r} e^{-r\theta} \leq \sum_{r=1}^{N^2} \frac{1}{r!} e^{-r\theta} \leq \exp(e^{-\theta}) - 1 \end{aligned}$$

2008-07-12 APTS-ASP 177
 ↳ 8: Cutoff
 ↳ A special case
 ↳ Cutoff (IV): upper bound in special case

Cutoff (IV): upper bound in special case
 Model for Gibb's sampler. Consider $N \times N$ array of ± 1 . At each step choose entry at random, flip sign.
 As above, identify $\binom{N^2}{r}$ eigenfunctions of eigenvalue $1 - \frac{2r}{N^2}$, for $0 \leq r \leq N^2$. Set $n = \frac{N^2}{4}(\log(N^2) + \theta)$.

$$\left\| \frac{P_x^{(n)}(\cdot)}{\pi(\cdot)} - 1 \right\|_{\pi}^2 = \sum_{r=1}^{N^2} \binom{N^2}{r} \left(1 - \frac{2r}{N^2}\right)^{2n}$$

$$\leq \sum_{r=1}^{N^2} \binom{N^2}{r} \exp\left(-\frac{2r}{N^2} \left(\frac{N^2}{2}(\log(N^2) + \theta)\right)\right)$$

$$= \sum_{r=1}^{N^2} \binom{N^2}{r} (N^2)^{-r} e^{-r\theta} \leq \sum_{r=1}^{N^2} \frac{1}{r!} e^{-r\theta} \leq \exp(e^{-\theta}) - 1$$

1. Eigenfunctions are just products $X_{i_1} \dots X_{i_k}$ of spin variables $X_r = \pm 1$.
2. Note, $1 - x \leq e^{-x}$ always.
3. Do the sums. In particular, note $PX_1 = \frac{1}{N^2}(-X_1) + (1 - \frac{1}{N^2})X_1 = (1 - \frac{2}{N^2})X_1, \dots$

Cutoff (V): lower bound in special case

The upper bound suggests a cutoff:

$$\text{dist}_{TV}(P_x^{(n)}, \pi) \leq \sqrt{\frac{\exp(e^{-\theta}) - 1}{2}}$$

Since $n = \frac{N^2}{4}(\log(N^2) + \theta)$, the cutoff occurs at around $\frac{N^2}{4} \log(N^2)$ and lasts of order $\frac{N^2}{4}$.

However to make **sure** this works, we also need a lower bound on $\text{dist}_{TV}(P_x^{(n)}, \pi)$. Achieve this by comparing means and variances of $Z \sum_{i=1}^{N^2} X_i$, where X_i is spin at site i . Simple estimates confirm that there is still substantial total variation distance at $\frac{N^2}{4} \log(N^2)$, so this is a real cutoff.

Moral: **effective** convergence can be much faster than one realizes, and occur over a fairly well defined period of time.

2008-07-12 APTS-ASP 179
 ↳ 8: Cutoff
 ↳ A special case
 ↳ Cutoff (V): lower bound in special case

Cutoff (V): lower bound in special case
 The upper bound suggests a cutoff:
 $\text{dist}_{TV}(P_x^{(n)}, \pi) \leq \sqrt{\frac{\exp(e^{-\theta}) - 1}{2}}$

Since $n = \frac{N^2}{4}(\log(N^2) + \theta)$, the cutoff occurs at around $\frac{N^2}{4} \log(N^2)$ and lasts of order $\frac{N^2}{4}$.
 However to make sure this works, we also need a lower bound on $\text{dist}_{TV}(P_x^{(n)}, \pi)$. Achieve this by comparing means and variances of $Z \sum_{i=1}^{N^2} X_i$, where X_i is spin at site i . Simple estimates confirm that there is still substantial total variation distance at $\frac{N^2}{4} \log(N^2)$, so this is a real cutoff.
 Moral: effective convergence can be much faster than one realizes, and occur over a fairly well defined period of time.

Calculations for other cases can be much harder.

In general, expect cutoff when there are large numbers of "second" eigenvalues. Should one expect cutoff for the case of an Ising model with weak interaction? Probably

The famous *Peres conjecture* says cutoff is to be expected for a chain with transitive symmetry if $(1 - \lambda_2)\tau \rightarrow \infty$, where λ_2 is the second largest eigenvalue (so $1 - \lambda_2$ is the "spectral gap"), and τ is the (deterministic) time at which the total variation distance to equilibrium becomes smaller than $\frac{1}{2}$. However there is a counterexample to Peres' conjecture as expressed above, (communication from Connor, PhD thesis 2007, which is communication of Diaconis, of work of which Diaconis knows ...). So the conjecture needs to be refined!

Use Markov's inequality to convert mean and variance comparisons into inequalities.

(Further practical remarks ...)

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.

Athreya, K. B. and P. Ney (1978). A new approach to the limit theory of recurrent Markov chains. *Trans. Amer. Math. Soc.* 245, 493-501.

Bayer, D. and P. Diaconis (1992). Trailing the dovetail shuffle to its lair. *Ann. Appl. Probab.* 2(2), 294-313.

Breiman, L. (1992). *Probability*, Volume 7 of *Classics in Applied Mathematics*. Philadelphia, PA: Society for Industrial and Applied Mathematics (SIAM). Corrected reprint of the 1968 original.

Brown, B. M. (1971). Martingale central limit theorems. *Ann. Math. Statist.* 42, 59-66.

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.

Fleming, I. (1953). *Casino Royale*. Jonathan Cape.

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.

Jerrum, M. (2003). *Counting, sampling and integrating: algorithms and complexity*. Lectures in Mathematics ETH Zürich. Basel: Birkhäuser Verlag.

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

Kendall, W. S. (2004). Geometric ergodicity and perfect simulation. *Electronic Communications in Probability* 9, 140-151.

Kendall, W. S., F. Liang, and J.-S. Wang (Eds.) (2005). *Markov chain Monte Carlo: Innovations and Applications*. Number 7 in IMS Lecture Notes. Singapore: World Scientific.

Kendall, W. S. and G. Montana (2002). Small sets and Markov transition densities. *Stochastic Process. Appl.* 99(2), 177-194.

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.

Meyn, S. P. and R. L. Tweedie (1993). *Markov chains and stochastic stability*. Communications and Control Engineering Series. London: Springer-Verlag London Ltd.

Murdoch, D. J. and P. J. Green (1998). Exact sampling from a continuous state space. *Scand. J. Statist.* 25(3), 483-502.

Norris, J. R. (1998). *Markov chains*, Volume 2 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge: Cambridge University Press. Reprint of 1997 original.

Nummelin, E. (1978).

A splitting technique for Harris recurrent Markov chains.
Z. Wahrsch. Verw. Gebiete 43(4), 309-318.

Nummelin, E. (1984).

General irreducible Markov chains and nonnegative operators, Volume 83 of
Cambridge Tracts in Mathematics.
Cambridge: Cambridge University Press.

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 (1995).
Stochastic geometry and its applications (Second ed.).
Chichester: John Wiley & Sons.
(First edition in 1987 joint with Akademie Verlag, Berlin).

Williams, D. (1991).

Probability with martingales.
Cambridge Mathematical Textbooks. Cambridge: Cambridge University Press.

Photographs used in text

- ▶ Police phone box en.wikipedia.org/wiki/Image:Earls_Court_Police_Box.jpg
- ▶ The standing martingale
en.wikipedia.org/wiki/Image:Hunterhorse.jpg
- ▶ Boat Race: en.wikipedia.org/wiki/Image:Boat_Race_Finish_2008_-_Oxford_winners.jpg
- ▶ Impact site of fragment G of Comet Shoemaker-Levy 9 on Jupiter
en.wikipedia.org/wiki/Image:Impact_site_of_fragment_G.gif
- ▶ The cardplayers en.wikipedia.org/wiki/Image:Paul_C%C3%A9zanne%2C_Les_joueurs_de_carte_%281892-95%29.jpg
- ▶ Chinese abacus en.wikipedia.org/wiki/Image:Boulier1.JPG
- ▶ Error function
en.wikipedia.org/wiki/Image/Error_Function.svg
- ▶ Boomerang en.wikipedia.org/wiki/Image:Boomerang.jpg
- ▶ Alexander Lyapunov en.wikipedia.org/wiki/Image:Alexander_Ljapunow_jung.jpg
- ▶ Riffle shuffle (photo by Johnny Blood)
en.wikipedia.org/wiki/Image:Riffle_shuffle.jpg