

Chapman & Hall/CRC
Mathematical and Computational Biology Series

Stochastic Modelling for Systems Biology

Published Titles

Differential Equations and Mathematical Biology, *D.S. Jones and B.D. Sleeman*

Modeling and Simulation of Capsules and Biological Cells, *C. Pozrikidis*

Cancer Modelling and Simulation, *Luigi Preziosi*

Data Analysis Tools for DNA Microarrays, *Sorin Draghici*

Computational Neuroscience: A Comprehensive Approach, *Jianfeng Feng*

The Ten Most Wanted Solutions in Protein Bioinformatics, *Anna Tramontano*

Exactly Solvable Models of Biological Invasion, *Sergei V. Petrovskii and Lian-Bai Li*

Knowledge Discovery in Proteomics, *Igor Jurisica and Dennis Wigle*

Normal Mode Analysis: Theory and Applications to Biological and Chemical Systems, *Qiang Cui and Ivet Bahar*

An Introduction to Systems Biology: Design Principles of Biological Circuits, *Uri Alon*

Stochastic Modelling for Systems Biology, *James Darren Wilkinson*

Forthcoming Titles

Computational Biology: A Statistical Mechanics Perspective, *Ralf Blossey*

Introduction to Bioinformatics, *Anna Tramontano*

Biological Sequence Analysis with Iterative Maps, *Jonas S. Almeida*

Practical Guide to Protein Bioinformatics: Sequence, Structure and Function, *Shoba Ranganathan*

CHAPMAN & HALL/CRC

Mathematical and Computational Biology Series

Aims and scope:

This series aims to capture new developments and summarize what is known over the whole spectrum of mathematical and computational biology and medicine. It seeks to encourage the integration of mathematical, statistical and computational methods into biology by publishing a broad range of textbooks, reference works and handbooks. The titles included in the series are meant to appeal to students, researchers and professionals in the mathematical, statistical and computational sciences, fundamental biology and bioengineering, as well as interdisciplinary researchers involved in the field. The inclusion of concrete examples and applications, and programming techniques and examples, is highly encouraged.

Series Editors

Alison M. Etheridge
Department of Statistics
University of Oxford

Louis J. Gross
Department of Ecology and Evolutionary Biology
University of Tennessee

Suzanne Lenhart
Department of Mathematics
University of Tennessee

Philip K. Maini
Mathematical Institute
University of Oxford

Shoba Ranganathan
Research Institute of Biotechnology
Macquarie University

Hershel M. Safer
Weizmann Institute of Science
Bioinformatics & Bio Computing

Eberhard O. Voit
The Wallace H. Coulter Department of Biomedical Engineering
Georgia Tech and Emory University

Proposals for the series should be submitted to one of the series editors above or directly to:

CRC Press, Taylor & Francis Group

24-25 Blades Court
Deodar Road
London SW15 2NU
UK

TAMPEREEN TEKNILLINEN YLIOPISTO
Signaalinkäsittelyn laitos

SGN / MARJA-LEENA LINNE

254 s

Published in 2006 by
CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2006 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group

No claim to original U.S. Government works
Printed in the United States of America on acid-free paper
10 9 8 7 6 5 4 3 2

International Standard Book Number-10: 1-58488-540-8 (Hardcover)
International Standard Book Number-13: 978-1-58488-540-5 (Hardcover)

This book contains information obtained from authentic and highly regarded sources. Reprinted material is quoted with permission, and sources are indicated. A wide variety of references are listed. Reasonable efforts have been made to publish reliable data and information, but the author and the publisher cannot assume responsibility for the validity of all materials or for the consequences of their use.

No part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC) 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Library of Congress Cataloging-in-Publication Data

Catalog record is available from the Library of Congress

informa

Taylor & Francis Group
is the Academic Division of Informa plc.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

APPLIED AND COMPUTATIONAL MATHEMATICS
A JOHN WILEY & SONS, INC., PUBLICATION

Chapman & Hall/CRC
Mathematical and Computational Biology Series

Stochastic Modelling for Systems Biology

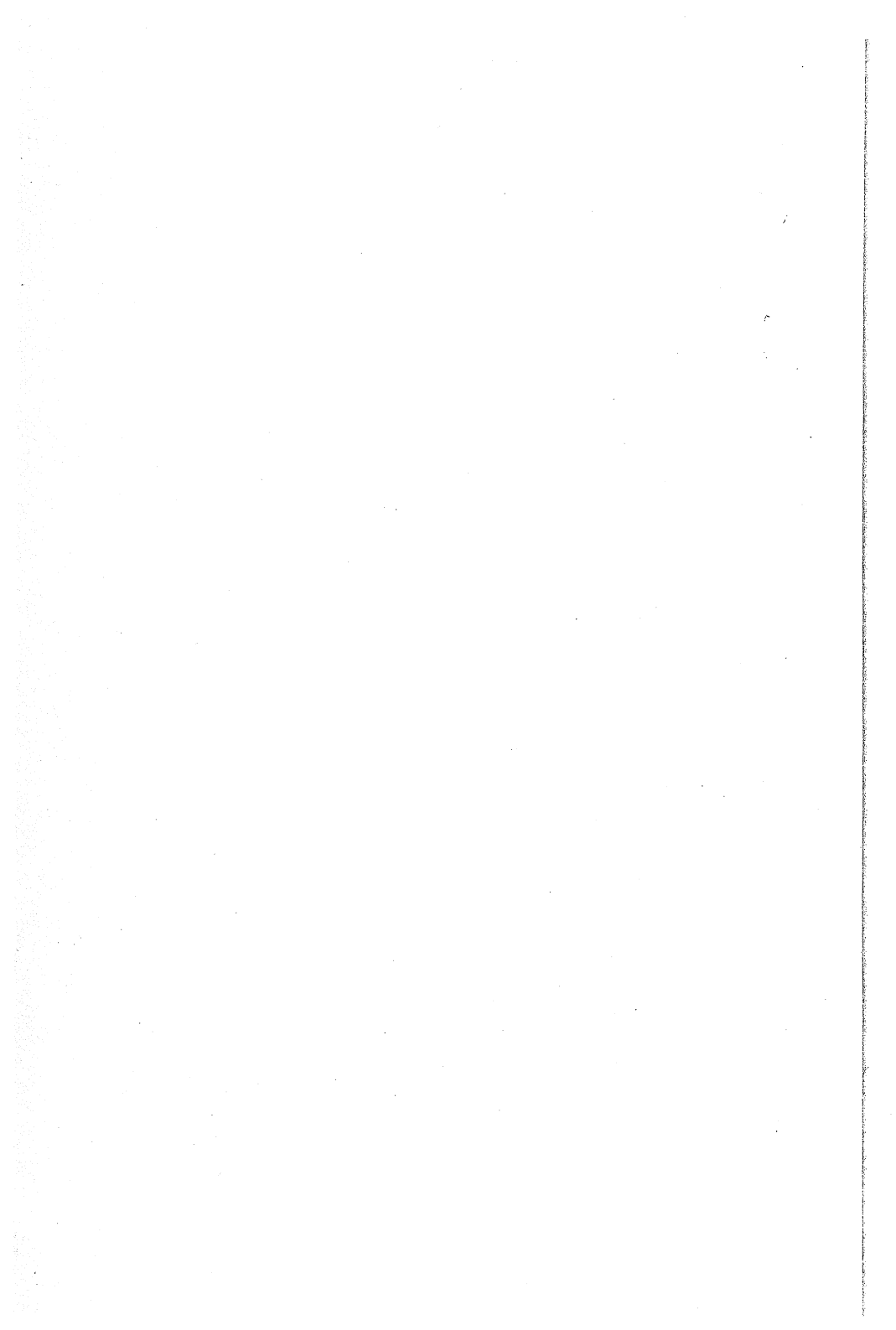
Darren James Wilkinson

Contents

1	Introduction to biological modelling	1
1.1	What is modelling?	1
1.2	Aims of modelling	2
1.3	Why is stochastic modelling necessary?	2
1.4	Chemical reactions	6
1.5	Modelling genetic and biochemical networks	8
1.6	Modelling higher-level systems	16
1.7	Exercises	17
1.8	Further reading	17
2	Representation of biochemical networks	19
2.1	Coupled chemical reactions	19
2.2	Graphical representations	19
2.3	Petri nets	21
2.4	Systems Biology Markup Language (SBML)	31
2.5	SBML-shorthand	36
2.6	Exercises	42
2.7	Further reading	43
3	Probability models	45
3.1	Probability	45
3.2	Discrete probability models	56
3.3	The discrete uniform distribution	64
3.4	The binomial distribution	64
3.5	The geometric distribution	65
3.6	The Poisson distribution	67

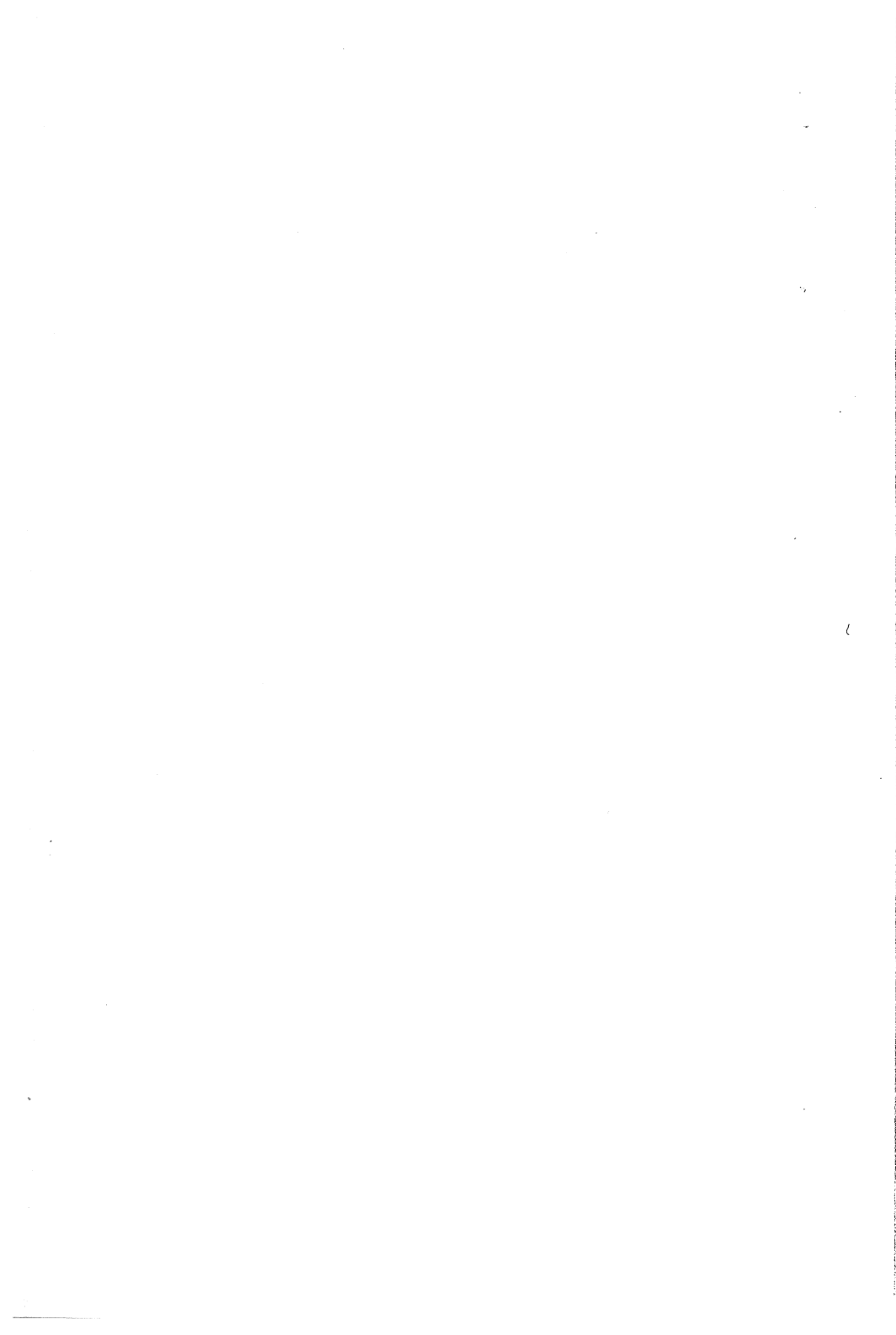
3.7	Continuous probability models	70
3.8	The uniform distribution	75
3.9	The exponential distribution	77
3.10	The normal/Gaussian distribution	82
3.11	The gamma distribution	86
3.12	Exercises	88
3.13	Further reading	89
4	Stochastic simulation	91
4.1	Introduction	91
4.2	Monte-Carlo integration	91
4.3	Uniform random number generation	92
4.4	Transformation methods	93
4.5	Lookup methods	97
4.6	Rejection samplers	98
4.7	The Poisson process	101
4.8	Using the statistical programming language, R	101
4.9	Analysis of simulation output	106
4.10	Exercises	107
4.11	Further reading	108
5	Markov processes	109
5.1	Introduction	109
5.2	Finite discrete time Markov chains	109
5.3	Markov chains with continuous state space	115
5.4	Markov chains in continuous time	121
5.5	Diffusion processes	133
5.6	Exercises	135
5.7	Further reading	137
6	Chemical and biochemical kinetics	139
6.1	Classical continuous deterministic chemical kinetics	139
6.2	Molecular approach to kinetics	145
6.3	Mass-action stochastic kinetics	147
6.4	The Gillespie algorithm	149
6.5	Stochastic Petri nets (SPNs)	150
6.6	Rate constant conversion	153
6.7	The Master equation	157
6.8	Software for simulating stochastic kinetic networks	160
6.9	Exercises	161
6.10	Further reading	161
7	Case studies	163
7.1	Introduction	163
7.2	Dimerisation kinetics	163

7.3	Michaelis-Menten enzyme kinetics	168
7.4	An auto-regulatory genetic network	172
7.5	The <i>lac</i> operon	176
7.6	Exercises	178
7.7	Further reading	179
8	Beyond the Gillespie algorithm	181
8.1	Introduction	181
8.2	Exact simulation methods	181
8.3	Approximate simulation strategies	186
8.4	Hybrid simulation strategies	190
8.5	Exercises	196
8.6	Further reading	196
9	Bayesian inference and MCMC	197
9.1	Likelihood and Bayesian inference	197
9.2	The Gibbs sampler	202
9.3	The Metropolis-Hastings algorithm	212
9.4	Hybrid MCMC schemes	216
9.5	Exercises	216
9.6	Further reading	217
10	Inference for stochastic kinetic models	219
10.1	Introduction	219
10.2	Inference given complete data	220
10.3	Discrete-time observations of the system state	223
10.4	Diffusion approximations for inference	230
10.5	Network inference	234
10.6	Exercises	234
10.7	Further reading	235
11	Conclusions	237
A	SBML Models	239
A.1	Auto-regulatory network	239
A.2	Lotka-Volterra reaction system	242
A.3	Dimerisation-kinetics model	243
	References	247
	Index	251



List of tables

- | | | |
|-----|---|----|
| 2.1 | The auto-regulatory system displayed in tabular (matrix) form (zero stoichiometries omitted for clarity) | 24 |
| 2.2 | Table representing the overall effect of each transition (reaction) on the marking (state) of the network | 25 |



List of figures

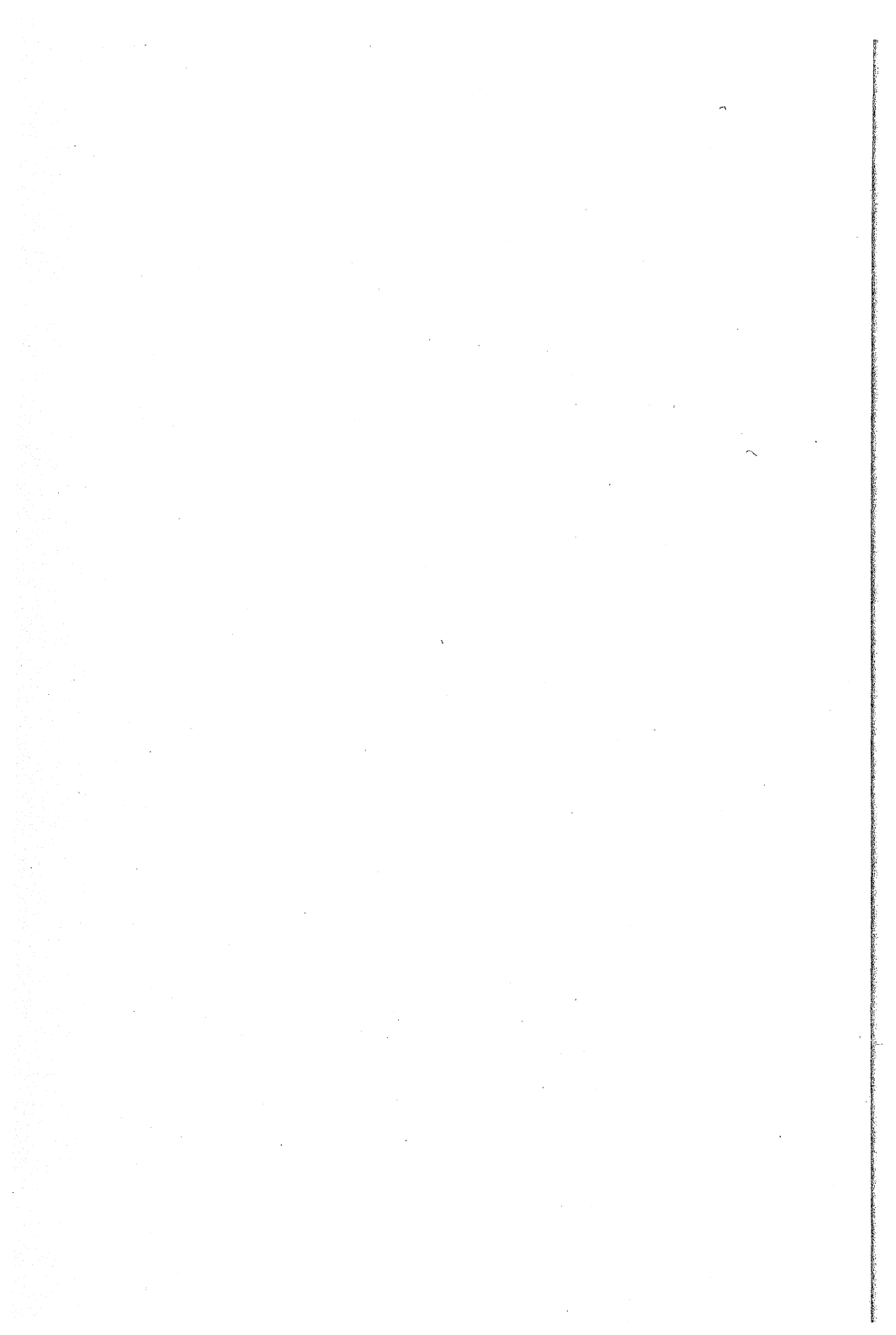
1.1	Five deterministic solutions of the linear birth-death process for values of $\lambda - \mu$ given in the legend ($x_0 = 50$)	4
1.2	Five realisations of a stochastic linear birth-death process together with the continuous deterministic solution ($x_0 = 50, \lambda = 3, \mu = 4$)	5
1.3	Five realisations of a stochastic linear birth-death process together with the continuous deterministic solution for four different (λ, μ) combinations, each with $\lambda - \mu = -1$ and $x_0 = 50$	6
1.4	Transcription of a single prokaryotic gene	9
1.5	A simple illustrative model of the transcription process in eukaryotic cells	11
1.6	A simple prokaryotic transcription repression mechanism	12
1.7	A very simple model of a prokaryotic auto-regulatory gene network. Here dimers of a protein P coded for by a gene g repress their own transcription by binding to a regulatory region q upstream of g and downstream of the promoter p .	14
1.8	Key mechanisms involving the <i>lac</i> operon. Here an inhibitor protein I can repress transcription of the <i>lac</i> operon by binding to the operator o . However, in the presence of lactose, the inhibitor preferentially binds to it, and in the bound state can no longer bind to the operator, thereby allowing transcription to proceed.	15
2.1	A simple graph of the auto-regulatory reaction network	20
2.2	A simple digraph	21
2.3	A Petri net for the auto-regulatory reaction network	22
2.4	A Petri net labelled with tokens	23
2.5	A Petri net with new numbers of tokens after reactions have taken place	23
3.1	CDF for the sum of a pair of fair dice	58
3.2	PMF and CDF for a $B(8, 0.7)$ distribution	65
3.3	PMF and CDF for a $Po(5)$ distribution	69
3.4	PDF and CDF for a $U(0, 1)$ distribution	76
3.5	PDF and CDF for an $Exp(1)$ distribution	78
3.6	PDF and CDF for a $N(0, 1)$ distribution	83
3.7	Graph of $\Gamma(x)$ for small positive values of x	87
3.8	PDF and CDF for a $\Gamma(3, 1)$ distribution	88

4.1	Density of $Y = \exp(X)$, where $X \sim N(2, 1)$.	108
5.1	An R function to simulate a sample path of length n from a Markov chain with transition matrix P and initial distribution π_0	115
5.2	A sample R session to simulate and analyse the sample path of a finite Markov chain. The last two commands show how to use R to directly compute the stationary distribution of a finite Markov chain.	116
5.3	SBML-shorthand for the simple gene activation process with $\alpha = 0.5$ and $\beta = 1$	124
5.4	A simulated realisation of the simple gene activation process with $\alpha = 0.5$ and $\beta = 1$	127
5.5	An R function to simulate a sample path with n events from a continuous time Markov chain with transition rate matrix Q and initial distribution π_0	128
5.6	SBML-shorthand for the immigration-death process with $\lambda = 1$ and $\mu = 0.1$	128
5.7	A single realisation of the immigration-death process with parameters $\lambda = 1$ and $\mu = 0.1$, initialised at $X(0) = 0$. Note that the stationary distribution of this process is Poisson with mean 10.	130
5.8	R function for discrete-event simulation of the immigration-death process	131
5.9	R function for simulation of a diffusion process using the Euler method	134
5.10	A single realisation of the diffusion approximation to the immigration-death process with parameters $\lambda = 1$ and $\mu = 0.1$, initialised at $X(0) = 0$	135
5.11	R code for simulating the diffusion approximation to the immigration-death process	136
6.1	Lotka-Volterra dynamics for $[Y_1](0) = 4$, $[Y_2](0) = 10$, $k_1 = 1$, $k_2 = 0.1$, $k_3 = 0.1$. Note that the equilibrium solution for this combination of rate parameters is $[Y_1] = 1$, $[Y_2] = 10$.	141
6.2	Lotka-Volterra dynamics in phase-space for rate parameters $k_1 = 1$, $k_2 = 0.1$, $k_3 = 0.1$. The dynamics for the initial condition $[Y_1](0) = 4$, $[Y_2](0) = 10$ are shown as the bold orbit. Note that the system moves around this orbit in an anti-clockwise direction. Orbits for other initial conditions are shown as dotted curves. Note that the equilibrium solution for this combination of rate parameters is $[Y_1] = 1$, $[Y_2] = 10$.	142
6.3	Dimerisation kinetics for $[P](0) = 1$, $[P_2](0) = 0$, $k_1 = 1$, $k_2 = 0.5$. This combination of parameters gives $K_{eq} = 2$, $c = 1$, and hence equilibrium concentrations of $[P] = 0.39$, $[P_2] = 0.30$.	143
6.4	An R function to numerically integrate a system of coupled ODEs using a simple first-order Euler method	145

6.5	An R function to implement the Gillespie algorithm for a stochastic Petri net representation of a coupled chemical reaction system	150
6.6	Some R code to set up the LV system as a SPN and then simulate it using the Gillespie algorithm. The state of the system is initialised to 50 prey and 100 predators, and the stochastic rate constants are $c = (1, 0.005, 0.6)'$.	152
6.7	A single realisation of a stochastic LV process. The state of the system is initialised to 50 prey and 100 predators, and the stochastic rate constants are $c = (1, 0.005, 0.6)'$.	152
6.8	A single realisation of a stochastic LV process in phase-space. The state of the system is initialised to 50 prey and 100 predators, and the stochastic rate constants are $c = (1, 0.005, 0.6)'$.	153
6.9	SBML-shorthand for the stochastic Lotka-Volterra system	154
6.10	An R function to discretise the output of <code>gillespie</code> onto a regular grid of time points. The result is returned as an R multivariate time series object.	154
6.11	An R function to implement the Gillespie algorithm for a SPN, recording the state on a regular grid of time points. The result is returned as an R multivariate time series object.	155
7.1	SBML-shorthand for the dimerisation kinetics model (continuous deterministic version)	164
7.2	Left: Simulated continuous deterministic dynamics of the dimerisation kinetics model. Right: A simulated realisation of the discrete stochastic dynamics of the dimerisation kinetics model.	164
7.3	SBML-shorthand for the dimerisation kinetics model (discrete stochastic version)	165
7.4	R code to build an SPN object representing the dimerisation kinetics model	166
7.5	Left: A simulated realisation of the discrete stochastic dynamics of the dimerisation kinetics model plotted on a concentration scale. Right: The trajectories for levels of P from 20 runs overlaid.	167
7.6	Left: The mean trajectory of P together with some approximate (point-wise) “confidence bounds” based on 1,000 runs of the simulator. Right: Density histogram of the simulated realisations of P at time $t = 10$ based on 10,000 runs, giving an estimate of the PMF for $P(10)$.	167
7.7	SBML-shorthand for the Michaelis-Menten kinetics model (continuous deterministic version)	169
7.8	Left: Simulated continuous deterministic dynamics of the Michaelis-Menten kinetics model. Right: Simulated continuous deterministic dynamics of the Michaelis-Menten kinetics model based on the two-dimensional representation.	170
7.9	SBML-shorthand for the Michaelis-Menten kinetics model (discrete stochastic version)	171

- 7.10 Left: A simulated realisation of the discrete stochastic dynamics of the Michaelis-Menten kinetics model. Right: A simulated realisation of the discrete stochastic dynamics of the reduced-dimension Michaelis-Menten kinetics model. 172
- 7.11 SBML-shorthand for the reduced dimension Michaelis-Menten kinetics model (discrete stochastic version) 172
- 7.12 Left: A simulated realisation of the discrete stochastic dynamics of the prokaryotic genetic auto-regulatory network model, for a period of 5,000 seconds. Right: A close-up on the first period of 250 seconds of the left plot. 174
- 7.13 Left: Close-up showing the time-evolution of the number of molecules of P over a 10-second period. Right: Empirical PMF for the number of molecules of P at time $t = 10$ seconds, based on 10,000 runs. 174
- 7.14 Left: Empirical PMF for the number of molecules of P at time $t = 10$ seconds when k_2 is changed from 0.01 to 0.02, based on 10,000 runs. Right: Empirical PMF for the prior predictive uncertainty regarding the observed value of P at time $t = 10$ based on the prior distribution $k_2 \sim U(0.005, 0.03)$. 176
- 7.15 SBML-shorthand for the *lac*-operon model (discrete stochastic version) 177
- 7.16 A simulated realisation of the discrete stochastic dynamics of the *lac*-operon model for a period of 50,000 seconds. An intervention is applied at time $t = 20,000$, when 10,000 molecules of lactose are added to the cell. 178
- 8.1 An R function to implement the first reaction method for a stochastic Petri net representation of a coupled chemical reaction system. It is to be used in the same way as the `gillespie` function from Figure 6.5. 183
- 8.2 An R function to implement the Poisson timestep method for a stochastic Petri net representation of a coupled chemical reaction system. It is to be used in the same way as the `gillespied` function from Figure 6.11. 187
- 8.3 An R function to integrate the CLE using an Euler method for a stochastic Petri net representation of a coupled chemical reaction system. It is to be used in the same way as the `pts` function from Figure 8.2. 190
- 9.1 Plot showing the prior and posterior for the Poisson rate example. Note how the prior is modified to give a posterior more consistent with the data (which has a sample mean of 3). 200
- 9.2 An R function to implement a Gibbs sampler for the simple normal random sample model. Example code for using this function is given in Figure 9.3. 207

- 9.3 Example R code illustrating the use of the function `normgibbs` from Figure 9.2. The plots generated by running this code are shown in Figure 9.4. In this example the prior took the form $\mu \sim N(10, 100)$, $\tau \sim \Gamma(3, 11)$, and the sufficient statistics for the data were $n = 15$, $\bar{x} = 25$, $s^2 = 20$. The sampler was run for 11,000 iterations with the first 1,000 discarded as burn-in, and the remaining 10,000 iterations used for the main monitoring run. 208
- 9.4 Figure showing the Gibbs sampler output resulting from running the example code in Figure 9.3. The top two plots give an indication of the bivariate posterior distribution. The second row shows trace plots of the marginal distributions of interest, indicating a rapidly mixing MCMC algorithm. The final row shows empirical marginal posterior distributions for the parameters of interest. 209
- 9.5 An R function to implement a Metropolis sampler for a standard normal random quantity based on $U(-\alpha, \alpha)$ innovations. So, `metrop(10000, 1)` will execute a run of length 10,000 with an α of 1. This α is close to optimal. Running with $\alpha = 0.1$ gives a chain that is too cold, and $\alpha = 100$ gives a chain that is too hot. 215



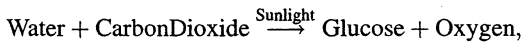
Introduction to biological modelling

1.1 What is modelling?

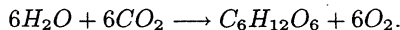
Modelling is an attempt to describe, in a precise way, an understanding of the elements of a system of interest, their states, and their interactions with other elements. The model should be sufficiently detailed and precise so that it can in principle be used to simulate the behaviour of the system on a computer. In the context of molecular cell biology, a model may describe (some of) the mechanisms involved in transcription, translation, gene regulation, cellular signalling, DNA damage and repair processes, homeostatic processes, the cell cycle, or apoptosis. Indeed any biochemical mechanism of interest can, in principle, be modelled. At a higher level, modelling may be used to describe the functioning of a tissue, organ, or even an entire organism. At still higher levels, models can be used to describe the behaviour and time evolution of populations of individual organisms.

The first issue to confront when embarking on a modelling project is to decide on exactly which features to include in the model, and in particular, the *level* of detail the model is intended to capture. So, a model of an entire organism is unlikely to describe the detailed functioning of every individual cell, but a model of a cell is likely to include a variety of very detailed descriptions of key cellular processes. Even then, however, a model of a cell is unlikely to contain details of every single gene and protein.

Fortunately, biologists are used to thinking about processes at different scales and different levels of detail. Consider, for example, the process of photosynthesis. When studying photosynthesis for the first time at school, it is typically summarised by a single chemical reaction mixing water with carbon dioxide to get glucose and oxygen (catalysed by sunlight). This could be written very simply as



or more formally by replacing the molecules by their chemical formulas and balancing to get



Of course, further study reveals that photosynthesis consists of many reactions, and that the single reaction was simply a summary of the overall effect of the process. However, it is important to understand that the above equation is not really *wrong*, it just represents the overall process at a higher level than the more detailed description that biologists often prefer to work with. Whether a single overall equation or a full breakdown into component reactions is necessary depends on whether intermediaries such as ADP and ATP are elements of interest to the modeller. Indeed, really accurate

modelling of the process would require a model far more detailed and complex than most biologists would be comfortable with, using molecular dynamic simulations that explicitly manage the position and momentum of every molecule in the system.

The “art” of building a good model is to capture the essential features of the biology without burdening the model with non-essential details. Every model is to some extent a simplification of the biology, but models are valuable because they take ideas that might have been expressed verbally or diagrammatically and make them more explicit, so that they can begin to be understood in a *quantitative* rather than purely *qualitative* way.

1.2 Aims of modelling

The features of a model depend very much on the aims of the modelling exercise. We therefore need to consider why people model and what they hope to achieve by so doing. Often the most basic aim is to make clear the current state of knowledge regarding a particular system, by attempting to be precise about the elements involved and the interactions between them. Doing this can be a particularly effective way of highlighting gaps in understanding. In addition, having a detailed model of a system allows people to *test* that their understanding of a system is correct, by seeing if the implications of their models are consistent with observed experimental data. However, this work will often represent only the initial stage of the modelling process. Once people have a model they are happy with, they often want to use their models *predictively*, by conducting “virtual experiments” that might be difficult, time-consuming, or impossible to do in the lab. Such experiments may uncover important indirect relationships between model components that would be hard to predict otherwise. An additional goal of modern biological modelling is to pool a number of small models of well-understood mechanisms into a large model in order to investigate the effect of interactions between the model components. Models can also be extremely useful for informing the design and analysis of complex biological experiments.

In summary, modelling and computer simulation are becoming increasingly important in post-genomic biology for integrating knowledge and experimental data and making testable predictions about the behaviour of complex biological systems.

1.3 Why is stochastic modelling necessary?

Ignoring quantum mechanical effects, current scientific wisdom views biological systems as essentially deterministic in character, with dynamics entirely predictable given sufficient knowledge of the state of the system (together with complete knowledge of the physics and chemistry of interacting biomolecules). At first this perhaps suggests that a deterministic approach to the modelling of biological systems is likely to be successful. However, despite the rapid advancements in computing technology, we are still a very long way away from a situation where we might expect to be able to model biological systems of realistic size and complexity over interesting time scales using such a molecular dynamic approach. We must therefore use models that leave out many details of the “state” of a system (such as the position, orientation,

and momentum of every single molecule under consideration), in favour of a higher level view. Viewed at this higher level, the dynamics of the system are not deterministic, but intrinsically stochastic, and consideration of statistical physics is necessary to uncover the precise nature of the stochastic processes governing the system dynamics. A more detailed discussion of this issue will have to be deferred until much later in the book, once the appropriate concepts and terminology have been established. In the meantime, it is helpful to highlight the issues using a very simple example that illustrates the importance of stochastic modelling, both for simulation and inference.

The example we will consider is known as the *linear birth-death process*. In the first instance it is perhaps helpful to view this as a model for the number of bacteria in a bacterial colony. It is assumed that each bacterium in the colony gives rise to new individuals at rate λ (that is, on average, each bacterium will produce λ offspring per unit time). Similarly, each bacterium dies at rate μ (that is, on average, the proportion of bacteria that die per unit time is μ). These definitions are not quite right, but we will define such things much more precisely later. Let the number of bacteria in the colony at time t be denoted $X(t)$. Assume that the number of bacteria in the colony at time zero is known to be x_0 . Viewed in a continuous deterministic manner, this description of the system leads directly to the ordinary differential equation

$$\frac{dX(t)}{dt} = (\lambda - \mu)X(t),$$

which can be solved analytically to give the complete dynamics of the system as

$$X(t) = x_0 \exp\{(\lambda - \mu)t\}.$$

So, predictably, in the case $\lambda > \mu$ the population size will increase exponentially as $t \rightarrow \infty$, and will decrease in size exponentially if $\lambda < \mu$. Similarly, it will remain at constant size x_0 if $\lambda = \mu$. Five such solutions are given in Figure 1.1. There are other things worth noting about this solution. In particular, the solution clearly only depends on $\lambda - \mu$ and not on the particular values that λ and μ take (so, for example, $\lambda = 0.5, \mu = 0$ will lead to exactly the same solution as $\lambda = 1, \mu = 0.5$). In some sense, therefore, $\lambda - \mu$ (together with x_0) is a “sufficient” description of the system dynamics. At first this might sound like a good thing, but it is clear that there is a flipside: namely that studying *experimental data* on bacteria numbers can only provide information about $\lambda - \mu$, and not on the particular values of λ and μ separately (as the data can only provide information about the “shape” of the curve, and the shape of the curve is determined by $\lambda - \mu$). Of course, this is not a problem if the continuous deterministic model is really appropriate, as then $\lambda - \mu$ is the only thing one needs to know and the precise values of λ and μ are not important for predicting system behaviour. Note, however, that the lack of identifiability of λ and μ has implications for *network inference* as well as inference for rate constants. It is clear that in this model we cannot know from experimental data if we have a pure birth or death process, or a process involving both births and deaths, as it is not possible to know if λ or μ is zero.*

The problem of course, is that bacteria don't vary in number continuously and

* This also illustrates another point that is not widely appreciated — the fact that reliable network in-

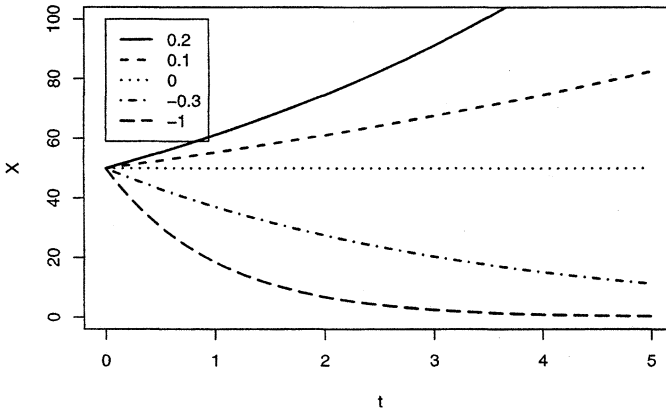


Figure 1.1 Five deterministic solutions of the linear birth-death process for values of $\lambda - \mu$ given in the legend ($x_0 = 50$)

deterministically. They vary discretely and stochastically. Using the techniques that will be developed later in this book, it is straightforward to understand the stochastic process associated with this model and to simulate it on a computer. By their very nature, such stochastic processes are random, and each time they are simulated they will look different. In order to understand the behaviour of such processes it is therefore necessary (in general) to study many realisations of the process. Five realisations are given in Figure 1.2, together with the corresponding deterministic solution.

It is immediately clear that the stochastic realisations exhibit much more interesting behaviour and match much better with the kind of experimental data one is likely to encounter. They also allow one to ask questions and get answers to issues that can't be addressed using a continuous deterministic model. For example, according to the deterministic model, the population size at time $t = 2$ is given by $X(2) = 50/e^2 \simeq 6.77$. Even leaving aside the fact that this is not an integer, we see from the stochastic realisations that there is considerable uncertainty for the value of $X(2)$, and stochastic simulation allows us to construct, *inter alia*, a likely range of values for $X(2)$. Another quantity of considerable practical interest is the "time to extinction" (the time, t , at which $X(t)$ first becomes zero). Under the deterministic model, $X(t)$ never reaches zero, but simply tends to zero as $t \rightarrow \infty$. We see from the stochastic realisations that these do go extinct, and that there is considerable randomness associated with the time that this occurs. Again, stochastic simulation will

ference is necessarily more difficult than rate-parameter inference, as determining the existence of a reaction is equivalent to deciding whether the rate of that reaction is zero.

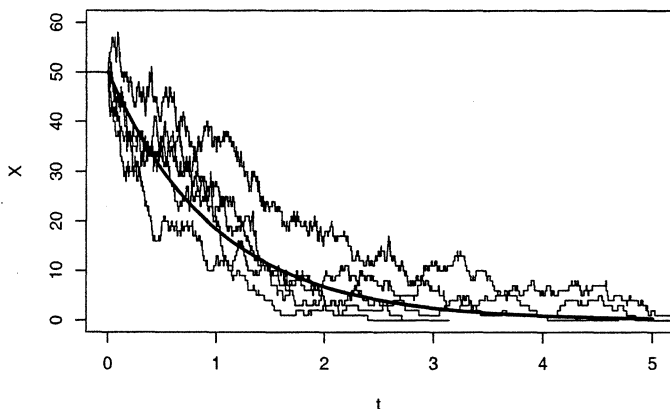


Figure 1.2 *Five realisations of a stochastic linear birth-death process together with the continuous deterministic solution ($x_0 = 50$, $\lambda = 3$, $\mu = 4$)*

allow us to understand the *distribution* associated with the time to extinction, something that simply isn't possible using a deterministic framework.

Another particularly noteworthy feature of the stochastic process representation is that it depends explicitly on both λ and μ , and not just on $\lambda - \mu$. This is illustrated in Figure 1.3. It is clear that although $\lambda - \mu$ controls the essential “shape” of the process, $\lambda + \mu$ controls the degree of “noise” or “volatility” in the system. This is a critically important point to understand — it tells us that if stochastic effects are present in the system, we cannot properly understand the system dynamics unless we know both λ and μ . Consequently, *we cannot simply fit a deterministic model to available experimental data and then use the inferred rate constants in a stochastic simulation*, as it is not possible to infer the stochastic rate constants using a deterministic model.

This has important implications for the use of stochastic models for inference from experimental data. It suggests that given some data on the variation in colony size over time, it ought to be possible to get information about $\lambda - \mu$ from the overall shape of the data, and information about $\lambda + \mu$ from the volatility of the data. If we know both $\lambda - \mu$ and $\lambda + \mu$, we can easily determine both λ and μ separately. Once we know both λ and μ , we can accurately simulate the dynamics of the system we are interested in (as well as inferring network structure, as we could also test to see if either λ or μ is zero). However, it is only possible to make satisfactory inferences for both λ and μ if the stochastic nature of the system is taken into account at the inference stage of the process.

Although we have here considered a trivial example, the implications are broad. In particular, they apply to the genetic and biochemical network models that much of this book will be concerned with. This is because genetic and biochemical networks

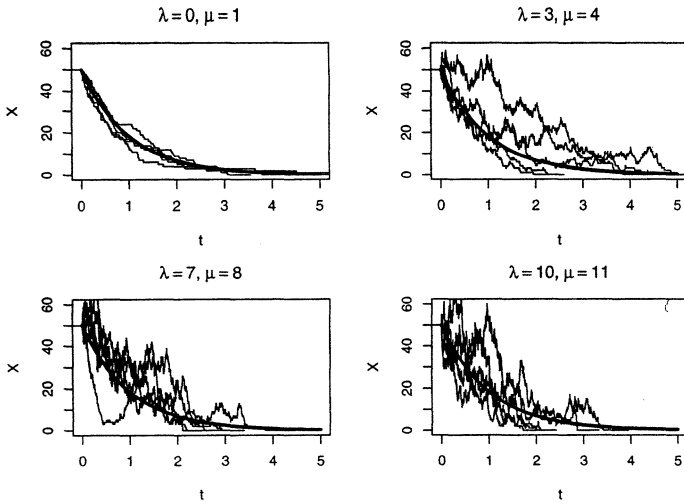


Figure 1.3 Five realisations of a stochastic linear birth-death process together with the continuous deterministic solution for four different (λ, μ) combinations, each with $\lambda - \mu = -1$ and $x_0 = 50$

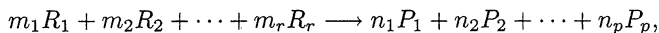
involve the interaction of integer numbers of molecules that react when they collide after random times, driven by Brownian motion. Although it is now becoming increasingly accepted that stochastic mechanisms are important in many (if not most) genetic and biochemical networks, routine use of stochastic simulation in order to understand system dynamics is still not widespread. This could be because inference methods regularly used in practice work by fitting continuous deterministic models to experimental data. We have just seen that such methods cannot in general give us reliable information about all of the parameters important for determining the stochastic dynamics of a system, and so stochastic simulation cannot be done reliably until we have good methods of inference for stochastic models. It turns out that it is possible to formalise the problem of inference for stochastic kinetic models from time-course experimental data, and this is the subject matter of Chapter 10. However, it should be pointed out at the outset that inference for stochastic models is at least an order of magnitude more difficult than inference for deterministic models (in terms of the mathematics required, algorithmic complexity, and computation time), and is still the subject of a great deal of on going research.

1.4 Chemical reactions

There are a number of ways one could represent a model of a biological system. Biologists have traditionally favoured diagrammatic schemas coupled with verbal explanations in order to convey qualitative information regarding mechanisms. At the other extreme, applied mathematicians traditionally prefer to work with systems of

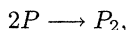
ordinary or partial differential equations (ODEs or PDEs). These have the advantage of being more precise and fully quantitative, but also have a number of disadvantages. In some sense differential equation models are too low level a description, as they not only encode the essential features of the model, but also a wealth of accompanying baggage associated with a particular interpretation of chemical kinetics that is not always well suited to application in the molecular biology context. Somewhere between these two extremes, the biochemist will tend to view systems as networks of coupled chemical reactions, and it appears that most of the best ways of representing biochemical mechanisms exist at this level of detail, though there are many ways of representing networks of this type. Networks of coupled chemical reactions are sufficiently general that they can be simulated in different ways using different algorithms depending on assumptions made about the underlying kinetics. On the other hand, they are sufficiently detailed and precise so that once the kinetics have been specified, they can be used directly to construct full dynamic simulations of the system behaviour on a computer.

A general chemical reaction takes the form

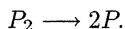


where r is the number of reactants and p is the number of products. R_i represents the i th reactant molecule and P_j is the j th product molecule. m_i is the number of molecules of R_i consumed in a single reaction step, and n_j is the number of molecules of P_j produced in a single reaction step. The coefficients m_i and n_j are known as *stoichiometries*. The stoichiometries are usually (though not always) assumed to be integers, and in this case it is assumed that there is no common factor of the stoichiometries. That is, it is assumed that there is no integer greater than one which exactly divides each stoichiometry on both the left and right sides. There is no assumption that the R_i and P_j are distinct, and it is perfectly reasonable for a given molecule to be both consumed and produced by a single reaction.[†] The reaction equation describes precisely which chemical species[‡] react together, and in what proportions, along with what is produced.

In order to make things more concrete, consider the dimerisation of a protein P . This is normally written



as two molecules of P react together to produce a single molecule of P_2 . Here P has a stoichiometry of 2 and P_2 has a stoichiometry of 1. Stoichiometries of 1 are not usually written explicitly. Similarly, the reaction for the dissociation of the dimer would be written

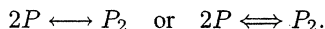


A reaction that can happen in both directions is known as *reversible*. Reversible re-

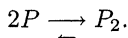
[†] Note that a chemical species that occurs on both the left and right hand sides with the same stoichiometry is somewhat special, and is sometimes referred to as a *modifier*. Clearly the reaction will have no effect on the amount of this species. Such a species is usually included in the reaction because the *rate* at which the reaction proceeds depends on the level of this species.

[‡] The use of the term "species" to refer to a particular type of molecule will be explained later in the chapter.

actions are quite common in biology and tend not to be written as two separate reactions. They can be written with a double-headed arrow such as



If one direction predominates over the other, this is sometimes emphasised in the notation. So, if the above protein prefers the dimerised state, this may be written something like



It is important to remember that the notation for a reversible reaction is simply a convenient shorthand for the two separate reaction processes taking place. In the context of the discrete stochastic models to be studied in this book, it will not usually be acceptable to replace the two separate reactions by a single reaction proceeding at some kind of combined rate.

1.5 Modelling genetic and biochemical networks

Before moving on to look at different ways of representing and working with systems of coupled chemical reactions in the next chapter, it will be helpful to end this chapter by looking in detail at some basic biochemical mechanisms and how their essential features can be captured with fairly simple systems of coupled chemical reactions. Although biological modelling can be applied to biological systems at a variety of different scales, it turns out that stochastic effects are particularly important and prevalent at the scale of genetic and biochemical networks, and these will therefore provide the main body of examples for this book.

1.5.1 Transcription (*prokaryotes*)

Transcription is a key cellular process, and control of transcription is a fundamental regulation mechanism. As a result, virtually any model of genetic regulation is likely to require some modelling of the transcription process. This process is much simpler in prokaryotic organisms, so it will be helpful to consider this in the first instance. Here, typically, a promoter region exists just upstream of the gene of interest. RNA-polymerase (RNAP) is able to bind to this promoter region and initiate the transcription process, which ultimately results in the production of an mRNA transcript and the release of RNAP back into the cell. The transcription process itself is complex, but whether it will be necessary to model this explicitly will depend very much on the modelling goals. If the modeller is primarily interested in control and the downstream effects of the transcription process, it may not be necessary to model transcription itself in detail.

The process is illustrated diagrammatically in Figure 1.4. Here, g is the gene of interest, p is the upstream promoter region, and r is the mRNA transcript of g . A very simple representation of this process as a system of coupled chemical reactions can

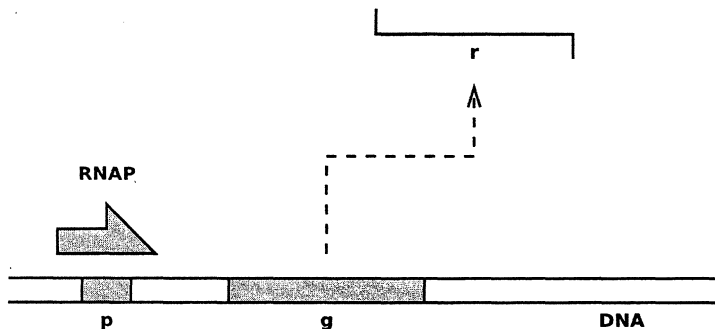
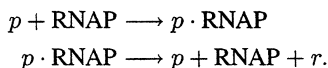


Figure 1.4 *Transcription of a single prokaryotic gene*

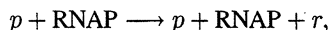
be written as follows:



As discussed, the second reaction is really the end result of a very large number of reactions. It is also worth emphasising that the reactions do not represent a closed system, as r appears to be produced out of thin air. In reality, it is created from other chemical species within the cell, but we have chosen here not to model at such a fine level of detail. One detail not included here that may be worth considering is the reversible nature of the binding of RNAP to the promoter region. It is also worth noting that these two reactions form a simple linear chain, whereby the product of the first reaction is the reactant for the second. Indeed, we could write the pair of reactions as



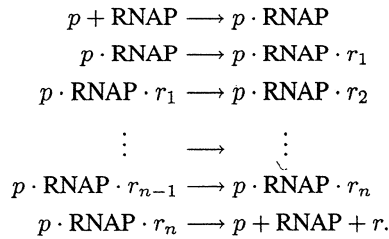
It is therefore tempting to summarise this chain of reactions by the single reaction



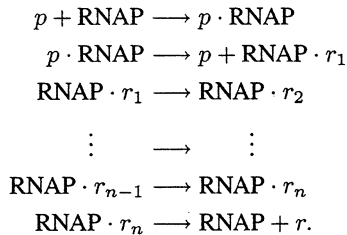
but this is likely to be inadequate for any model of regulation or control where the intermediary compound $p \cdot \text{RNAP}$ is important, such as any model for competitive binding of RNAP and a repressor in the promoter region.

If modelling the production of the entire RNA molecule in a single step is felt to be an oversimplification, it is relatively straightforward to model the explicit elongation of the molecule. As a first attempt, consider the following model for the transcription

of an RNA molecule consisting of n nucleotides.



This still does not model the termination process in detail; see Arkin, Ross, & McAdams (1998) for details of how this could be achieved. One problem with the above model is that the gene is blocked in the first reaction and is not free for additional RNAP binding until it is released again after the last reaction. This prevents concurrent transcription from occurring. An alternative would be to model the process as follows:



This model frees the gene for further transcription as soon as the transcription process starts. In fact, it is probably more realistic to free the gene once a certain number of nucleotides have been transcribed. Another slightly undesirable feature of the above model is that it does not prevent one RNAP from “overtaking” another during concurrent transcription (but this is not particularly important or easy to fix).

1.5.2 Eukaryotic transcription (a very simple case)

The transcription process in eukaryotic cells is rather more complex than in prokaryotes. This book is not an appropriate place to explore the many and varied mechanisms for control and regulation of eukaryotic transcription, so we will focus on a simple illustrative example, shown in Figure 1.5. In this model there are two transcription factor (TF) binding sites upstream of a gene, g . Transcription factor TF1 reversibly binds to site $tf1$, and TF2 reversibly binds to $tf2$, but is only able to bind if TF1 is already in place. Also, TF1 cannot dissociate if TF2 is in place. The transcription process cannot initiate (starting with RNAP binding) unless both TFs are in place.

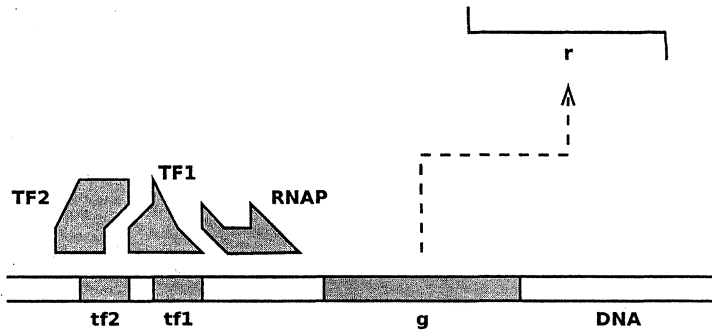
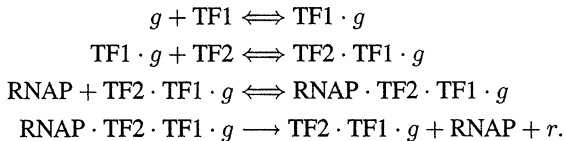


Figure 1.5 A simple illustrative model of the transcription process in eukaryotic cells

We can model this as follows:



Note that we have not explicitly included *tf1*, *tf2*, and *g* separately in the model, as they are all linked on a DNA strand and hence are a single entity from a modelling perspective. Instead we use *g* to represent the gene of interest together with its regulatory region (including *tf1* and *tf2*). Note that this system, like the previous, also forms a linear progression of ordered reactions and does not involve a “feedback” loop of any sort.

1.5.3 Gene repression (prokaryotes)

Regulation and control are fundamental to biological systems. These necessarily involve feedback and a move away from a simple ordered set of reactions (hence the term biochemical *network*). Sometimes such systems are large and complex, but feedback, and its associated non-linearity, can be found in small and apparently simple systems. We will look here at a simple control mechanism (repression of a prokaryotic gene) and see how this can be embedded into a regulatory feedback system in a later example.

Figure 1.6 illustrates a model where a repressor protein *R* can bind to regulatory site *q*, downstream of the RNAP binding site *p* but upstream of the gene *g*, thus preventing transcription of *g*. We can formulate a set of reactions for this process in

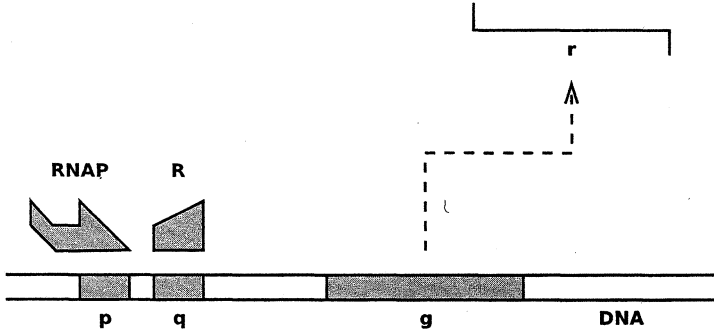
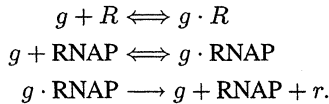


Figure 1.6 A simple prokaryotic transcription repression mechanism

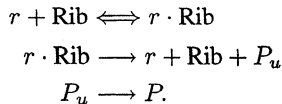
the following way:



This set of equations no longer has a natural ordering and hence cannot be read from top to bottom to go from reactants to products in an obvious way. Each reaction represents a possible direction for the system to move in. Also note that there are actually five reactions represented here, as two of the three listed are reversible. The crucial thing to appreciate here is that from a modelling perspective, $g \cdot R$ is a different species to g , and so the fact that RNAP can bind to g does not suggest that RNAP can bind to $g \cdot R$. Thus, this set of reactions precisely captures the mechanism of interest; namely that RNAP can bind to g when it is free but not when it is repressed by R .

1.5.4 Translation

Translation (like transcription) is a complex process involving several hundred reactions to produce a single protein from a single mRNA transcript. Again, however, it will not always be necessary to model every aspect of the translation process — just those features pertinent to system features of interest. The really key stages of the translation process are the binding of a ribosome (Rib) to the mRNA, the translation of the mRNA, and the folding of the resulting polypeptide chain into a functional protein. These stages are easily coded as a set of reactions.



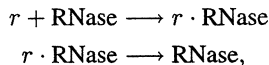
Here, P_u denotes unfolded protein and P denotes the folded protein. In some situations it will also be necessary to model various post-translational modifications. Clearly the second and third reactions are gross simplifications of the full translation process. Elongation can be modelled in more detail using an approach similar to that adopted for transcription; see Arkin, Ross & McAdams (1998) for further details. Folding could also be modelled similarly if necessary.

1.5.5 Degradation

The simplest model for degradation is just

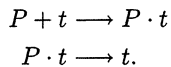


where \emptyset is the “empty set” symbol, meaning that r is transformed to nothing (as far as the model is concerned). A more appropriate model for RNA degradation would be



where RNase denotes ribonuclease (an RNA-degrading enzyme). Modelling in this way is probably only important if there is limited RNase availability, but is interesting in conjunction with a translation model involving Rib, as it will then capture the competitive binding of Rib and RNase to r .

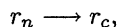
Models for protein degradation can be handled similarly. Here one would typically model the tagging of the protein with a cell signalling molecule, t , and then subsequent degradation in a separate reaction. A minimal model would therefore look like the following:



In fact, cell protein degradation machinery is rather complex; see Proctor et al. (2005) for a more detailed treatment of this problem.

1.5.6 Transport

In eukaryotes, mRNA is transported out of the cell nucleus before the translation process can begin. Often, the modelling of this process will be unnecessary, but could be important if the transportation process itself is of interest, or if the number of available transportation points is limited. A model for this could be as simple as



where r_n denotes the mRNA pool within the nucleus, and r_c the corresponding pool in the cytoplasm. However, this would not take into account the limited number of

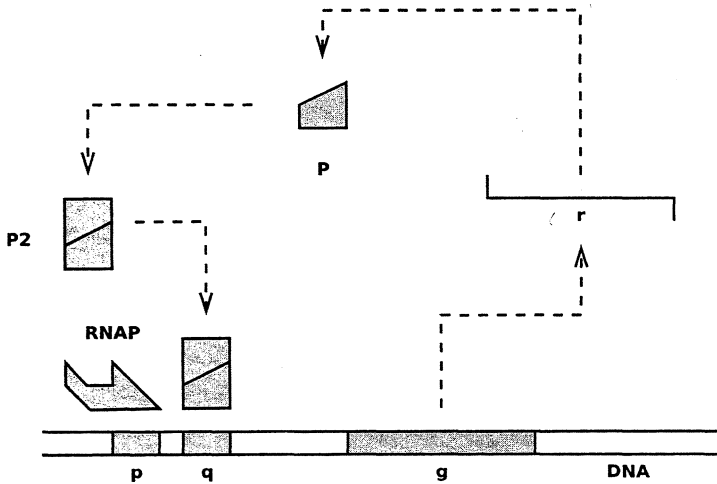
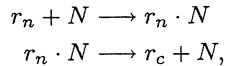


Figure 1.7 A very simple model of a prokaryotic auto-regulatory gene network. Here dimers of a protein P coded for by a gene g repress their own transcription by binding to a regulatory region q upstream of g and downstream of the promoter p .

transport points. A more realistic model would therefore be



where N denotes the set of available mRNA transport points embedded in the outer shell of the cell nucleus. In fact, this system is very closely related to the Michaelis-Menten enzyme kinetic system that will be examined in more detail later in the book. Here, the transport points behave like an enzyme whose abundance limits the flow from r_n to r_c .

1.5.7 Prokaryotic auto-regulation

Now that we have seen how to generate very simple models of key processes involved in gene expression and regulation, we can put them together in the form of a simple prokaryotic auto-regulatory network.

Figure 1.7 illustrates a simple gene expression auto-regulation mechanism often present in prokaryotic gene networks. Here, dimers of the protein P coded for by the gene g repress their own transcription by binding to a (repressive) regulatory region upstream of g .

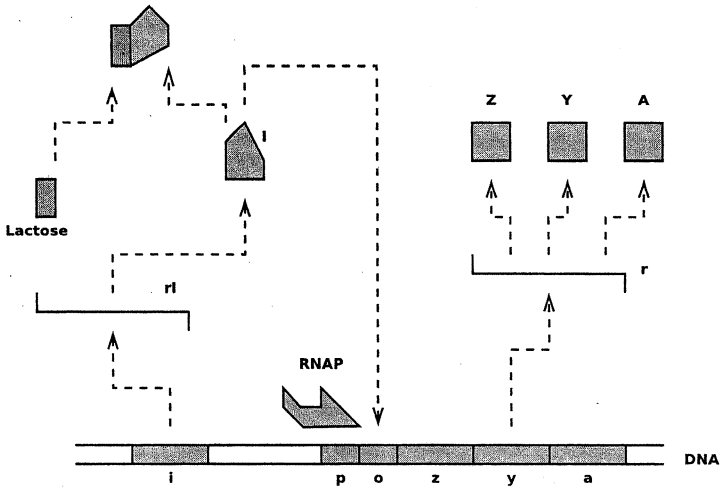
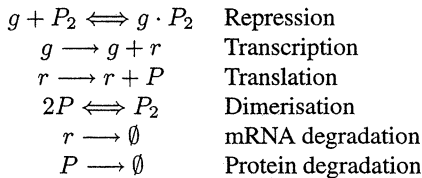


Figure 1.8 Key mechanisms involving the lac operon. Here an inhibitor protein I can repress transcription of the lac operon by binding to the operator o . However, in the presence of lactose, the inhibitor preferentially binds to it, and in the bound state can no longer bind to the operator, thereby allowing transcription to proceed.

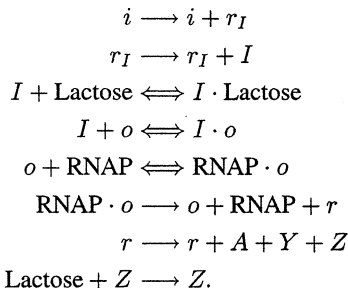


Notice that this model is minimal in terms of the level of detail included. In particular, the transcription part ignores RNAP binding, the translation/mRNA degradation parts ignore Rib/RNase competitive binding, and so on. However, as we will see later, this model contains many of the interesting features of an auto-regulatory feedback network.

1.5.8 lac operon

We will finish this section by looking briefly at a classic example of prokaryotic gene regulation — probably the first well-understood genetic regulatory network. The genes in the operon code for enzymes required for the respiration of lactose (Figure 1.8). That is, the enzymes convert lactose to glucose, which is then used as the “fuel” for respiration in the usual way. These enzymes are only required if there is a shortage of glucose and an abundance of lactose, and so there is a transcription control mechanism regulating their production. Upstream of the *lac* operon there is a gene coding for a protein which represses transcription of the operon by binding

to the DNA just downstream of the RNAP binding site. Under normal conditions (absence of lactose), transcription of the *lac* operon is turned off. However, in the presence of lactose, the inhibitor protein preferentially binds to lactose, and in the bound state can no longer bind to the DNA. Consequently, the repression of transcription is removed, and production of the required enzymes can take place. We can represent this with the following simple set of reactions:

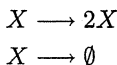


Here, i represents the gene for the inhibitor protein, r_I the associated mRNA, and I the inhibitor protein itself. The *lac* operon is denoted o and is treated as a single entity from a modelling viewpoint. The mRNA transcript from the operon is denoted by r , and this codes for all three *lac* proteins. The final reaction represents the transformation of lactose to something not directly relevant to the regulation mechanism.

Again, this system is fairly minimal in terms of the detail included, and all of the degradation reactions have been omitted, along with what happens to lactose once it has been acted on by β -galactosidase (Z). In fact, there is also another mechanism we have not considered here that ensures that transcription of the operon will only occur when there is a shortage of glucose (as respiration of glucose is always preferred).

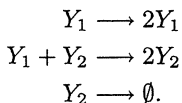
1.6 Modelling higher-level systems

We have concentrated so far on fairly low-level biochemical models where the concept of modelling with “chemical reactions” is perhaps most natural. However, it is important to recognise that we use the notation of chemical reactions simply to describe things that combine and the things that they produce, and that this framework can be used to model higher-level phenomena in a similar way. In Section 1.3 the linear birth-death process was introduced as a model for the number of bacteria present in a colony. We can use our chemical reaction notation to capture the qualitative structure of this model.



The first equation represents “birth” of new bacteria and the second “death.” There are many possible extensions of this simple model, including the introduction of immigration of new bacteria from another source and emigration of bacteria to another source.

The above model represents a model for a “population” of individuals (here the individuals are bacteria), and it is possible to extend such models to populations involving more than one “species.” Consider the Lotka-Volterra predator-prey model for two interacting species:



Again, this is not a real reaction system in the strictest sense, but it is interesting and useful, as it is the simplest model exhibiting the kind of non-linear auto-regulatory feedback behaviour considered earlier. Also, as it only involves two species and three reactions, it is relatively easy to work with without getting lost in detail. Here, Y_1 represents a “prey” species (such as rabbits) and Y_2 represents a “predator” species (such as foxes).[§] The first reaction is a simple representation of prey reproduction. The second reaction is an attempt to capture predator-prey interaction (consumption of prey by predator, in turn influencing predator reproduction rate). The third reaction represents death of predators due to natural causes. We will revisit this model in greater detail in later chapters.

1.7 Exercises

1. Write out a more detailed and realistic model for the simple auto-regulatory network considered in Section 1.5.7. Include RNAP binding, Rib/RNase competitive binding, and so on.
2. Consider the *lac* operon model from Section 1.5.8.
 - (a) First add more detail to the model, as in the previous exercise.
 - (b) Look up the β -galactosidase pathway and add detail from this to the model.
 - (c) Find details of the additional regulation mechanism mentioned that ensures lactose is only respired in an absence of glucose and try to incorporate that into the model.

1.8 Further reading

See Bower & Bolouri (2000) for more detailed information on modelling, and the different possible approaches to modelling genetic and biochemical networks. Kitano (2001) gives a more general overview of biological modelling and systems biology. McAdams & Arkin (1997) and Arkin et al. (1998) explore biological modelling in the context of the discrete stochastic models we will consider later. The *lac* operon is discussed in many biochemistry texts, including Stryer (1988). The original references for the Lotka-Volterra predator-prey models are Lotka (1925) and Volterra (1926).

[§] Note that the use of reactions to model the interaction of “species” in a population dynamics context explains the use of the term “species” to refer to a particular type of chemical molecule in a set of coupled chemical reactions.

The website associated with this book[¶] contains a range of links to online information of relevance to the various chapters of the book. Now would probably be a good time to have a quick look at it and “bookmark” it for future reference.

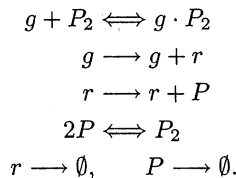
[¶] URL: <http://www.staff.ncl.ac.uk/d.j.wilkinson/smfsb/>

Representation of biochemical networks

2.1 Coupled chemical reactions

As was illustrated in the first chapter, a powerful and flexible way to specify a model is to simply write down a list of reactions corresponding to the system of interest. Note, however, that the reactions themselves specify only the qualitative structure of a model and must be augmented with additional information before they can be used to carry out a dynamic simulation on a computer. The model is completed by specifying the *rate* of every reaction, together with initial amounts of each reacting species.

Reconsider the auto-regulation example from Section 1.5.7:



Although only six reactions are listed, there are actually eight, as two are reversible. Each of those eight reactions must have a rate law associated with it. We will defer discussion of rate laws until Chapter 6. For now, it is sufficient to know that the rate laws quantify the propensity of particular reactions to take place and are likely to depend on the *current* amounts of available reactants. In addition there must be an *initial* amount for each of the five chemical species involved: $g \cdot P_2$, g , r , P , and P_2 . Given the reactions, the rate laws, and the initial amounts (together with some assumptions regarding the underlying kinetics, which are generally not regarded as part of the model), the model is specified and can in principle be simulated dynamically on a computer.

The problem is that even this short list of reactions is hard to understand on its own, whereas the simple biologist's diagram (Figure 1.7) is not sufficiently detailed and explicit to completely define the model. What is needed is something between the biologist's diagram and the list of reactions.

2.2 Graphical representations

2.2.1 Introduction

One way to begin to understand a reaction network is to display it as a pathway diagram of some description. The diagram in Figure 2.1 is similar to that used by

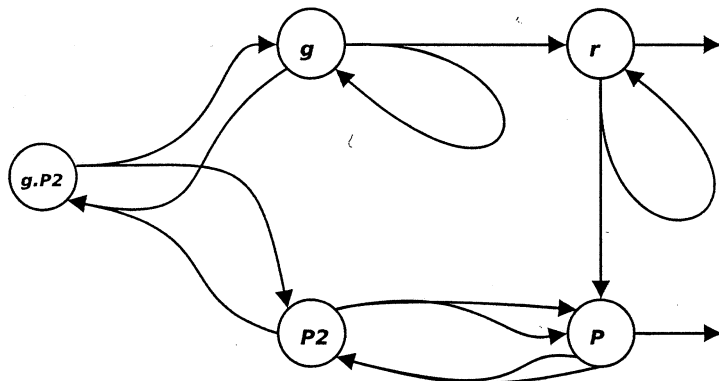


Figure 2.1 A simple graph of the auto-regulatory reaction network

some biological model-building tools such as JDesigner, which is part of the Systems Biology Workbench (SBW).*

Such a diagram is easier to understand than the reaction list, yet it contains the same amount of information and hence could be used to generate a reaction list. Note that auto-regulation by its very nature implies a “loop” in the reaction network, which is very obvious and explicit in the associated diagram. One possible problem with diagrams such as this, however, is the fact that it can sometimes be hard to distinguish which are the reactants and which are the products in any given reaction (particularly in large complex networks), and this can make it difficult to understand the flow of species through the network. Also, the presence of “branching” and “looping” arcs makes them slightly unnatural to work with directly in a mathematical sense.

Such problems are easily overcome by formalising the notion of pathway diagrams using the concept of a *graph* (here we mean the mathematical notion of a graph, not the idea of the graph of a function), where each *node* represents either a chemical species or a reaction, and arcs are used to indicate reaction pathways. In order to make this explicit, some elementary graph theoretic notation is helpful.

2.2.2 Graph theory

Definition 2.1 A directed graph or digraph, \mathcal{G} is a tuple (V, E) , where $V = \{v_1, \dots, v_n\}$ is a set of nodes (or vertices) and $E = \{(v_i, v_j) | v_i, v_j \in V, v_i \rightarrow v_j\}$ is a set of directed edges (or arcs), where we use the notation $v_i \rightarrow v_j$ if and only if there is a directed edge from node v_i to v_j .

So the graph shown in Figure 2.2 has mathematical representation

$$\mathcal{G} = (\{a, b, c\}, \{(a, c), (c, b)\}).$$

* Software web links tend to go out of date rather quickly, so a regularly updated list is available from the book's web page.

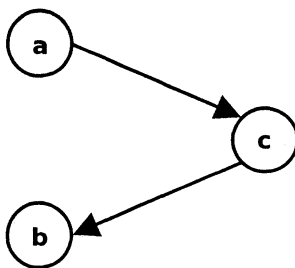


Figure 2.2 A simple digraph

Definition 2.2 A graph is described as simple if there do not exist edges of the form (v_i, v_i) and there are no repeated edges. A bipartite graph is a simple graph where the nodes are partitioned into two distinct subsets V_1 and V_2 (so that $V = V_1 \cup V_2$ and $V_1 \cap V_2 = \emptyset$) such that there are no arcs joining nodes from the same subset.

Referring back to the previous example, the partition $V_1 = \{a, b\}$, $V_2 = \{c\}$ gives a bipartite graph (\mathcal{G} is said to be bipartite over the partition), and the partition $V_1 = \{a\}$, $V_2 = \{b, c\}$ does not (as the edge (c, b) would then be forbidden). A *weighted* graph is a graph which has (typically positive) numerical values associated with each edge.

2.2.3 Reaction graphs

It turns out that it is very natural to represent sets of coupled chemical reactions using weighted bipartite graphs where the nodes are partitioned into two sets representing the species and reactions. An arc from a species node to a reaction node indicates that the species is a reactant for that reaction, and an arc from a reaction node to a species node indicates that the species is a product of the reaction. The weights associated with the arcs represent the stoichiometries associated with the reactants and products. There is a very strong correspondence between reaction graphs modelled this way, and the theory of Petri nets, which are used extensively in computing science for a range of modelling problems. A particular advantage of Petri net theory is that it is especially well suited to the discrete-event stochastic simulation models this book is mainly concerned with. It is therefore helpful to have a basic familiarity with Petri nets and their application to biological modelling.

2.3 Petri nets

2.3.1 Introduction

Petri nets are a mathematical framework for systems modelling together with an associated graphical representation. Goss & Peccoud (1998) were among the first to use stochastic Petri nets for biological modelling. Recent reviews of the use of Petri nets for biological modelling include Pinney, Westhead & McConkey (2003) and

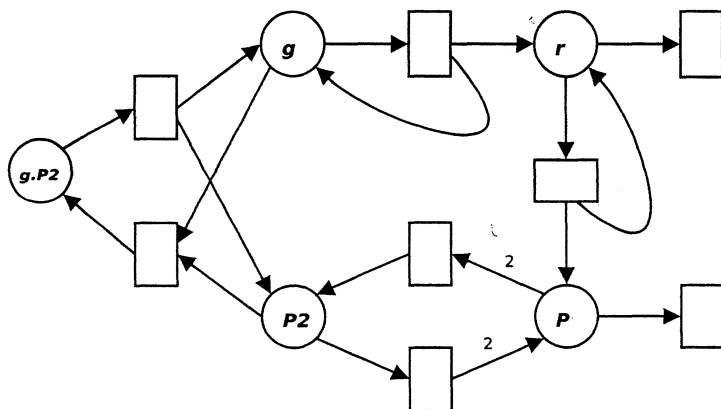


Figure 2.3 A Petri net for the auto-regulatory reaction network

Hardy & Robillard (2004). A brief introduction to key elements of Petri net theory will be outlined here — see Reisig (1985) and Murata (1989) for further details.

By way of an informal introduction, the Petri net corresponding to the network shown in Figure 2.1 is shown in Figure 2.3. The rectangular boxes in Figure 2.3 represent individual reactions. The arcs into each box denote reactants, and the arcs out of each box denote products. Numbers on arcs denote the *weight* of the arc (un-numbered arcs are assumed to have a weight of 1). The weights represent reaction stoichiometries. The Petri net graph is only a slight refinement of the basic graph considered earlier, but it is easier to comprehend visually and more convenient to deal with mathematically (it is a bipartite graph). It is easy to see how to work through the graph and generate the full list of chemical reactions.

Traditionally, each place (species) node of a Place/Transition (P/T) Petri net has an integer number of “tokens” associated with it, representing the abundance of that “species.” This fits in particularly well with the discrete stochastic molecular kinetics models we will consider in more detail later. Here, the number of tokens at a given node may be interpreted as the number of molecules of that species in the model at a given time (Figure 2.4). The collection of all token numbers at any given point in time is known as the current *marking* of the net (which corresponds here to the *state* of the reaction system). The Petri net shows what happens when particular transitions “fire” (reactions occur). For example, in the above state, if two reactions occur, one a repression binding and the other a translation, the new Petri net will be as given in Figure 2.5. This is because the repression binding has the effect of increasing $g \cdot P_2$ by 1, and decreasing g and P_2 by 1, and the translation has the effect of increasing P by 1 (as r is both increased and decreased by 1, the net effect is that it remains unchanged). So the old and new Petri net markings can be written as

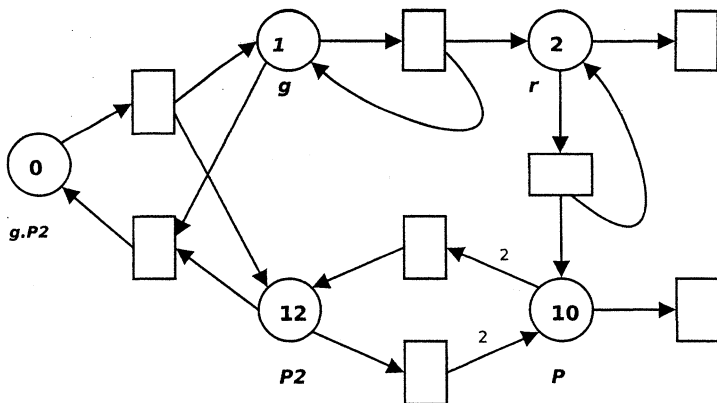


Figure 2.4 A Petri net labelled with tokens

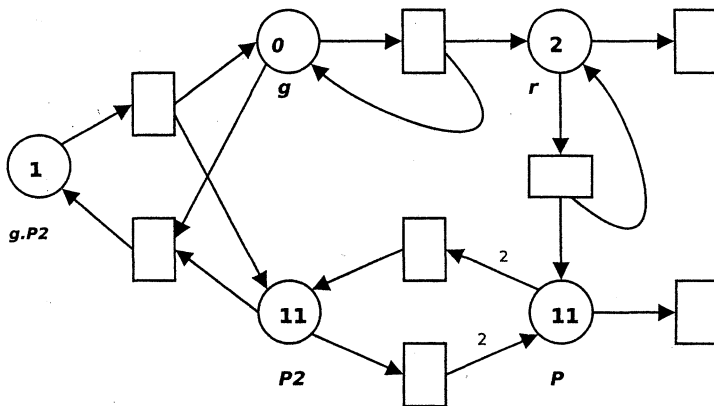


Figure 2.5 A Petri net with new numbers of tokens after reactions have taken place

Species	No. tokens
$g \cdot P_2$	0
g	1
r	2
P	10
P_2	12

and

Species	No. tokens
$g \cdot P_2$	1
g	0
r	2
P	11
P_2	11

Table 2.1 *The auto-regulatory system displayed in tabular (matrix) form (zero stoichiometries omitted for clarity)*

Species	Reactants (<i>Pre</i>)					Products (<i>Post</i>)				
	$g \cdot P_2$	g	r	P	P_2	$g \cdot P_2$	g	r	P	P_2
Repression		1			1	1				
Reverse repression	1								1	1
Transcription		1					1	1		
Translation			1					1	1	
Dimerisation				2						1
Dissociation					1				2	
mRNA degradation		1								
Protein degradation				1						

respectively. A transition (reaction) can only fire (take place) if there are sufficiently many tokens (molecules) associated with each input place (reactant species). We are now in a position to consider Petri nets more formally.

2.3.2 Petri net formalism and matrix representations

Definition 2.3 A Petri net, N , is an n -tuple $(P, T, Pre, Post, M)$, where $P = \{p_1, \dots, p_u\}$, ($u > 0$) is a finite set of places, $T = \{t_1, \dots, t_v\}$, ($v > 0$) is a finite set of transitions, and $P \cap T = \emptyset$. Pre is a $v \times u$ integer matrix containing the weights of the arcs going from places to transitions (the (i, j) th element of this matrix is the weight of the arc going from place j to transition i), and $Post$ is a $v \times u$ integer matrix containing the weights of arcs from transitions to places (the (i, j) th element of this matrix is the weight of the arc going from transition i to place j).[†] Note that Pre and $Post$ will both typically be sparse matrices.[‡] M is a u -dimensional integer vector representing the current marking of the net (i.e. the current state of the system).

The initial marking of the net is typically denoted M_0 . Note that the form of Pre and $Post$ ensure that arcs only exist between nodes of different types, so the resulting network is a bipartite graph. A particular transition, t_i can only fire if $M_j \geq Pre_{ij}$, $j = 1, \dots, u$.

In order to make this concrete, let us now write out the reaction list for Figure 2.3 in the form of a table, shown in Table 2.1. We can then use this to give a formal Petri net specification of the system.

[†] Non-existent arcs have a weight of zero.

[‡] A sparse matrix is a matrix consisting mainly of zeros. There are special algorithms for working with sparse matrices that are much more efficient than working with the full matrix directly. It is hard to be precise about how sparse a matrix has to be before it is worth treating as a sparse matrix, but for an $n \times n$ matrix, if the number of non-zero elements is closer to order n than order n^2 , it is likely to be worthwhile using sparse matrix algorithms.

Table 2.2 Table representing the overall effect of each transition (reaction) on the marking (state) of the network

Species	$g \cdot P_2$	g	r	P	P_2
Repression	1	-1	0	0	-1
Reverse repression	-1	1	0	0	1
Transcription	0	0	1	0	0
Translation	0	0	0	1	0
Dimerisation	0	0	0	-2	1
Dissociation	0	0	0	2	-1
mRNA degradation	0	0	-1	0	0
Protein degradation	0	0	0	-1	0

$$N = (P, T, Pre, Post, M), P = \begin{pmatrix} g \cdot P_2 \\ g \\ r \\ P \\ P_2 \end{pmatrix}, T = \begin{pmatrix} \text{Repression} \\ \text{Reverse repression} \\ \text{Transcription} \\ \text{Translation} \\ \text{Dimerisation} \\ \text{Dissociation} \\ \text{mRNA degradation} \\ \text{Protein degradation} \end{pmatrix}$$

$$Pre = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}, Post = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, M = \begin{pmatrix} 0 \\ 1 \\ 2 \\ 10 \\ 12 \end{pmatrix}$$

Now, when a particular transition (reaction) occurs (given by a particular row of Table 2.1), the numbers of tokens associated with each place (species) will decrease according to the numbers on the LHS (Pre) and increase according to the numbers on the RHS ($Post$). So, it is the *difference* between the RHS and the LHS that is important for calculating the change in state associated with a given transition (or reaction). We can write this matrix out in the form of a table, shown in Table 2.2, or more formally as a matrix

$$A = Post - Pre = \begin{pmatrix} 1 & -1 & 0 & 0 & -1 \\ -1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -2 & 1 \\ 0 & 0 & 0 & 2 & -1 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & -1 & 0 \end{pmatrix}$$

This “net effect” matrix, A , is of fundamental importance in the theory and application of Petri nets and chemical reaction networks. Unfortunately there is no agreed standard notation for this matrix within either the Petri net or biochemical network literature. Within the Petri net literature, it is usually referred to as the *incidence matrix* and is often denoted by the letter A , C or I .[§] Within the biochemical network field, it is usually referred to as the *reaction* or *stoichiometry matrix* and often denoted by the letter A or S . As if this were not confusing enough, the matrix is often (but by no means always) defined to be the *transpose* of the matrix we have called A ,[¶] as this is often more convenient to work with. Suffice to say that care needs to be taken when exploring and interpreting the wider literature in this area. For clarity and convenience throughout this book, A (the reaction matrix) will represent the matrix as we have already defined it, and S (the stoichiometry matrix), will be used to denote its transpose.

Definition 2.4 *The reaction matrix*

$$A = Post - Pre$$

is the $v \times u$ -dimensional matrix whose rows represent the effect of individual transitions (reactions) on the marking (state) of the network. Similarly, the stoichiometry matrix

$$S = A'$$

is the $u \times v$ -dimensional matrix whose columns represent the effect of individual transitions on the marking of the network.^{||}

However, it must be emphasised that this is not a universally adopted notation.

Now, suppose that we have some reactions. For example, suppose we have one repression binding reaction and one translation reaction. We could write this list of reactions in a table as

[§] I is a particularly bad choice, as this is typically used in linear algebra to denote the *identity matrix*, which has 1s along the diagonal and zeros elsewhere.

[¶] The transpose of a matrix is the matrix obtained by interchanging the rows and columns of a matrix. So in this case it would be a matrix where the rows represented places (species) and the columns represented transitions (reactions).

^{||} Here and elsewhere, the notation $'$ is used to denote the transpose of a vector or matrix.

Reaction	No. transitions
Repression	1
Reverse repression	0
Transcription	0
Translation	1
Dimerisation	0
Dissociation	0
mRNA degradation	0
Protein degradation	0

or more neatly as a vector

$$r = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

Note that in order to save space, column vectors such as r are sometimes written as the transpose of row vectors, e.g. $r = (1, 0, 0, 1, 0, 0, 0, 0)'$. We can now use matrix algebra to update the marking (state) of the network.

Proposition 2.1 *If r represents the transitions that have taken place subsequent to the marking M , the new marking \tilde{M} is related to the old marking via the matrix equation***

$$\tilde{M} = M + Sr. \tag{2.1}$$

Note that this equation (2.1) is of fundamental importance both for the mathematical analysis of Petri nets and biochemical networks and also for the development of simulation and inference algorithms. We will use this equation extensively in a variety of different ways throughout the book. Before we justify it, it is probably helpful to see a simple and direct application of it in practice.

In the context of our example, we can compute the new marking from the old

** Note that in matrix equations, addition is defined element-wise, but multiplication is defined in a special way. The product of the $n \times m$ matrix A and the $m \times p$ matrix B is the $n \times p$ matrix whose (i, j) th element is $\sum_{k=1}^m a_{ik}b_{kj}$. A vector is treated as a matrix with one column.

marking as

$$\begin{aligned}
 \tilde{M} &= M + Sr \\
 &= M + A'r \\
 &= \begin{pmatrix} 0 \\ 1 \\ 2 \\ 10 \\ 12 \end{pmatrix} + \begin{pmatrix} 1 & -1 & 0 & 0 & -1 \\ -1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -2 & 1 \\ 0 & 0 & 0 & 2 & -1 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & -1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \\
 &= \begin{pmatrix} 0 \\ 1 \\ 2 \\ 10 \\ 12 \end{pmatrix} + \begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 & -2 & 2 & 0 & -1 \\ -1 & 1 & 0 & 0 & 1 & -1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \\
 &= \begin{pmatrix} 0 \\ 1 \\ 2 \\ 10 \\ 12 \end{pmatrix} + \begin{pmatrix} 1 \\ -1 \\ 0 \\ 1 \\ -1 \end{pmatrix} \\
 &= \begin{pmatrix} 1 \\ 0 \\ 2 \\ 11 \\ 11 \end{pmatrix}.
 \end{aligned}$$

We can see that (2.1) appears to have worked in this particular example, but we need to establish why it is true in general.

Proof. The i th row of A represents the effect on the marking of the i th reaction. Similarly the i th column of S , which we denote s^i . Clearly r_i is the number of type i reactions that take place, so the change in marking, $\tilde{M} - M$ is given by

$$\begin{aligned}
 \tilde{M} - M &= s^1 r_1 + \cdots + s^v r_v \\
 &= \sum_{i=1}^v s^i r_i
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{i=1}^v S e_i r_i \\
 &= S \sum_{i=1}^v e_i r_i \\
 &= S r,
 \end{aligned}$$

where e_i is the i th unit v -vector.^{††} □

2.3.3 Network invariants and conservation laws

There are two Petri net concepts that are of particular relevance to biochemical networks: P - and T -invariants.

Definition 2.5 A P -invariant (sometimes referred to in the literature as an S -invariant) is a non-zero v -vector y that is a solution to the matrix equation $Ay = 0$. That is, y is any non-zero vector in the null-space of A .^{‡‡}

The null-space of A therefore characterises the set of P -invariants. These P -invariants are interesting because they correspond to *conservation laws* of the network.

In the example we have been studying, it is clear that the vector $y = (1, 1, 0, 0, 0)'$ is a P -invariant (as $Ay = 0$). This vector corresponds to the fairly obvious conservation law

$$g \cdot P_2 + g = \text{Constant.}$$

That is, the total number of copies of the gene does not change. It is true in general that if y is a P -invariant then the linear combination of states, $y'M$, is conserved by the reaction network. To see why this works, we can evaluate the current linear combination by computing $y'M$, where M is the current marking. Similarly, the value of the linear combination when the marking is \tilde{M} is $y'\tilde{M}$. So the change in the linear combination is

$$\begin{aligned}
 y'\tilde{M} - y'M &= y'(\tilde{M} - M) \\
 &= y'Sr \\
 &= (S'y)'r \\
 &= (Ay)'r \\
 &= 0
 \end{aligned}$$

where the second line follows from (2.1) and the last line follows from the fact that y is a P -invariant.^{§§}

Definition 2.6 A T -invariant is a non-zero, non-negative (integer-valued) v -vector x that is a solution to the matrix equation $Sx = 0$. That is, x is in the null-space of S .

^{††} The i th unit vector, e_i , is the vector with a 1 in the i th position and zeros elsewhere. Multiplying a matrix by e_i has the effect of picking out the i th column.

^{‡‡} The null-space of a matrix (sometimes known as the kernel) is defined to be the set of all vectors that get mapped to zero by the matrix.

^{§§} We also used the fact that for arbitrary (conformable) matrices A and B , we have $(AB)' = B'A'$.

These invariants are of interest because they correspond to sequences of transitions (reactions) that return the system to its original marking (state). This is clear immediately from equation (2.1). If the primary concern is continuous deterministic modelling, then any non-negative solution is of interest. However, if the main interest is in discrete stochastic models of biochemical networks, only non-negative integer vector solutions correspond to sequences of transitions (reactions) that can actually take place. Hence, we will typically want to restrict our attention to these.

In the example we have been considering it is easily seen that the vectors $x = (1, 1, 0, 0, 0, 0, 0, 0, 0)'$ and $\tilde{x} = (0, 0, 0, 0, 1, 1, 0, 0, 0)'$ are both T -invariants of our network. The first corresponds to a repression binding and its reverse reaction, and the second corresponds to a dimerisation and corresponding dissociation. However, not all T -invariants are associated with reversible reactions.

Although it is trivial to verify whether a given vector is a P - or T -invariant, it is perhaps less obvious how to systematically find such invariants and classify the set of all such invariants for a given Petri net. In fact, the Singular Value Decomposition (SVD) is a classic matrix algorithm (Golub & Van Loan 1996) that completely characterises the null-space of a matrix and its transpose and hence helps considerably in this task. However, if we restrict our attention to positive integer solutions, there is still more work to be done even once we have the SVD. We will revisit this issue in Chapter 10 when the need to find invariants becomes more pressing.

Before leaving the topic of invariants, it is worth exploring the relationship between the number of (linearly independent) P - and T -invariants. The *column-rank* of a matrix is the dimension of the image-space of the matrix (the space spanned by the columns of the matrix). The *row-rank* is the column-rank of the transpose of the matrix. It is a well-known result from linear algebra that the row and column ranks are the same, and so we can refer unambiguously to the *rank* of a matrix. It is fairly clear that the dimension of the image-space and null-space must sum to the dimension of the space being mapped into, which is the number of rows of the matrix. So, if we fix on S , which has dimension $u \times v$, suppose the rank of the matrix is k . Let the dimension of the null-space of S be p and the dimension of the null-space of $A (= S')$ be t . Then we have $k + p = u$ and $k + t = v$. This leads immediately to the following result.

Proposition 2.2 *The number of linearly independent P -invariants, p , and the number of linearly independent T -invariants, t , are related by*

$$t - p = v - u. \quad (2.2)$$

In the context of our example, we have $u = 5$ and $v = 8$. As we have found a P -invariant, we know that $p \geq 1$. Now using (2.2) we can deduce that $t \geq 4$. So there are at least four linearly independent T -invariants for this network, and we have so far found only two.

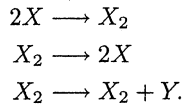
2.3.4 Reachability

Another Petri net concept of considerable relevance is that of *reachability*.

Definition 2.7 A marking \tilde{M} is reachable from marking M if there exists a finite sequence of transitions leading from M to \tilde{M} .

If such a sequence of transitions is summarised in the v -vector r , it is clear that r will be a non-negative integer solution to (2.1). However, it is important to note that the converse does not follow: the existence of a non-negative integer solution to (2.1) does not guarantee the reachability of \tilde{M} from M . This is because the markings have to be non-negative. A transition cannot occur unless the number of tokens at each place is at least that required for the transition to take place. It can happen that there exists a set of non-negative integer transitions between two valid markings M and \tilde{M} , but all possible sequences corresponding to this set are impossible.

This issue is best illustrated with an example. Suppose we have the reaction network



So, X can dimerise and dissociate, and dimers of X somehow catalyse the production of Y . If we formulate this as a Petri net with $P = (X, X_2, Y)$, and the transitions in the above order, we get the stoichiometry matrix

$$S = \begin{pmatrix} -2 & 2 & 0 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

If we begin by thinking about getting from initial marking $M = (2, 0, 0)'$ to marking $\tilde{M} = (2, 0, 1)'$, we see that it is possible and can be achieved with the sequence of transitions t_1, t_3 , and t_2 , giving reaction vector $r = (1, 1, 1)'$. We can now ask if marking $\tilde{M} = (1, 0, 1)'$ is reachable from $M = (1, 0, 0)'$. If we look for a non-negative integer solution to (2.1), we again find that $r = (1, 1, 1)'$ is appropriate, as $\tilde{M} - M$ is the same in both scenarios. However, this does not correspond to any legal sequence of transitions, as *no* transitions are legal from the initial marking $M = (1, 0, 0)'$. In fact, $r = (0, 0, 1)'$ is another solution, since $(1, 1, 0)'$ is a T-invariant. This solution is forbidden despite the fact that firing of t_3 will not cause the marking to go negative.

This is a useful warning that although the matrix representation of Petri net (and biochemical network) theory is powerful and attractive, it is not a complete characterisation — discrete-event systems analysis is a delicate matter, and there is much that systems biologists interested in discrete stochastic models can learn from the Petri net literature.

2.4 Systems Biology Markup Language (SBML)

Different representations of biochemical networks are useful for different purposes. Graphical representations (including Petri nets) are useful both for visualisation and analysis, and matrix representations are useful for mathematical and computational analysis. The Systems Biology Markup Language (SBML), described in Hucka et

al. (2003), is a way of representing biochemical networks that is intended to be convenient for computer software to generate and parse, thereby enabling communication of biochemical network models between disparate modelling and simulation tools. It is essentially an eXtensible Markup Language (XML) encoding (DuCharme 1999) of the reaction list, together with the additional information required for quantitative modelling and simulation. It is intended to be independent of particular kinetic theories and should be as appropriate for discrete stochastic models as for continuous deterministic ones (in fact, there are a few minor problems using SBML for discrete stochastic models, and these will be discussed as they become relevant).

We will concentrate here on the version of SBML known as “Level 2 (version 1),” as this is the current specification (at the time of writing) and contains sufficient features for the biochemical network models considered in this book. Further details regarding SBML, including the specification and XML Schema can be obtained from the SBML.org website. Note that SBML should perhaps not be regarded as an alternative to other representations, but simply as an electronic format which could in principle be used in conjunction with any of the representations we have considered. Also note that it is not intended that SBML models should be generated and manipulated “by hand” using a text editor, but rather by software tools which present to the user a more human-oriented representation. It is also worth bearing in mind that SBML continues to evolve. At the time of writing, SBML Level 2 (version 1) is the current specification, but Level 2 version 2 is in preparation, and many proposals are already in place for Level 3. The principal differences between Level 1 and Level 2 are that Level 2 supports the notion of “events” and encodes all mathematical formulae using MathML (an XML encoding for mathematical notation) rather than as strings containing algebraic expressions. However, there is another more subtle difference to be examined later which makes it difficult to correctly and unambiguously define models intended for discrete stochastic simulation in Level 1. This problem was addressed for Level 2, and this is the main reason for advocating that Level 2 is used in preference to Level 1 for the encoding of discrete stochastic models.

2.4.1 Basic document structure

An SBML (Level 2) model consists of lists of *functions*, *units*, *compartments*, *species*, *parameters*, *rules*, *reactions*, and *events*. Each of these lists is optional. We will concentrate here on units, compartments, species, parameters, and reactions, as these are sufficient for adequately describing most simple discrete stochastic models. This basic structure is encoded in SBML as follows.

```
<?xml version="1.0" encoding="UTF-8"?>
<sbml xmlns="http://www.sbml.org/sbml/level2"
      level="2" version="1">
  <model id="MyBiochemicalNetwork"
        name="My biochemical network">
    <listOfUnitDefinitions>
      ...
    </listOfUnitDefinitions>
    <listOfCompartments>
```

```

    ...
  </listOfCompartments>
  <listOfSpecies>
    ...
  </listOfSpecies>
  <listOfParameters>
    ...
  </listOfParameters>
  <listOfReactions>
    ...
  </listOfReactions>
</model>
</sbml>

```

2.4.2 Units

The (optional) units list allows definition and redefinition of the units used by the model. Discrete stochastic models will often have the following units declaration.

```

<listOfUnitDefinitions>
  <unitDefinition id="substance">
    <listOfUnits>
      <unit kind="item"/>
    </listOfUnits>
  </unitDefinition>
</listOfUnitDefinitions>

```

This declaration has the effect of changing the default substance units from the default value (*mole*) to *item*. The effect of this is that subsequent (unqualified) specifications of (and references to) amounts will be assumed to be in the unit of *item*. That is, amounts will be interpreted as numbers of molecules rather than the default of numbers of moles. Units turn out to be a rather delicate issue. We will examine units in some detail later in the book when we look at kinetics and rate laws. For further details on using units with SBML, see the specification document. Note that most simulators in current use ignore the units section of the SBML document. In practice, this means that deterministic simulators will assume units of mole, and most stochastic simulators will assume units of *item* irrespective of the content of this section. However, it is important to ensure that models are encoded correctly so that they are not misinterpreted later.

2.4.3 Compartments

The compartment list simply states the compartments in the model. So for a model with two nested compartments, the declaration might be as follows.

```

<listOfCompartments>
  <compartment id="Cell" size="1"/>
  <compartment id="Nucleus" size="1" outside="Cell"/>
</listOfCompartments>

```

A model must have at least one compartment, and each compartment should be given an *id*. You may also specify a size (or volume), in the current units (the default size units are *litres*). A compartment that is contained inside another compartment should declare the compartment that is “outside.”

2.4.4 Species

The species list simply states all species in the model. So, for the auto-regulatory network model we have been considering throughout this chapter, these could be declared using

```
<listOfSpecies>
  <species id="Gene" compartment="Cell" initialAmount="10"
           hasOnlySubstanceUnits="true"/>
  <species id="P2Gene" name="P2.Gene" compartment="Cell"
           initialAmount="0" hasOnlySubstanceUnits="true"/>
  <species id="Rna" compartment="Cell" initialAmount="0"
           hasOnlySubstanceUnits="true"/>
  <species id="P" compartment="Cell" initialAmount="0"
           hasOnlySubstanceUnits="true"/>
  <species id="P2" compartment="Cell" initialAmount="0"
           hasOnlySubstanceUnits="true"/>
</listOfSpecies>
```

There are several things worth pointing out about this declaration. First, the initial amounts are assumed to be in the default substance units, unless explicit units are specified (see the specification for how to do this). This default will be moles unless substance units have been redefined (say, to item). Also note that each species is declared with the non-default attribute *hasOnlySubstanceUnits*. This has the effect of ensuring that wherever the species is referred to elsewhere in the model (for example, in rate laws), it will be interpreted as an amount (in the appropriate substance units), and not a concentration (which is the SBML default). Most stochastic simulators will make this assumption anyway, but it is important to encode models correctly so that they will not be misinterpreted later. One of the main problems with SBML Level 1 from the viewpoint of discrete stochastic modellers is that it does not have the attribute *hasOnlySubstanceUnits*, and therefore references to species necessarily refer to concentrations rather than amounts, which makes specification of reaction rate laws very unnatural.

2.4.5 Parameters

The parameters section can be used to declare names for numeric values to be used in algebraic formulae. They are most often used to declare rate constants for the kinetic laws of biochemical reactions, but can be used for other variables as well. An example parameter section might be as follows.

```
<listOfParameters>
  <parameter id="k1" value="0.01"/>
```

```

<parameter id="k2" value="0.1"/>
</listOfParameters>

```

Parameters defined here are “global” to the whole model. In contrast, any parameters defined in the context of a particular reaction will be local to the kinetic law for that reaction only.

2.4.6 Reactions

The reaction list consists of a list of reactions. A reaction in turn consists of a list of reactants, a list of products and a rate law. A reaction may also declare “modifier” species — species that are not created or destroyed by the reaction, but figure in the rate law. For example, consider the following encoding of a dimerisation reaction.

```

<reaction id="Dimerisation" reversible="false">
  <listOfReactants>
    <speciesReference species="P" stoichiometry="2"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="P2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k4 </ci>
        <cn> 0.5 </cn>
        <ci> P </ci>
      <apply>
        <minus/>
        <ci> P </ci>
        <cn type="integer"> 1 </cn>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k4" value="1"/>
    </listOfParameters>
  </kineticLaw>
</reaction>

```

There are several things to note about this declaration. One thing that perhaps appears strange is that the reaction is declared to be not reversible when we know that dimerisation is typically a reversible process. However, reactions should only be flagged as reversible when the associated kinetic law represents the combined effects of both forward and backward reactions. It turns out that while this is fine for continuous deterministic models, there is no satisfactory way to do this for discrete stochastic models. As a result, when developing a model for discrete stochastic simulation, the forward and backward reactions must be specified separately (and declared not reversible) along with their separate kinetic laws. We will examine the meaning and

specification of reaction rates and kinetic laws later in this book. The specification of reactants and products and their associated stoichiometries is fairly self-explanatory. The kinetic law itself is a MathML (W3C 2000) encoding of the simple algebraic formula $k_4 * 0.5 * P * (P - 1)$. Rate laws will be discussed in more detail later, but it is important to know that the units of this law are of the form *substance / time*, using the default substance and time units (the default time unit is *second*). However, by default, any reference to a species within a rate law will refer to the *concentration* of that species (with units *substance / size*), unless the species has declared the attribute *hasOnlySubstanceUnits* (or is in a compartment with zero dimensions), in which case it will refer to the *amount* of that species (with units *substance*). The kinetic law given above uses a local parameter k_4 . This constant will always be used in the formula of the kinetic law, masking any global parameter of the same name. To use a global parameter called k_4 , the entire section

```
<listOfParameters>
  <parameter name="k4" value="1"/>
</listOfParameters>
```

should be removed from the kinetic law. Kinetic laws can use a mixture of global and local parameters. Any reference to a compartment will be replaced by the size (volume) of the compartment. A list of reactions should be included in the SBML file between `<listOfReactions>` and `</listOfReactions>` tags.

2.4.7 The full SBML model

The various model components are embedded into the basic model structure. For completeness, Appendix A.1 lists a full SBML model for the simple auto-regulatory network we have been considering. Note that as this model uses locally specified parameters rather than globally defined parameters, there is no `<listOfParameters>` section in the model definition. This model can also be downloaded from the book's website.

2.5 SBML-shorthand

2.5.1 Introduction

SBML is rapidly becoming the *lingua franca* for electronic representation of models of interest in systems biology. Dozens of different software tools provide SBML support to varying degrees, many providing both SBML import and export provisions, and some using SBML as their native format. However, while SBML is a good format for computers to parse and generate, its verbosity and low signal-to-noise ratio make it rather inconvenient for humans to read and write. I have found it helpful to develop a shorthand notation for SBML that is much easier for humans to read and write, and can easily be "compiled" into SBML for subsequent import into other SBML-aware software tools. The notation can be used as a partial substitute for the numerous GUI-based model-building tools that are widely used for systems biology model development. An additional advantage of the notation is that it is much more

suitable than raw SBML for presentation in a book such as this (because it is more concise and readable). Many of the examples discussed in subsequent chapters will be presented using the shorthand notation, so it is worth presenting the essential details here. Here we describe a particular version of the shorthand notation, known as 2.1.1. A compiler for translating the shorthand notation into full SBML is freely available (see the book's website for details).

2.5.2 Basic structure

The description format is plain ASCII text. The suggested file extension is `.mod`, but this is not required. All whitespace other than carriage returns is insignificant (unless it is contained within a quoted "name" element). Carriage returns are significant. The description is case-sensitive. Blank lines are ignored. The comment character is `#` — all text from a `#` to the end of the line is ignored.

The model description must begin with the characters `@model:2.1.1=` (the 2.1.1 corresponds to the version number of the specification). The text following the `=` on the first line is the model identification string (ID). An optional model name may also be specified, following the ID, enclosed in double quotes. The model is completed with the specification of the five sections, `@units`, `@compartments`, `@species`, `@parameters`, and `@reactions`, corresponding to the SBML sections, `<listOfUnitDefinitions>`, `<listOfCompartments>`, `<listOfSpecies>`, `<listOfParameters>`, and `<listOfReactions>`, respectively. The sections must occur in the stated order. Sections are optional, but if present, may not be empty. These are the only sections covered by this specification.

2.5.3 Units

The format of the individual sections will be explained mainly by example. The following SBML-shorthand

```
@units
substance=item
fahrenheit=celsius:m=1.8,o=32
mmols=mole:s=-3; litre:e=-1; second:e=-1
```

would be translated to

```
<listOfUnitDefinitions>
  <unitDefinition id="substance">
    <listOfUnits>
      <unit kind="item"/>
    </listOfUnits>
  </unitDefinition>
  <unitDefinition id="fahrenheit">
    <listOfUnits>
      <unit kind="Celsius" multiplier="1.8" offset="32"/>
    </listOfUnits>
  </unitDefinition>
```

```

<unitDefinition id="mmls">
  <listOfUnits>
    <unit kind="mole" scale="-3"/>
    <unit kind="litre" exponent="-1"/>
    <unit kind="second" exponent="-1"/>
  </listOfUnits>
</unitDefinition>
</listOfUnitDefinitions>

```

The unit attributes exponent, multiplier, scale, and offset are denoted by the letters e, m, s, and o respectively. Note that because there is no way to refer to units elsewhere in SBML-shorthand, the only function for this section is to redefine default units such as substance and size.

2.5.4 Compartments

The following SBML-shorthand

```

@compartments
cell=1
cytoplasm<cell=0.8
nucleus<cell=0.1
mito<cytoplasm "Mitochondria"
cell2

```

would be translated to

```

<listOfCompartments>
  <compartment id="cell" size="1"/>
  <compartment id="cytoplasm" size="0.8" outside="cell"/>
  <compartment id="nucleus" size="0.1" outside="cell"/>
  <compartment id="mito" name="Mitochondria"
                                outside="cytoplasm"/>
  <compartment id="cell2"/>
</listOfCompartments>

```

Note that if a name attribute is to be specified, it should be specified at the end of the line in double quotes. This is true for other SBML elements too.

2.5.5 Species

The following shorthand

```

@species
cell:Gene = 10b "The Gene"
cell:P2=0
cell:S1=100 s
cell:[S2]=20 sc
cell:[S3]=1000 bc
mito:S4=0 b

```

would be translated to

```

<listOfSpecies>
  <species id="Gene" name="The Gene" compartment="cell"
    initialAmount="10" boundaryCondition="true"/>
  <species id="P2" compartment="cell" initialAmount="0"/>
  <species id="S1" compartment="cell" initialAmount="100"
    hasOnlySubstanceUnits="true"/>
  <species id="S2" compartment="cell"
    initialConcentration="20"
    hasOnlySubstanceUnits="true" constant="true"/>
  <species id="S3" compartment="cell"
    initialConcentration="1000"
    boundaryCondition="true" constant="true"/>
  <species id="S4" compartment="mito" initialAmount="0"
    boundaryCondition="true"/>
</listOfSpecies>

```

Compartments are compulsory. An `initialConcentration` (as opposed to an `initialAmount`) is flagged by enclosing the species id in brackets. The boolean attributes `hasOnlySubstanceUnits`, `boundaryCondition`, and `constant` can be set to true by appending the letters s, b, and c, respectively. The order of the flags is not important.

2.5.6 Parameters

The section

```

@parameters
  k1=1
  k2=10

```

would be translated to

```

<listOfParameters>
  <parameter name="k1" value="1"/>
  <parameter name="k2" value="10"/>
</listOfParameters>

```

2.5.7 Reactions

Each reaction is specified by exactly two or three lines of text. The first line declares the reaction name and whether the reaction is reversible (`@rr=` for reversible and `@r=` otherwise). The second line specifies the reaction itself using a fairly standard notation. The (optional) third line specifies the full rate law for the kinetics. If local parameters are used, they should be declared on the same line in a comma-separated list (separated from the rate law using a `:`)

So, for example,

```

@reactions
@r=RepressionBinding "Repression Binding"
  Gene + 2P -> P2Gene

```



```

k2*Gene
@rr=Reverse
P2Gene -> Gene+2P
k1r*P2Gene : k1r=1,k2=3
@r=NoKL
Harry->Jim
@r=Test
Fred -> Fred2
k4*Fred : k4=1

```

would translate to

```

<listOfReactions>
  <reaction id="RepressionBinding" name="Repression Binding"
    reversible="false">
    <listOfReactants>
      <speciesReference species="Gene"/>
      <speciesReference species="P" stoichiometry="2"/>
    </listOfReactants>
    <listOfProducts>
      <speciesReference species="P2Gene"/>
    </listOfProducts>
    <kineticLaw>
      <math xmlns="http://www.w3.org/1998/Math/MathML">
        <apply>
          <times/>
          <ci> k2 </ci>
          <ci> Gene </ci>
        </apply>
      </math>
    </kineticLaw>
  </reaction>
  <reaction id="Reverse">
    <listOfReactants>
      <speciesReference species="P2Gene"/>
    </listOfReactants>
    <listOfProducts>
      <speciesReference species="Gene"/>
      <speciesReference species="P" stoichiometry="2"/>
    </listOfProducts>
    <kineticLaw>
      <math xmlns="http://www.w3.org/1998/Math/MathML">
        <apply>
          <times/>
          <ci> k1r </ci>
          <ci> P2Gene </ci>
        </apply>
      </math>
    <listOfParameters>
      <parameter id="k1r" value="1"/>
      <parameter id="k2" value="3"/>
    </listOfParameters>
  </reaction>
</listOfReactions>

```

```

    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="NoKL" reversible="false">
  <listOfReactants>
    <speciesReference species="Harry"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="Jim"/>
  </listOfProducts>
</reaction>
<reaction id="Test" reversible="false">
  <listOfReactants>
    <speciesReference species="Fred"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="Fred2"/>
  </listOfProducts>
<kineticLaw>
  <math xmlns="http://www.w3.org/1998/Math/MathML">
    <apply>
      <times/>
      <ci> k4 </ci>
      <ci> Fred </ci>
    </apply>
  </math>
  <listOfParameters>
    <parameter id="k4" value="1"/>
  </listOfParameters>
</kineticLaw>
</reaction>
</listOfReactions>

```

2.5.8 Example

The auto-regulatory network whose SBML is given in Appendix A.1 can be represented in SBML-shorthand in the following way.

```

@model:2.1.1=AutoRegulatoryNetwork "Auto-regulatory network"
@units
  substance=item
@compartments
  Cell
@species
  Cell:Gene=10 s
  Cell:P2Gene=0 s "P2.Gene"
  Cell:Rna=0 s
  Cell:P=0 s
  Cell:P2=0 s
@reactions

```

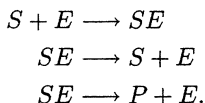
```

@r=RepressionBinding "Repression binding"
  Gene+P2 -> P2Gene
  k1*Gene*P2 : k1=1
@r=ReverseRepressionBinding "Reverse repression binding"
  P2Gene -> Gene+P2
  k1r*P2Gene : k1r=10
@r=Transcription
  Gene -> Gene+Rna
  k2*Gene : k2=0.01
@r=Translation
  Rna -> Rna+P
  k3*Rna : k3=10
@r=Dimerisation
  2P -> P2
  k4*0.5*P*(P-1) : k4=1
@r=Dissociation
  P2 -> 2P
  k4r*P2 : k4r=1
@r=RnaDegradation "RNA Degradation"
  Rna ->
  k5*Rna : k5=0.1
@r=ProteinDegradation "Protein degradation"
  P ->
  k6*P : k6=0.01

```

2.6 Exercises

1. Go to the book's website^{¶¶} and follow the Chapter 2 links to explore the world of SBML and SBML-aware software tools. In particular, download and read the SBML Level 1 and Level 2 specification documents.
2. Consider this simple model of Michaelis-Menten enzyme kinetics



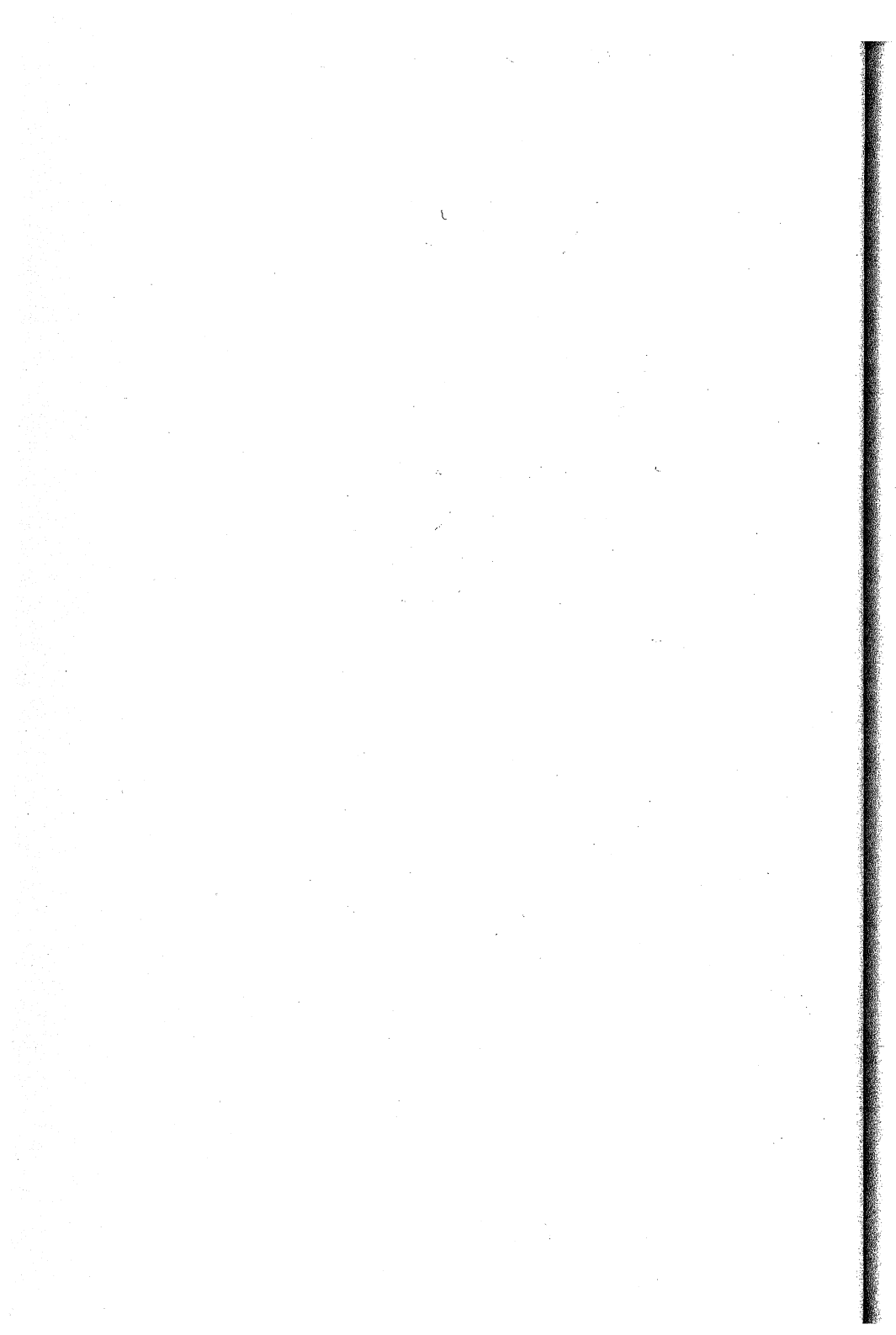
- (a) Represent this reaction network graphically using a Petri net style diagram.
- (b) Represent it mathematically as a Petri net, $N = (P, T, Pre, Post, M)$, assuming that there are currently 100 molecules of substrate S , 20 molecules of the enzyme E , and no molecules of the substrate-enzyme complex SE or the product P .
- (c) Calculate the reaction and stoichiometry matrices A and S .
- (d) If the first reaction occurs 20 times, the second 10 and the last 5, what will be the new state of the system?
- (e) Can you find a different set of transitions that will lead to the same state?

^{¶¶} URL: <http://www.staff.ncl.ac.uk/d.j.wilkinson/smf/sb/>

- (f) Can you identify any P - or T -invariants for this system?
- (g) Write the model using SBML-shorthand (do not attempt to specify any kinetic laws yet).
- (h) Hand-translate the SBML-shorthand into SBML, then validate it using the on-line validation tool at the the SBML.org website.
- (i) Download and install the SBML-shorthand compiler and use it to translate SBML-shorthand into SBML.
- (j) Download and install some SBML-aware model-building tools. Try loading your valid SBML model into the tools, and also try building it from scratch (at the time of writing, JDesigner, Jigcell, Cellware, and Cell Designer are all popular tools — there should be links to each of these from the SBML.org website).

2.7 Further reading

Murata (1989) provides a good tutorial introduction to the general theory of P/T Petri nets. Information on SBML and all things related can be found at the SBML.org website. In particular, background papers on SBML, the various SBML specifications, SBML models, tools for model validation and visualisation, and other software supporting SBML can all be found there. Pinney et al. (2003) and Hardy & Robillard (2004) give introductions to the use of Petri nets in systems biology and Goss & Peccoud (1998) explore the use of Stochastic Petri Nets (SPNs) for discrete-event simulation of stochastic kinetic models.



Probability models

3.1 Probability

3.1.1 Sample spaces, events, and sets

The models and representations considered in the previous chapter provide a framework for thinking about the state of a biochemical network, the reactions that can take place, and the change in state that occurs as a result of particular chemical reactions. As yet, however, little has been said about *which* reactions are likely to occur or *when*. The state of a biochemical network evolves continuously through time with discrete changes in state occurring as the result of reaction events. These reaction events are *random*, governed by *probabilistic laws*. It is therefore necessary to have a fairly good background in probability theory in order to properly understand these processes. In a short text such as this it is impossible to provide complete coverage of all of the necessary material. However, this chapter is meant to provide a quick summary of the essential concepts in a form that should be accessible to anyone with a high school mathematics education who has ever studied some basic probability and statistics. Readers with a strong background in probability will want to skip through this chapter. Note, however, that particular emphasis is placed on the properties of the exponential distribution, as these turn out to be central to understanding the various stochastic simulation algorithms that will be examined in detail in later chapters.

Any readers finding this chapter difficult should go back to a classic introductory text such as Ross (2003). The material in this chapter should provide sufficient background for the next few chapters (concerned with stochastic processes and simulation of biochemical networks). However, it does not cover sufficient statistical theory for the later chapters concerned with inference from data. Suitable additional reading matter for those chapters will be discussed at an appropriate point in the text.

Probability theory is used as a model for situations for which the outcomes occur randomly. Generically, such situations are called *experiments*, and the set of all possible outcomes of the experiment is known as the *sample space* corresponding to an experiment. The sample space is usually denoted by S , and a generic element of the sample space (a possible outcome) is denoted by s . The sample space is chosen so that exactly one outcome will occur. The size of the sample space is *finite*, *countably infinite*, or *uncountably infinite*.

Definition 3.1 A subset of the sample space (a collection of possible outcomes) is known as an event. We write $E \subset S$ if E is a subset of S . Events may be classified into four types:

the null event is the empty subset of the sample space;

an atomic event is a subset consisting of a single element of the sample space;

a compound event is a subset consisting of more than one element of the sample space;

the sample space itself is also an event.

Definition 3.2

The union of two events E and F (written $E \cup F$) is the event that at least one of E and F occurs. The union of the events can be obtained by forming the union of the sets.

The intersection of two events E and F (written $E \cap F$) is the event that both E and F occur. The intersection of two events can be obtained by forming the intersection of the sets.

The complement of an event, A , denoted A^c or \bar{A} , is the event that A does not occur, and hence consists of all those elements of the sample space that are not in A .

Two events A and B are disjoint or mutually exclusive if they cannot both occur. That is, their intersection is empty

$$A \cap B = \emptyset.$$

Note that for any event A , the events A and A^c are disjoint, and their union is the whole of the sample space:

$$A \cap A^c = \emptyset \quad \text{and} \quad A \cup A^c = S.$$

The event A is true if the outcome of the experiment, s , is contained in the event A ; that is, if $s \in A$. We say that the event A implies the event B , and write $A \Rightarrow B$, if the truth of B automatically follows from the truth of A . If A is a subset of B , then occurrence of A necessarily implies occurrence of the event B . That is

$$(A \subseteq B) \iff (A \cap B = A) \iff (A \Rightarrow B).$$

In order to carry out more sophisticated manipulations of events, it is helpful to know some basic rules of set theory.

Proposition 3.1

Commutative laws:

$$A \cup B = B \cup A$$

$$A \cap B = B \cap A$$

Associative laws:

$$(A \cup B) \cup C = A \cup (B \cup C)$$

$$(A \cap B) \cap C = A \cap (B \cap C)$$

Distributive laws:

$$(A \cup B) \cap C = (A \cap C) \cup (B \cap C)$$

$$(A \cap B) \cup C = (A \cup C) \cap (B \cup C)$$

DeMorgan's laws:

$$(A \cup B)^c = A^c \cap B^c$$

$$(A \cap B)^c = A^c \cup B^c$$

Disjoint union:

$$A \cup B = (A \cap B^c) \cup (A^c \cap B) \cup (A \cap B)$$

and $A \cap B^c$, $A^c \cap B$, and $A \cap B$ are disjoint.

3.1.2 Probability axioms

Once a suitable mathematical framework for understanding events in terms of sets has been established, it is possible to construct a corresponding framework for understanding probabilities of events in terms of sets.

Definition 3.3 *The real valued function $P(\cdot)$ is a probability measure if it acts on subsets of S and obeys the following axioms:*

I. $P(S) = 1$.

II. If $A \subseteq S$ then $P(A) \geq 0$.

III. If A and B are disjoint ($A \cap B = \emptyset$) then

$$P(A \cup B) = P(A) + P(B).$$

Repeated use of Axiom III gives the more general result that if A_1, A_2, \dots, A_n are mutually disjoint, then

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i).$$

Indeed, we will assume further that the above result holds even if we have a countable infinite collection of disjoint events ($n = \infty$).

These axioms seem to fit well with most people's intuitive understanding of probability, but there are a few additional comments worth making.

1. Axiom I says that one of the possible outcomes must occur. A probability of 1 is assigned to the event "something occurs." This fits in exactly with the definition of sample space. Note, however, that the implication does not go the other way! When dealing with infinite sample spaces, there are often events of probability 1 which are not the sample space, and events of probability zero, which are not the empty set.
2. Axiom II simply states that we wish to work only with positive probabilities, because in some sense, probability measures the *size* of the set (event).

3. Axiom III says that probabilities “add up” in a natural way. Allowing this result to hold for countably infinite unions is slightly controversial, but it makes the mathematics much easier, so it will be assumed throughout this text.

These axioms are (almost) all that is needed in order to develop a theory of probability, but there are a collection of commonly used properties which follow directly from these axioms and are used extensively when carrying out probability calculations.

Proposition 3.2

1. $P(A^c) = 1 - P(A)$
2. $P(\emptyset) = 0$
3. If $A \subseteq B$, then $P(A) \leq P(B)$
4. (Addition Law) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

Proof. For 1, since $A \cap A^c = \emptyset$, and $S = A \cup A^c$, Axiom III tells us that $P(S) = P(A \cup A^c) = P(A) + P(A^c)$. From axiom I we know that $P(S) = 1$. Re-arranging gives the result. Property 2 follows from property 1 as $\emptyset = S^c$. It simply says that the probability of no outcome is zero, which again fits in with our definition of a sample space. For property 3 write B as

$$B = A \cup (B \cap A^c),$$

where A and $B \cap A^c$ are disjoint. Then, from the third axiom,

$$P(B) = P(A) + P(B \cap A^c)$$

so that

$$P(A) = P(B) - P(B \cap A^c) \leq P(B).$$

Property 4 is one of the exercises at the end of the chapter. \square

3.1.3 Interpretations of probability

Most people have an intuitive feel for the notion of probability, and the axioms seem to capture its essence in a mathematical form. However, for probability theory to be anything other than an interesting piece of abstract pure mathematics, it must have an interpretation that in some way connects it to reality. If you wish only to study probability as a mathematical theory, there is no need to have an interpretation. However, if you are to use probability theory as your foundation for a theory which makes probabilistic statements about the world around us, then there must be an interpretation of probability which makes some connection between the mathematical theory and reality.

While there is (almost) unanimous agreement about the mathematics of probability, the axioms, and their consequences, there is considerable disagreement about the interpretation of probability. The three most common interpretations are given below.

Classical interpretation

The classical interpretation of probability is based on the assumption of underlying equally likely events. That is, for any events under consideration, there is always a sample space which can be considered where all atomic events are equally likely. If this sample space is given, then the probability axioms may be deduced from set-theoretic considerations.

This interpretation is fine when it is obvious how to partition the sample space into equally likely events and is in fact entirely compatible with the other two interpretations to be described in that case. The problem with this interpretation is that for many situations it is not at all obvious what the partition into equally likely events is. For example, consider the probability that it will rain in a particular location tomorrow. This is clearly a reasonable event to consider, but it is not at all clear what sample space we should construct with equally likely outcomes. Consequently, the classical interpretation falls short of being a good interpretation for real-world problems. However, it provides a good starting point for a mathematical treatment of probability theory and is the interpretation adopted by many mathematicians and theoreticians.

Frequentist interpretation

An interpretation of probability widely adopted by statisticians is the relative frequency interpretation. This interpretation makes a much stronger connection with reality than the previous one and fits in well with traditional statistical methodology. Here probability only has meaning for events from experiments which could in principle be repeated arbitrarily many times under essentially identical conditions. Here, the probability of an event is simply the "long-run proportion" of times that the event occurs under many repetitions of the experiment. It is reasonable to suppose that this proportion will settle down to some limiting value eventually, which is the probability of the event. In such a situation, it is possible to derive the axioms of probability from consideration of the long run frequencies of various events. The probability p , of an event E , is defined by

$$p = \lim_{n \rightarrow \infty} \frac{r}{n}$$

where r is the number of times E occurred in n repetitions of the experiment.

Unfortunately it is hard to know precisely why such a limiting frequency should exist. A bigger problem, however, is that the interpretation only applies to outcomes of repeatable experiments, and there are many "one-off" events, such as "rain here tomorrow," on which we would like to be able to attach probabilities.

Subjective interpretation

This final common interpretation of probability is somewhat controversial, but does not suffer from the problems the other interpretations do. It suggests that the association of probabilities to events is a personal (subjective) process, relating to your *degree of belief* in the likelihood of the event occurring. It is controversial because it accepts that *different* people will assign *different* probabilities to the *same event*.

While in some sense it gives up on an objective notion of probability, it is in no sense arbitrary. It can be defined in a precise way, from which the axioms of probability may be derived as requirements of self-consistency.

A simple way to define *your* subjective probability that some event E will occur is as follows. Your probability is the number p such that you consider $\pounds p$ to be a *fair price* for a gamble which will pay you $\pounds 1$ if E occurs and nothing otherwise.

So, if you consider 40 pence to be a fair price for a gamble which pays you $\pounds 1$ if it rains tomorrow, then 0.4 is your subjective probability for the event.* The subjective interpretation is sometimes known as the *degree of belief interpretation*, which is the interpretation of probability underlying the theory of *Bayesian Statistics* — a powerful theory of statistical inference named after Thomas Bayes, the 18th-century Presbyterian minister who first proposed it. Consequently, this interpretation of probability is sometimes also known as the *Bayesian interpretation*.

Summary

While the interpretation of probability is philosophically very important, all interpretations lead to the same set of axioms, from which the rest of probability theory is deduced. Consequently, for much of this text, it will be sufficient to adopt a fairly classical approach, taking the axioms as given and investigating their consequences independently of the precise interpretation adopted. However, the inferential theory considered in the later chapters is distinctly Bayesian in nature, and Bayesian inference is most naturally associated with a subjective interpretation of probability.

3.1.4 Classical probability

Classical probability theory is concerned with carrying out probability calculations based on *equally likely outcomes*. That is, it is assumed that the sample space has been constructed in such a way that every subset of the sample space consisting of a single element has the same probability. If the sample space contains n possible outcomes ($\#S = n$), we must have for all $s \in S$,

$$P(\{s\}) = \frac{1}{n}$$

and hence for all $E \subseteq S$

$$P(E) = \frac{\#E}{n}.$$

More informally, we have

$$P(E) = \frac{\text{number of ways } E \text{ can occur}}{\text{total number of outcomes}}.$$

* For non-British readers, replace \pounds with $\$$ and pence with cents.

Example

Suppose that a fair coin is thrown twice, and the results are recorded. The sample space is

$$S = \{HH, HT, TH, TT\}.$$

Let us assume that each outcome is equally likely — that is, each outcome has a probability of $1/4$. Let A denote the event *head on the first toss*, and B denote the event *head on the second toss*. In terms of sets

$$A = \{HH, HT\}, \quad B = \{HH, TH\}.$$

So

$$P(A) = \frac{\#A}{n} = \frac{2}{4} = \frac{1}{2}$$

and similarly $P(B) = 1/2$. If we are interested in the event $C = A \cup B$ we can work out its probability from the set definition as

$$P(C) = \frac{\#C}{4} = \frac{\#(A \cup B)}{4} = \frac{\#\{HH, HT, TH\}}{4} = \frac{3}{4}$$

or by using the addition formula

$$P(C) = P(A \cup B) = P(A) + P(B) - P(A \cap B) = \frac{1}{2} + \frac{1}{2} - P(A \cap B).$$

Now $A \cap B = \{HH\}$, which has probability $1/4$, so

$$P(C) = \frac{1}{2} + \frac{1}{2} - \frac{1}{4} = \frac{3}{4}.$$

In this simple example, it seems easier to work directly with the definition. However, in more complex problems, it is usually much easier to work out how many elements there are in an intersection than in a union, making the addition law very useful.

The multiplication principle

In the above example we saw that there were two distinct experiments — *first throw* and *second throw*. There were two equally likely outcomes for the first throw and two equally likely outcomes for the second throw. This leads to a combined experiment with $2 \times 2 = 4$ possible outcomes. This is an example of the *multiplication principle*.

Proposition 3.3 *If there are p experiments and the first has n_1 equally likely outcomes, the second has n_2 equally likely outcomes, and so on until the p th experiment has n_p equally likely outcomes, then there are*

$$n_1 \times n_2 \times \cdots \times n_p = \prod_{i=1}^p n_i$$

equally likely possible outcomes for the p experiments.

3.1.5 Conditional probability and the multiplication rule

We now have a way of understanding the probabilities of events, but so far we have no way of *modifying* those probabilities when certain events occur. For this, we need an extra axiom which can be justified under any of the interpretations of probability.

Definition 3.4 The conditional probability of A given B , written $P(A|B)$ is defined by

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad \text{for } P(B) > 0.$$

Note that we can only condition on events with positive probability.

Under the classical interpretation of probability, we can see that if we are told that B has occurred, then all outcomes in B are equally likely, and all outcomes not in B have zero probability — so B is the new sample space. The number of ways that A can occur is now just the number of ways $A \cap B$ can occur, and these are all equally likely. Consequently we have

$$P(A|B) = \frac{\#(A \cap B)}{\#B} = \frac{\#(A \cap B)/\#S}{\#B/\#S} = \frac{P(A \cap B)}{P(B)}.$$

Because conditional probabilities really just correspond to a new probability measure defined on a smaller sample space, they obey all of the properties of “ordinary” probabilities. For example, we have

$$\begin{aligned} P(B|B) &= 1 \\ P(\emptyset|B) &= 0 \\ P(A \cup C|B) &= P(A|B) + P(C|B), \quad \text{for } A \cap C = \emptyset \end{aligned}$$

and so on.

The definition of conditional probability simplifies when one event is a special case of the other. If $A \subseteq B$, then $A \cap B = A$ so

$$P(A|B) = \frac{P(A)}{P(B)}, \quad \text{for } A \subseteq B.$$

Example

A die is rolled and the number showing recorded. Given that the number rolled was even, what is the probability that it was a six?

Let E denote the event “even” and F denote the event “a six.” Clearly $F \subseteq E$, so

$$P(F|E) = \frac{P(F)}{P(E)} = \frac{1/6}{1/2} = \frac{1}{3}.$$

The formula for conditional probability is useful when we want to calculate $P(A|B)$ from $P(A \cap B)$ and $P(B)$. However, more commonly we want to know $P(A \cap B)$ and we know $P(A|B)$ and $P(B)$. A simple rearrangement gives us the multiplication rule.

Proposition 3.4 (multiplication rule)

$$P(A \cap B) = P(B) \times P(A|B)$$

Example

Two cards are dealt from a deck of 52 cards. What is the probability that they are both Aces?

Let A_1 be the event “first card an Ace” and A_2 be the event “second card an Ace.” $P(A_2|A_1)$ is the probability of a second Ace. Given that the first card has been drawn and was an Ace, there are 51 cards left, 3 of which are Aces, so $P(A_2|A_1) = 3/51$. So,

$$\begin{aligned} P(A_1 \cap A_2) &= P(A_1) \times P(A_2|A_1) \\ &= \frac{4}{52} \times \frac{3}{51} \\ &= \frac{1}{221}. \end{aligned}$$

The multiplication rule generalises to more than two events. For example, for three events we have

$$P(A_1 \cap A_2 \cap A_3) = P(A_1) P(A_2|A_1) P(A_3|A_1 \cap A_2).$$

3.1.6 Independent events, partitions and Bayes Theorem

Recall the multiplication rule (Proposition 3.4):

$$P(A \cap B) = P(B) P(A|B).$$

For some events A and B , knowing that B has occurred will not alter the probability of A , so that $P(A|B) = P(A)$. When this is so, the multiplication rule becomes

$$P(A \cap B) = P(A) P(B),$$

and the events A and B are said to be *independent events*. Independence is a very important concept in probability theory, and it is often used to build up complex events from simple ones. Do not confuse the independence of A and B with the exclusivity of A and B — they are entirely different concepts. If A and B both have positive probability, then they cannot be both independent and exclusive (this is an end-of-chapter exercise).

When it is clear that the occurrence of B can have no influence on A , we will *assume* independence in order to calculate $P(A \cap B)$. However, if we can calculate $P(A \cap B)$ directly, we can check the independence of A and B by seeing if it is true that

$$P(A \cap B) = P(A) P(B).$$

We can generalise independence to collections of events as follows.

Definition 3.5 *The set of events $A = \{A_1, A_2, \dots, A_n\}$ is mutually independent if for any subset, $B \subseteq A$, $B = \{B_1, B_2, \dots, B_r\}$, $r \leq n$ we have*

$$P(B_1 \cap \dots \cap B_r) = P(B_1) \times \dots \times P(B_r).$$

Note that mutual independence is much stronger than *pair-wise* independence, where

we only require independence of subsets of size $r = 2$. That is, pair-wise independence *does not* imply mutual independence.

Definition 3.6 A partition of a sample space is simply the decomposition of the sample space into a collection of mutually exclusive events with positive probability. That is, $\{B_1, \dots, B_n\}$ forms a partition of S if

- $S = B_1 \cup B_2 \cup \dots \cup B_n = \bigcup_{i=1}^n B_i$,
- $B_i \cap B_j = \emptyset, \forall i \neq j$,
- $P(B_i) > 0, \forall i$.

Theorem 3.1 (Theorem of total probability) Suppose that we have a partition $\{B_1, \dots, B_n\}$ of a sample space, S . Suppose further that we have an event A . Then

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

Proof. A can be written as the disjoint union

$$A = (A \cap B_1) \cup \dots \cup (A \cap B_n),$$

and so the probability of A is given by

$$\begin{aligned} P(A) &= P((A \cap B_1) \cup \dots \cup (A \cap B_n)) \\ &= P(A \cap B_1) + \dots + P(A \cap B_n) && \text{(by Axiom III)} \\ &= P(A|B_1)P(B_1) + \dots + P(A|B_n)P(B_n) && \text{(by Proposition 3.4)} \\ &= \sum_{i=1}^n P(A|B_i)P(B_i). \end{aligned}$$

□

Theorem 3.2 (Bayes Theorem) For all events A, B such that $P(B) > 0$ we have

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

This is a very important result as it tells us how to “turn conditional probabilities around” — that is, it tells us how to work out $P(A|B)$ from $P(B|A)$, and this is often very useful.

Proof. By Definition 3.4 we have

$$\begin{aligned} P(A|B) &= \frac{P(A \cap B)}{P(B)} \\ &= \frac{P(A)P(B|A)}{P(B)}. && \text{(by Proposition 3.4)} \end{aligned}$$

□

Example

A clinic offers you a free test for a very rare, but hideous disease. The test they offer is very reliable. If you have the disease it has a 98% chance of giving a positive result, and if you do not have the disease, it has only a 1% chance of giving a positive result. Despite having no *a priori* reason to suppose that you have the disease, you nevertheless decide to take the test and find that you test positive — what is the probability that you have the disease?

Let P be the event “test positive” and D be the event “you have the disease.” We know that

$$P(P|D) = 0.98 \text{ and that } P(P|D^c) = 0.01.$$

We want to know $P(D|P)$, so we use Bayes Theorem.

$$\begin{aligned} P(D|P) &= \frac{P(P|D)P(D)}{P(P)} \\ &= \frac{P(P|D)P(D)}{P(P|D)P(D) + P(P|D^c)P(D^c)} && \text{(using Theorem 3.1)} \\ &= \frac{0.98P(D)}{0.98P(D) + 0.01(1 - P(D))}. \end{aligned}$$

So we see that the probability you have the disease given the test result depends on the probability that you had the disease in the first place. This is a rare disease, affecting only one in ten thousand people, so that $P(D) = 0.0001$. Substituting this in gives

$$P(D|P) = \frac{0.98 \times 0.0001}{0.98 \times 0.0001 + 0.01 \times 0.9999} \simeq 0.01.$$

So, your probability of having the disease has increased from 1 in 10,000 to 1 in 100, but still is not that much to get worried about! Note the *crucial* difference between $P(P|D)$ and $P(D|P)$.

Another important thing to notice about the above example is the use of the theorem of total probability in order to expand the bottom line of Bayes Theorem. In fact, this is done so often that Bayes Theorem is often stated in this form.

Corollary 3.1 *Suppose that we have a partition $\{B_1, \dots, B_n\}$ of a sample space S . Suppose further that we have an event A , with $P(A) > 0$. Then, for each B_j , the probability of B_j given A is*

$$\begin{aligned} P(B_j|A) &= \frac{P(A|B_j)P(B_j)}{P(A)} \\ &= \frac{P(A|B_j)P(B_j)}{P(A|B_1)P(B_1) + \dots + P(A|B_n)P(B_n)} \\ &= \frac{P(A|B_j)P(B_j)}{\sum_{i=1}^n P(A|B_i)P(B_i)}. \end{aligned}$$

In particular, if the partition is simply $\{B, B^c\}$, then this simplifies to

$$P(B|A) = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|B^c)P(B^c)}.$$

3.2 Discrete probability models

3.2.1 Introduction, mass functions, and distribution functions

We have seen how to relate events to sets and how to calculate probabilities for events by working with the sets that represent them. So far, however, we have not developed any special techniques for thinking about *random quantities*. *Discrete probability models* provide a framework for thinking about *discrete random quantities*, and *continuous probability models* (to be considered in the next section) form a framework for thinking about *continuous random quantities*.

Example

Consider the sample space for tossing a fair coin twice:

$$S = \{HH, HT, TH, TT\}.$$

These outcomes are equally likely. There are several random quantities we could associate with this experiment. For example, we could count the number of heads or the number of tails.

Definition 3.7 A random quantity is a real valued function which acts on elements of the sample space (outcomes). That is, to each outcome, the random variable assigns a real number.

Random quantities (sometimes known as *random variables*) are always denoted by upper case letters.

In our example, if we let X be the number of heads, we have

$$\begin{aligned} X(HH) &= 2, \\ X(HT) &= 1, \\ X(TH) &= 1, \\ X(TT) &= 0. \end{aligned}$$

The observed value of a random quantity is the number corresponding to the actual outcome. That is, if the outcome of an experiment is $s \in S$, then $X(s) \in \mathbb{R}$ is the observed value. This observed value is always denoted with a lower case letter — here x . Thus $X = x$ means that the observed value of the random quantity X is the number x . The set of possible observed values for X is

$$S_X = \{X(s) | s \in S\}.$$

For the above example we have

$$S_X = \{0, 1, 2\}.$$

Clearly here the values are not all equally likely.

Example

Roll one die and call the random number which is uppermost Y . The sample space for the *random quantity* Y is

$$S_Y = \{1, 2, 3, 4, 5, 6\}$$

and these outcomes are all equally likely. Now roll two dice and call their sum Z . The sample space for Z is

$$S_Z = \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$$

and these outcomes are *not* equally likely. However, we know the probabilities of the events corresponding to each of these outcomes, and we could display them in a table as follows.

Outcome	2	3	4	5	6	7	8	9	10	11	12
Probability	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36

This is essentially a tabulation of the *probability mass function* for the random quantity Z .

Definition 3.8 (probability mass function) For any discrete random variable X , we define the probability mass function (PMF) to be the function which gives the probability of each $x \in S_X$. Clearly we have

$$P(X = x) = \sum_{\{s \in S | X(s) = x\}} P(\{s\}).$$

That is, the probability of getting a particular number is the sum of the probabilities of all those outcomes which have that number associated with them. Also $P(X = x) \geq 0$ for each $x \in S_X$, and $P(X = x) = 0$ otherwise.

Definition 3.9 The set of all pairs $\{(x, P(X = x)) | x \in S_X\}$ is known as the probability distribution of X .

Example

For the example above concerning the sum of two dice, the probability distribution is

$$\{(2, 1/36), (3, 2/36), (4, 3/36), (5, 4/36), (6, 5/36), (7, 6/36), (8, 5/36), (9, 4/36), (10, 3/36), (11, 2/36), (12, 1/36)\}$$

and the probability mass function can be tabulated as follows.

x	2	3	4	5	6	7	8	9	10	11	12
$P(X = x)$	1/36	2/36	3/36	4/36	5/36	6/36	5/36	4/36	3/36	2/36	1/36

For any discrete random quantity, X , we clearly have

$$\sum_{x \in S_X} P(X = x) = 1$$

as every outcome has some number associated with it. It can often be useful to know the probability that your random number is no greater than some particular value.

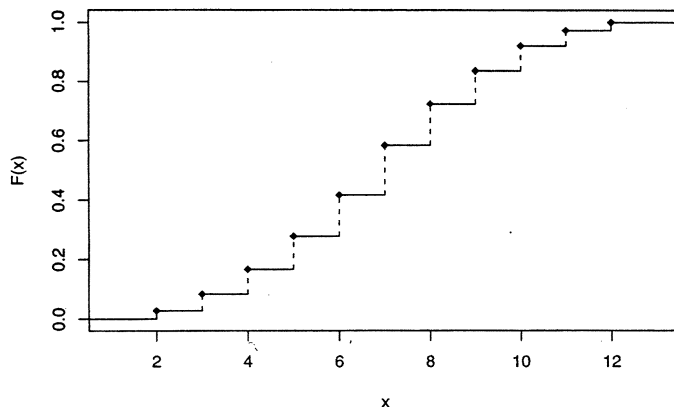


Figure 3.1 CDF for the sum of a pair of fair dice

Definition 3.10 (cumulative distribution function) The cumulative distribution function (CDF), is defined by

$$F_X(x) = P(X \leq x) = \sum_{\{y \in S_X | y \leq x\}} P(X = y).$$

Example

For the sum of two dice, the CDF can be tabulated for the outcomes as

x	2	3	4	5	6	7	8	9	10	11	12
$F_X(x)$	1/36	3/36	6/36	10/36	15/36	21/36	26/36	30/36	33/36	35/36	36/36

but it is important to note that the CDF is defined for all real numbers — not just the possible values. In our example we have

$$F_X(-3) = P(X \leq -3) = 0,$$

$$F_X(4.5) = P(X \leq 4.5) = P(X \leq 4) = 6/36,$$

$$F_X(25) = P(X \leq 25) = 1.$$

We may plot the CDF for our example as shown in Figure 3.1.

It is clear that for any random variable X , for all $x \in \mathbb{R}$, $F_X(x) \in [0, 1]$ and that $F_X(x) \rightarrow 0$ as $x \rightarrow -\infty$ and $F_X(x) \rightarrow 1$ as $x \rightarrow +\infty$.

3.2.2 Expectation and variance for discrete random quantities

It is useful to be able to summarise the distribution of random quantities. A location measure often used to summarise random quantities is the expectation of the random quantity. It is the “centre of mass” of the probability distribution.

Definition 3.11 The expectation of a discrete random quantity X , written $E(X)$ is defined by

$$E(X) = \sum_{x \in S_X} x P(X = x).$$

The expectation is often denoted by μ_X or even just μ . Note that the expectation is a known function of the probability distribution. It is *not* a random quantity, and it is *not* the sample mean of a set of data (random or otherwise). In some sense (to be made precise in Proposition 3.14), it represents the value of the random variable that you expect to get on average.

Example

For the sum of two dice, X , we have

$$E(X) = 2 \times \frac{1}{36} + 3 \times \frac{2}{36} + 4 \times \frac{3}{36} + \cdots + 12 \times \frac{1}{36} = 7.$$

By looking at the symmetry of the mass function, it is clear that in some sense 7 is the “central” value of the probability distribution.

As well as a method for summarising the location of a given probability distribution, it is also helpful to have a summary for the *spread*.

Definition 3.12 For a discrete random quantity X , the variance of X is defined by

$$\text{Var}(X) = \sum_{x \in S_X} \{(x - E(X))^2 P(X = x)\}.$$

The variance is often denoted σ_X^2 , or even just σ^2 . Again, this is a known function of the probability distribution. It is not random, and it is not the *sample* variance of a set of data. The variance can be rewritten as

$$\text{Var}(X) = \sum_{x_i \in S_X} x_i^2 P(X = x_i) - [E(X)]^2,$$

and this expression is usually a bit easier to work with. We also define the *standard deviation* of a random quantity by

$$\text{SD}(X) = \sqrt{\text{Var}(X)},$$

and this is usually denoted by σ_X or just σ . The variance represents the average squared distance of X from $E(X)$. The units of variance are therefore the square of the units of X (and $E(X)$). Some people prefer to work with the standard deviation because it has the same units as X and $E(X)$.

Example

For the sum of two dice, X , we have

$$\sum_{x_i \in S_X} x_i^2 P(X = x_i) = 2^2 \times \frac{1}{36} + 3^2 \times \frac{2}{36} + 4^2 \times \frac{3}{36} + \cdots + 12^2 \times \frac{1}{36} = \frac{329}{6}$$

and so

$$\text{Var}(X) = \frac{329}{6} - 7^2 = \frac{35}{6}, \quad \text{and} \quad \text{SD}(X) = \sqrt{\frac{35}{6}}.$$

3.2.3 Properties of expectation and variance

One of the reasons that expectation is widely used as a measure of location for probability distributions is the fact that it has many desirable mathematical properties which make it elegant and convenient to work with. Indeed, many of the nice properties of expectation lead to corresponding nice properties for variance, which is one of the reasons why variance is widely used as a measure of spread.

Suppose that X is a discrete random quantity, and that Y is another random quantity that is a known function of X . That is, $Y = g(X)$ for some function $g(\cdot)$. What is the expectation of Y ?

Example

Throw a die, and let X be the number showing. We have

$$S_X = \{1, 2, 3, 4, 5, 6\}$$

and each value is equally likely. Now suppose that we are actually interested in the square of the number showing. Define a new random quantity $Y = X^2$. Then

$$S_Y = \{1, 4, 9, 16, 25, 36\}$$

and clearly each of these values is equally likely. We therefore have

$$E(Y) = 1 \times \frac{1}{6} + 4 \times \frac{1}{6} + \cdots + 36 \times \frac{1}{6} = \frac{91}{6}.$$

The above example illustrates the more general result, that for $Y = g(X)$, we have

$$E(Y) = \sum_{x \in S_X} g(x) P(X = x).$$

Note that in general $E(g(X)) \neq g(E(X))$. For the above example, $E(X^2) = 91/6 \simeq 15.2$, and $E(X)^2 = 3.5^2 = 12.25$.

We can use this more general notion of expectation in order to redefine variance purely in terms of expectation as follows:

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

Having said that $E(g(X)) \neq g(E(X))$ in general, it does in fact hold in the (very) special, but important case where $g(\cdot)$ is a *linear* function.

Lemma 3.1 (expectation of a linear transformation) *If we have a random quantity X , and a linear transformation, $Y = aX + b$, where a and b are known real constants, then we have that*

$$E(aX + b) = aE(X) + b.$$

Proof.

$$\begin{aligned}
 E(aX + b) &= \sum_{x \in S_X} (ax + b) P(X = x) \\
 &= \sum_{x \in S_X} ax P(X = x) + \sum_{x \in S_X} b P(X = x) \\
 &= a \sum_{x \in S_X} x P(X = x) + b \sum_{x \in S_X} P(X = x) \\
 &= a E(X) + b.
 \end{aligned}$$

□

Lemma 3.2 (expectation of a sum) For two random quantities X and Y , the expectation of their sum is given by

$$E(X + Y) = E(X) + E(Y).$$

Note that this result is true irrespective of whether X and Y are independent. Let us see why.

Proof. First,

$$S_{X+Y} = \{x + y | (x \in S_X) \cap (y \in S_Y)\},$$

and so

$$\begin{aligned}
 E(X + Y) &= \sum_{(x+y) \in S_{X+Y}} (x + y) P((X = x) \cap (Y = y)) \\
 &= \sum_{x \in S_X} \sum_{y \in S_Y} (x + y) P((X = x) \cap (Y = y)) \\
 &= \sum_{x \in S_X} \sum_{y \in S_Y} x P((X = x) \cap (Y = y)) \\
 &\quad + \sum_{x \in S_X} \sum_{y \in S_Y} y P((X = x) \cap (Y = y)) \\
 &= \sum_{x \in S_X} \sum_{y \in S_Y} x P(X = x) P(Y = y | X = x) \\
 &\quad + \sum_{y \in S_Y} \sum_{x \in S_X} y P(Y = y) P(X = x | Y = y) \\
 &= \sum_{x \in S_X} x P(X = x) \sum_{y \in S_Y} P(Y = y | X = x) \\
 &\quad + \sum_{y \in S_Y} y P(Y = y) \sum_{x \in S_X} P(X = x | Y = y) \\
 &= \sum_{x \in S_X} x P(X = x) + \sum_{y \in S_Y} y P(Y = y) \\
 &= E(X) + E(Y).
 \end{aligned}$$

□

We can now put together the result for the expectation of a sum, and for the expectation of a linear transformation in order to give the following important property, commonly referred to as the *linearity of expectation*.

Proposition 3.5 (expectation of a linear combination) For any random quantities X_1, \dots, X_n and scalar constants a_0, \dots, a_n we have

$$\begin{aligned} E(a_0 + a_1 X_1 + a_2 X_2 + \dots + a_n X_n) &= a_0 + a_1 E(X_1) + a_2 E(X_2) + \dots \\ &\quad + a_n E(X_n). \end{aligned}$$

Lemma 3.3 (expectation of an independent product) If X and Y are independent random quantities, then

$$E(XY) = E(X)E(Y).$$

Proof.

$$S_{XY} = \{xy | (x \in S_X) \cap (y \in S_Y)\},$$

and so

$$\begin{aligned} E(XY) &= \sum_{xy \in S_{XY}} xy P((X = x) \cap (Y = y)) \\ &= \sum_{x \in S_X} \sum_{y \in S_Y} xy P(X = x) P(Y = y) \\ &= \sum_{x \in S_X} x P(X = x) \sum_{y \in S_Y} y P(Y = y) \\ &= E(X)E(Y) \end{aligned}$$

□

Note that here it is *vital* that X and Y are independent or the result does not hold.

Lemma 3.4 (variance of a linear transformation) If X is a random quantity with finite variance $\text{Var}(X)$, then

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

Proof.

$$\begin{aligned} \text{Var}(aX + b) &= E([aX + b]^2) - [E(aX + b)]^2 \\ &= E(a^2 X^2 + 2abX + b^2) - [aE(X) + b]^2 \\ &= a^2 [E(X^2) - E(X)^2] \\ &= a^2 \text{Var}(X). \end{aligned}$$

□

Lemma 3.5 (variance of an independent sum) If X and Y are independent random quantities, then

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

Proof.

$$\begin{aligned}
 \text{Var}(X + Y) &= \text{E}([X + Y]^2) - [\text{E}(X + Y)]^2 \\
 &= \text{E}(X^2 + 2XY + Y^2) - [\text{E}(X) + \text{E}(Y)]^2 \\
 &= \text{E}(X^2) + 2\text{E}(XY) + \text{E}(Y^2) - \text{E}(X)^2 - 2\text{E}(X)\text{E}(Y) \\
 &\quad - \text{E}(Y)^2 \\
 &= \text{E}(X^2) + 2\text{E}(X)\text{E}(Y) + \text{E}(Y^2) - \text{E}(X)^2 - 2\text{E}(X)\text{E}(Y) \\
 &\quad - \text{E}(Y)^2 \\
 &= \text{E}(X^2) - \text{E}(X)^2 + \text{E}(Y^2) - \text{E}(Y)^2 \\
 &= \text{Var}(X) + \text{Var}(Y)
 \end{aligned}$$

□

Again, it is vital that X and Y are independent or the result does not hold. Notice that this implies a slightly less attractive result for the standard deviation of the sum of two independent random quantities,

$$\text{SD}(X + Y) = \sqrt{\text{SD}(X)^2 + \text{SD}(Y)^2},$$

which is why it is often more convenient to work with variances.[†]

Putting together previous results we get the following proposition.

Proposition 3.6 (variance of a linear combination) *For mutually independent X_1, X_2, \dots, X_n we have*

$$\begin{aligned}
 \text{Var}(a_0 + a_1X_1 + a_2X_2 + \dots + a_nX_n) \\
 = a_1^2 \text{Var}(X_1) + a_2^2 \text{Var}(X_2) + \dots + a_n^2 \text{Var}(X_n).
 \end{aligned}$$

Before moving on to look at some interesting families of probability distributions, it is worth emphasising the link between expectation and the theorem of total probability.

Proposition 3.7 *If F is an event and X is a discrete random quantity with outcome space S_X , then by the theorem of total probability (Theorem 3.1) we have*

$$\begin{aligned}
 \text{P}(F) &= \sum_{x \in S_X} \text{P}(F|X = x) \text{P}(X = x) \\
 &= \text{E}(\text{P}(F|X)),
 \end{aligned}$$

where $\text{P}(F|X)$ is the random quantity which takes the value $\text{P}(F|X = x)$ when X takes the value x .

[†] Note the similarity to Pythagoras' Theorem for the lengths of the sides of a triangle. Viewed in the correct way, this result is seen to be a special case of Pythagoras' Theorem, as the standard deviation of a random quantity can be viewed as a length, and independence leads to orthogonality in an appropriately constructed inner-product space.

3.3 The discrete uniform distribution

The theory of discrete random quantities covered in the previous section provides a generic set of tools for studying distributions and their properties. However, there are several important “families” of discrete probability models that occur frequently and are therefore worthy of special consideration. The first is the “discrete uniform distribution,” which corresponds to a generalisation of the number obtained by rolling a single die. A random quantity X is discrete-uniform on the numbers from 1 to n , written

$$X \sim DU(n)$$

if each of the integers from 1 to n is equally likely to be the observed value of X . We therefore have outcome space

$$S_X = \{1, 2, \dots, n\}$$

and PMF

$$P(X = k) = \frac{1}{n}, \quad k \in S_X.$$

The CDF at the points in S_X is therefore clearly given by

$$P(X \leq k) = \frac{k}{n}, \quad k \in S_X.$$

It is straightforward to compute the expectation of X as

$$\begin{aligned} E(X) &= \sum_{k=1}^n k \times \frac{1}{n} \\ &= \frac{1}{n} \sum_{k=1}^n k \\ &= \frac{1}{n} \frac{n(n+1)}{2} \\ &= \frac{n+1}{2}. \end{aligned}$$

The variance can be calculated similarly to be

$$\text{Var}(X) = \frac{n^2 - 1}{12}.$$

Although this distributional family seems somewhat contrived, it is especially useful as it typically forms the starting point for the development of any stochastic simulation algorithm.

3.4 The binomial distribution

The binomial distribution is the distribution of the number of “successes” in a series of n independent “trials,” each of which results in a “success” (with probability p) or a “failure” (with probability $1 - p$). If the number of successes is X , we would write

$$X \sim B(n, p)$$

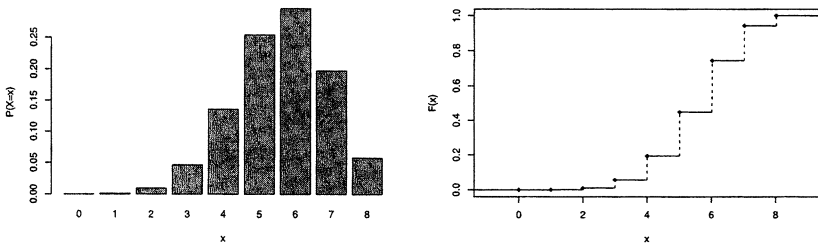


Figure 3.2 PMF and CDF for a $B(8, 0.7)$ distribution

to indicate that X is a binomial random quantity based on n independent trials, each occurring with probability p .

Let us now derive the probability mass function for $X \sim B(n, p)$. Clearly X can take on any value from 0 up to n , and no other. Therefore, we simply have to calculate $P(X = k)$ for $k = 0, 1, 2, \dots, n$. The probability of k successes followed by $n - k$ failures is clearly $p^k(1 - p)^{n-k}$. Indeed, this is the probability of any particular sequence involving k successes. There are $\binom{n}{k}$ such sequences, so by the multiplication principle, we have

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, 2, \dots, n.$$

Now, using the binomial theorem, we have

$$\sum_{k=0}^n P(X = k) = \sum_{k=0}^n \binom{n}{k} p^k (1 - p)^{n-k} = (p + [1 - p])^n = 1^n = 1,$$

and so this does define a valid probability distribution. There is no neat analytic expression for the CDF of a binomial distribution, but it is straightforward to compute and tabulate. A plot of the probability mass function and cumulative distribution function of a particular binomial distribution is shown in Figure 3.2.

The expectation and variance of the binomial distribution can be computed straightforwardly, and are found to be

$$E(X) = np \quad \text{and} \quad \text{Var}(X) = np(1 - p).$$

3.5 The geometric distribution

The geometric distribution is the distribution of the number of independent success/fail trials until the first success is encountered. If X is the number of trials until a success is encountered, and each independent trial has probability p of being a success, we write

$$X \sim \text{Geom}(p).$$

Clearly X can take on any positive integer ($S_X = \{1, 2, \dots\}$), so to deduce the PMF, we need to calculate $P(X = k)$ for $k = 1, 2, 3, \dots$. In order to have $X = k$, we must have an ordered sequence of $k - 1$ failures followed by one success. By the multiplication rule, therefore,

$$P(X = k) = (1 - p)^{k-1}p, \quad k = 1, 2, 3, \dots$$

For the geometric distribution, it is possible to calculate an analytic form for the CDF as follows. If $X \sim \text{Geom}(p)$, then

$$\begin{aligned} F_X(k) &= P(X \leq k) \\ &= \sum_{j=1}^k (1 - p)^{j-1}p \\ &= p \sum_{j=1}^k (1 - p)^{j-1} \\ &= p \times \frac{1 - (1 - p)^k}{1 - (1 - p)} \quad (\text{geometric series}) \\ &= 1 - (1 - p)^k. \end{aligned}$$

Consequently, there is no need to tabulate the CDF of the geometric distribution. Also note that the CDF tends to 1 as k increases. This confirms that the PMF we defined does determine a valid probability distribution.

3.5.1 Useful series in probability

Notice that we used the sum of a geometric series in the derivation of the CDF. There are many other series that crop up in the study of probability. A few of the more commonly encountered series are listed below.

$$\begin{aligned} \sum_{i=1}^n a^{i-1} &= \frac{1 - a^n}{1 - a} && (a > 0) \\ \sum_{i=1}^{\infty} a^{i-1} &= \frac{1}{1 - a} && (0 < a < 1) \\ \sum_{i=1}^{\infty} i a^{i-1} &= \frac{1}{(1 - a)^2} && (0 < a < 1) \\ \sum_{i=1}^{\infty} i^2 a^{i-1} &= \frac{1 + a}{(1 - a)^3} && (0 < a < 1) \\ \sum_{i=1}^n i &= \frac{n(n + 1)}{2} \\ \sum_{i=1}^n i^2 &= \frac{1}{6}n(n + 1)(2n + 1). \end{aligned}$$

We will use two of these in the derivation of the expectation and variance of the geometric distribution.

3.5.2 Expectation and variance of geometric random quantities

Suppose that $X \sim \text{Geom}(p)$. Then

$$\begin{aligned} E(X) &= \sum_{i=1}^{\infty} i P(X = i) \\ &= \sum_{i=1}^{\infty} i(1-p)^{i-1}p \\ &= p \times \frac{1}{(1-[1-p])^2} \\ &= \frac{1}{p}. \end{aligned}$$

Similarly,

$$\begin{aligned} E(X^2) &= \sum_{i=1}^{\infty} i^2 P(X = i) \\ &= \sum_{i=1}^{\infty} i^2(1-p)^{i-1}p \\ &= p \times \frac{1 + [1-p]}{(1-[1-p])^3} \\ &= \frac{2-p}{p^2}, \end{aligned}$$

and so

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

3.6 The Poisson distribution

The Poisson distribution is a very important discrete probability distribution, which arises in many different contexts in probability and statistics. For example, Poisson random quantities are often used in place of binomial random quantities in situations where n is large, p is small, and the expectation np is stable. The Poisson distribution is particularly important in the context of stochastic modelling of biochemical networks, as the number of reaction events occurring in a short time interval is approximately Poisson.

A Poisson random variable, X with parameter λ is written as

$$X \sim Po(\lambda).$$

3.6.1 Poisson as the limit of a binomial

Let $X \sim B(n, p)$. Put $\lambda = E(X) = np$ and let n increase and p decrease so that λ remains constant.

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, 2, \dots, n.$$

Replacing p by λ/n gives

$$\begin{aligned} P(X = k) &= \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{n!}{k!(n-k)!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \\ &= \frac{\lambda^k}{k!} \frac{n!}{(n-k)! n^k} \frac{(1 - \lambda/n)^n}{(1 - \lambda/n)^k} \\ &= \frac{\lambda^k}{k!} \frac{n}{n} \frac{(n-1)}{n} \frac{(n-2)}{n} \dots \frac{(n-k+1)}{n} \frac{(1 - \lambda/n)^n}{(1 - \lambda/n)^k} \\ &\rightarrow \frac{\lambda^k}{k!} \times 1 \times 1 \times 1 \times \dots \times 1 \times \frac{e^{-\lambda}}{1}, \quad \text{as } n \rightarrow \infty \\ &= \frac{\lambda^k}{k!} e^{-\lambda}. \end{aligned}$$

To see the limit, note that $(1 - \lambda/n)^n \rightarrow e^{-\lambda}$ as n increases (compound interest formula).

3.6.2 PMF

If $X \sim Po(\lambda)$, then the PMF of X is

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, 2, 3, \dots$$

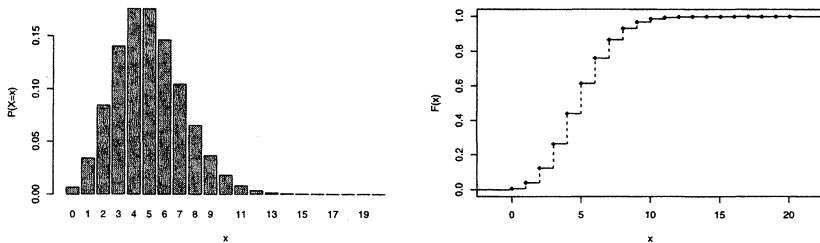
The PMF and CDF of $X \sim Po(5)$ are given in Figure 3.3. It is easy to verify that the PMF we have adopted for $X \sim Po(\lambda)$ does indeed define a valid probability distribution, as the probabilities sum to 1.

3.6.3 Expectation and variance of Poisson

If $X \sim Po(\lambda)$, the expectation and variance can be computed directly from the PMF, but it is clear from the binomial limit that we must have

$$E(X) = \lambda \quad \text{and} \quad \text{Var}(X) = \lambda.$$

That is, the mean and variance are both λ .

Figure 3.3 PMF and CDF for a $Po(5)$ distribution

3.6.4 Sum of Poisson random quantities

One of the particularly convenient properties of the Poisson distribution is that the sum of two independent Poisson random quantities is also a Poisson random quantity. If $X \sim Po(\lambda)$ and $Y \sim Po(\mu)$ and X and Y are independent, then $Z = X + Y \sim Po(\lambda + \mu)$. Clearly this result extends to the sum of many independent Poisson random variables. The proof is straightforward, but omitted.

This property of the Poisson distribution is fundamental to its usefulness for modelling “counts.” Consider, for example, modelling the number of counts recorded by a Geiger counter of relatively stable radioactive isotope (with a very long half-life). If we consider the number of counts in one second, the binomial limit argument suggests that since we are dealing with a huge number of molecules, and each independently has a constant but tiny probability of decaying and triggering a count, then the number of counts in that one-second interval should be Poisson. Suppose that the rate of decay is such that the count is $Po(\lambda)$, for some λ . Then in each one-second interval the count will be the same, independently of all other intervals. So, using the above additive property, the number of counts in a two-second interval will be $Po(2\lambda)$, and in a three-second interval it will be $Po(3\lambda)$, etc. So in an interval of length t seconds, the number of counts will be $Po(\lambda t)$.[‡] This is the Poisson process.

3.6.5 The Poisson process

A sequence of timed observations is said to follow a *Poisson process* with *rate* λ if the number of observations, X , in any interval of length t is such that

$$X \sim Po(\lambda t)$$

[‡] Of course we are ignoring the fact that once a molecule has decayed, it cannot decay again. But we are assuming that the number of available molecules is huge, and the number of decay events over the time scale of interest is so small that this effect can be ignored. For an isotope with a much shorter half-life, a *pure death process* would provide a much more suitable model. However, a formal discussion of processes of this nature will have to wait until Chapter 5.

and the numbers of events in disjoint intervals are independent of one another. We see that this is a simple model for discrete events occurring continuously in time. It turns out that the Poisson process is too simple to realistically model the time-evolution of biochemical networks, but a detailed understanding of this simple process is key to understanding more sophisticated stochastic processes. It is interesting to wonder about the time between successive events of the Poisson process, but it is clear that such random times do not have a discrete distribution. This, therefore, motivates the study of continuous random quantities.

3.7 Continuous probability models

3.7.1 Introduction, PDF and CDF

We now turn to techniques for handling continuous random quantities. These are random quantities with a sample space which is neither finite nor countably infinite. The sample space is usually taken to be the real line, or a part thereof. Continuous probability models are appropriate if the result of an experiment is a continuous *measurement*, rather than a *count* of a discrete set. In the context of systems biology, such measurements often (but not always) correspond to measurements of time.

If X is a continuous random quantity with sample space S_X , then for any particular $a \in S_X$, we generally have that

$$P(X = a) = 0.$$

This is because the sample space is so “large” and every possible outcome so “small” that the probability of any *particular* value is vanishingly small. Therefore, the probability mass function we defined for discrete random quantities is inappropriate for understanding continuous random quantities. In order to understand continuous random quantities, a little calculus is required.

Definition 3.13 (probability density function) *If X is a continuous random quantity, then there exists a function $f_X(x)$, called the probability density function (PDF), which satisfies the following:*

1. $f_X(x) \geq 0$, $\forall x$; (the symbol “ \forall ” means “for all”)

2. $\int_{-\infty}^{\infty} f_X(x) dx = 1$;

3. $P(a \leq X \leq b) = \int_a^b f_X(x) dx$ for any $a \leq b$.

Consequently we have

$$\begin{aligned} P(x \leq X \leq x + \delta x) &= \int_x^{x+\delta x} f_X(y) dy \\ &\simeq f_X(x)\delta x, && \text{(for small } \delta x) \\ \Rightarrow f_X(x) &\simeq \frac{P(x \leq X \leq x + \delta x)}{\delta x} \end{aligned}$$

and so we may interpret the PDF as

$$f_X(x) = \lim_{\delta x \rightarrow 0} \frac{P(x \leq X \leq x + \delta x)}{\delta x}.$$

Note that PDFs are *not* probabilities. For example, the density can take values greater than 1 in some regions as long as it still integrates to 1. Also note that because $P(X = a) = 0$, we have $P(X \leq k) = P(X < k)$ for continuous random quantities.

Definition 3.14 (cumulative distribution function) *Earlier in this chapter we defined the cumulative distribution function of a random variable X to be*

$$F_X(x) = P(X \leq x), \quad \forall x.$$

This definition works just as well for continuous random quantities, and is one of the many reasons why the distribution function is so useful. For a discrete random quantity we had

$$F_X(x) = P(X \leq x) = \sum_{\{y \in S_X | y \leq x\}} P(X = y),$$

but for a continuous random quantity we have the continuous analogue

$$\begin{aligned} F_X(x) &= P(X \leq x) \\ &= P(-\infty \leq X \leq x) \\ &= \int_{-\infty}^x f_X(z) dz. \end{aligned}$$

Just as in the discrete case, the distribution function is defined for all $x \in \mathbb{R}$, even if the sample space S_X is not the whole of the real line.

Proposition 3.8

1. *Since it represents a probability, $F_X(x) \in [0, 1]$.*
2. *$F_X(-\infty) = 0$ and $F_X(\infty) = 1$.*
3. *If $a < b$, then $F_X(a) \leq F_X(b)$. i.e. $F_X(\cdot)$ is a non-decreasing function.*
4. *When X is continuous, $F_X(x)$ is continuous. Also, by the Fundamental Theorem of Calculus, we have*

$$\frac{d}{dx} F_X(x) = f_X(x),$$

and so the slope of the CDF $F_X(x)$ is the PDF $f_X(x)$.

Proposition 3.9 *The median of a random quantity is the value m which is the “middle” of the distribution. That is, it is the value m such that*

$$P(X \leq m) = \frac{1}{2}.$$

Equivalently, it is the value, m such that

$$F_X(m) = 0.5.$$

Similarly, the lower quartile of a random quantity is the value l such that

$$F_X(l) = 0.25,$$

and the upper quartile is the value u such that

$$F_X(u) = 0.75.$$

3.7.2 Properties of continuous random quantities

Proposition 3.10 The expectation or mean of a continuous random quantity X is given by

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

which is just the continuous analogue of the corresponding formula for discrete random quantities. Similarly, the variance is given by

$$\begin{aligned} \text{Var}(X) &= \int_{-\infty}^{\infty} [x - E(X)]^2 f_X(x) dx \\ &= \int_{-\infty}^{\infty} x^2 f_X(x) dx - [E(X)]^2. \end{aligned}$$

Note that the expectation of $g(X)$ is given by

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

and so the variance is just

$$\text{Var}(X) = E([X - E(X)]^2) = E(X^2) - [E(X)]^2$$

as in the discrete case.

Note that all of the properties of expectation and variance derived for discrete random quantities also hold true in the continuous case. It is also worth explicitly noting the continuous version of Proposition 3.7.

Proposition 3.11 If F is an event and X is a continuous random quantity, then

$$\begin{aligned} P(F) &= \int_{-\infty}^{\infty} P(F|X=x) f(x) dx \\ &= E(P(F|X)), \end{aligned}$$

where $P(F|X)$ is the random quantity which takes the value $P(F|X=x)$ when X takes the value x .

Proposition 3.12 (PDF of a linear transformation) Let X be a continuous random quantity with PDF $f_X(x)$ and CDF $F_X(x)$, and let $Y = aX + b$. The PDF of Y is given by

$$f_Y(y) = \left| \frac{1}{a} \right| f_X\left(\frac{y-b}{a}\right).$$

Proof. First assume that $a > 0$. It turns out to be easier to work out the CDF first:

$$\begin{aligned} F_Y(y) &= P(Y \leq y) \\ &= P(aX + b \leq y) \\ &= P\left(X \leq \frac{y-b}{a}\right) && \text{(since } a > 0\text{)} \\ &= F_X\left(\frac{y-b}{a}\right). \end{aligned}$$

So,

$$F_Y(y) = F_X\left(\frac{y-b}{a}\right),$$

and by differentiating both sides with respect to y we get

$$f_Y(y) = \frac{1}{a} f_X\left(\frac{y-b}{a}\right).$$

If $a < 0$, the inequality changes sign, introducing a minus sign in the expression for the PDF. Both cases can therefore be summarised by the result. \square

Proposition 3.13 (PDF of a differentiable 1-1 transformation) *Using a similar argument it is straightforward to deduce the PDF of an arbitrary differentiable invertible transformation $Y = g(X)$ as*

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|.$$

Note that the term $\left| \frac{d}{dy} g^{-1}(y) \right|$ is known as the ‘‘Jacobian’’ of the transformation.

3.7.3 The law of large numbers

In many scenarios involving stochastic simulation, it is desirable to approximate the expectation of a random quantity, X , by the sample mean of a collection of independent observations of the random quantity X . Suppose we have a computer program that can simulate independent realisations of X (we will see in the next chapter how to do this); X_1, X_2, \dots, X_n . The *sample mean* is given by

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

If in some appropriate sense the sample mean converges to the expectation, $E(X)$, then for large n , we can use \bar{X} as an estimate of $E(X)$. It turns out that for random quantities with finite variance, this is indeed true, but a couple of lemmas are required before this can be made precise.

Lemma 3.6 *If X has finite mean and variance, $E(X) = \mu$ and $\text{Var}(X) = \sigma^2$, and X_1, X_2, \dots, X_n is a collection of independent realisations of X used to form the*

sample mean \bar{X} , then

$$E(\bar{X}) = \mu \quad \text{and} \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}.$$

Proof. The result follows immediately from Propositions 3.5 and 3.6. \square

Lemma 3.7 (Markov's inequality) If $X \geq 0$ is a non-negative random quantity with finite expectation $E(X) = \mu$ we have

$$P(X \geq a) \leq \frac{\mu}{a}, \quad \forall a > 0.$$

Proof. Note that this result is true for arbitrary non-negative random quantities, but we present here the proof for the continuous case.

$$\begin{aligned} \frac{\mu}{a} &= \frac{E(X)}{a} \\ &= \frac{1}{a} \int_0^{\infty} x f(x) dx \\ &= \int_0^{\infty} \frac{x}{a} f(x) dx \\ &\geq \int_a^{\infty} \frac{x}{a} f(x) dx \\ &\geq \int_a^{\infty} \frac{a}{a} f(x) dx \\ &= \int_a^{\infty} f(x) dx \\ &= P(X \geq a). \end{aligned}$$

\square

Lemma 3.8 (Chebyshev's inequality) If X is a random quantity with finite mean $E(X) = \mu$ and variance $\text{Var}(X) = \sigma^2$ we have

$$P(|X - \mu| < k\sigma) \geq 1 - \frac{1}{k^2}, \quad \forall k > 0.$$

Proof. Since $(X - \mu)^2$ is positive with expectation σ^2 , Markov's inequality gives

$$P([X - \mu]^2 \geq a) \leq \frac{\sigma^2}{a}.$$

Putting $a = k^2\sigma^2$ then gives

$$\begin{aligned} P([X - \mu]^2 \geq k^2\sigma^2) &\leq \frac{1}{k^2} \\ \Rightarrow P(|X - \mu| \geq k\sigma) &\leq \frac{1}{k^2} \\ \Rightarrow P(|X - \mu| < k\sigma) &\geq 1 - \frac{1}{k^2}. \end{aligned}$$

\square

We are now in a position to state the main result.

Proposition 3.14 (weak law of large numbers, WLLN) For X with finite mean $E(X) = \mu$ and variance $\text{Var}(X) = \sigma^2$, if X_1, X_2, \dots, X_n is an independent sample from X used to form the sample mean \bar{X} , we have

$$P(|\bar{X} - \mu| < \varepsilon) \geq 1 - \frac{\sigma^2}{n\varepsilon^2} \xrightarrow{\infty} 1, \quad \forall \varepsilon > 0.$$

In other words, the WLLN states that no matter how small one chooses the positive constant ε , the probability that \bar{X} is within a distance ε of μ tends to 1 as the sample size n increases. This is a precise sense in which the sample mean “converges” to the expectation.

Proof. Using Lemmas 3.6 and 3.8 we have

$$P\left(|\bar{X} - \mu| < k \frac{\sigma}{\sqrt{n}}\right) \geq 1 - \frac{1}{k^2}.$$

Substituting $k = \varepsilon\sqrt{n}/\sigma$ gives the result. \square

This result is known as the *weak* law, as there is a corresponding *strong* law, which we state without proof.

Proposition 3.15 (strong law of large numbers, SLLN) For X with finite mean $E(X) = \mu$ and variance $\text{Var}(X) = \sigma^2$, if X_1, X_2, \dots, X_n is an independent sample from X used to form the sample mean \bar{X} , we have

$$P\left(\bar{X} \xrightarrow{\infty} \mu\right) = 1.$$

3.8 The uniform distribution

Now that we understand the basic properties of continuous random quantities, we can look at some of the important standard continuous probability models. The simplest of these is the uniform distribution. This distribution turns out to be central to the theory of stochastic simulation that will be developed in the next chapter.

The random quantity X has a *uniform distribution* over the range $[a, b]$, written

$$X \sim U(a, b)$$

if the PDF is given by

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b, \\ 0, & \text{otherwise.} \end{cases}$$

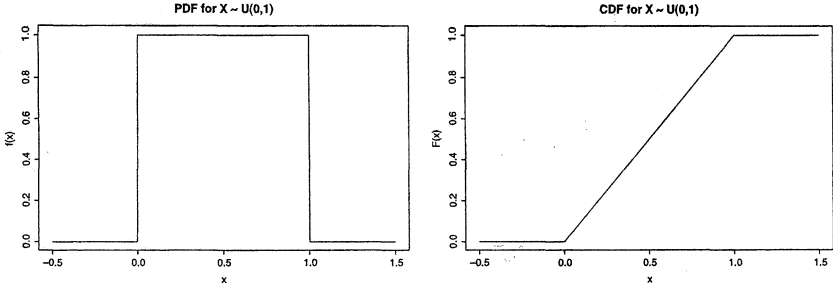


Figure 3.4 PDF and CDF for a $U(0, 1)$ distribution

Thus if $x \in [a, b]$, then

$$\begin{aligned}
 F_X(x) &= \int_{-\infty}^x f_X(y) dy \\
 &= \int_{-\infty}^a f_X(y) dy + \int_a^x f_X(y) dy \\
 &= 0 + \int_a^x \frac{1}{b-a} dy \\
 &= \frac{x-a}{b-a}.
 \end{aligned}$$

Therefore,

$$F_X(x) = \begin{cases} 0, & x < a, \\ \frac{x-a}{b-a}, & a \leq x \leq b, \\ 1, & x > b. \end{cases}$$

We can plot the PDF and CDF in order to see the “shape” of the distribution. Plots for $X \sim U(0, 1)$ are shown in Figure 3.4.

Clearly the lower quartile, median, and upper quartile of the uniform distribution are

$$\frac{3}{4}a + \frac{1}{4}b, \quad \frac{a+b}{2}, \quad \frac{1}{4}a + \frac{3}{4}b,$$

respectively. The expectation of a uniform random quantity is

$$\begin{aligned}
 E(X) &= \int_{-\infty}^{\infty} x f_X(x) dx \\
 &= \int_{-\infty}^a x f_X(x) dx + \int_a^b x f_X(x) dx + \int_b^{\infty} x f_X(x) dx \\
 &= 0 + \int_a^b \frac{x}{b-a} dx + 0 \\
 &= \left[\frac{x^2}{2(b-a)} \right]_a^b \\
 &= \frac{a+b}{2}.
 \end{aligned}$$

We can also calculate the variance of X . First we calculate $E(X^2)$ as follows:

$$\begin{aligned}
 E(X^2) &= \int_a^b \frac{x^2}{b-a} dx \\
 &= \frac{b^2 + ab + a^2}{3}.
 \end{aligned}$$

Now,

$$\begin{aligned}
 \text{Var}(X) &= E(X^2) - E(X)^2 \\
 &= \frac{b^2 + ab + a^2}{3} - \frac{(a+b)^2}{4} \\
 &= \frac{(b-a)^2}{12}.
 \end{aligned}$$

The uniform distribution is too simple to realistically model actual experimental data, but is very useful for computer simulation, as random quantities from many different distributions can be obtained from $U(0, 1)$ random quantities.

3.9 The exponential distribution

For reasons still to be explored, it turns out that the exponential distribution is the most important continuous distribution in the theory of discrete-event stochastic simulation. It is therefore vital to have a good understanding of this distribution and its many useful properties. We will begin by introducing this distribution in the abstract, but we will then go on to see why it arises so naturally by exploring its relationship with the Poisson process.

The random variable X has an *exponential distribution* with parameter $\lambda > 0$, written

$$X \sim \text{Exp}(\lambda)$$

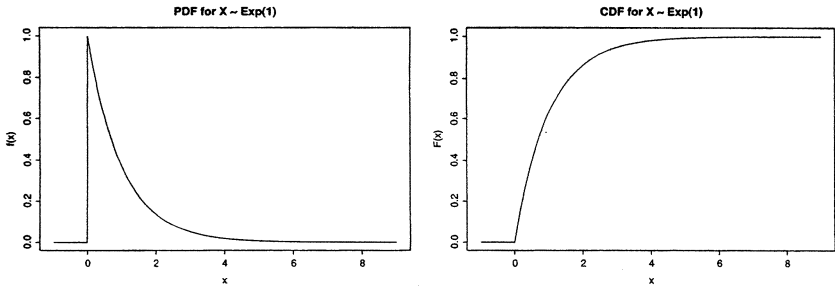


Figure 3.5 PDF and CDF for an $Exp(1)$ distribution

if it has PDF

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & \text{otherwise.} \end{cases}$$

The distribution function, $F_X(x)$ is therefore given by

$$F_X(x) = \begin{cases} 0, & x < 0, \\ 1 - e^{-\lambda x}, & x \geq 0. \end{cases}$$

The PDF and CDF for an $Exp(1)$ are shown in Figure 3.5.

The expectation of the exponential distribution is

$$\begin{aligned} E(X) &= \int_0^{\infty} x \lambda e^{-\lambda x} dx \\ &= [-x e^{-\lambda x}]_0^{\infty} + \int_0^{\infty} e^{-\lambda x} dx && \text{(by parts)} \\ &= 0 + \left[\frac{e^{-\lambda x}}{-\lambda} \right]_0^{\infty} \\ &= \frac{1}{\lambda}. \end{aligned}$$

Also,

$$E(X^2) = \int_0^{\infty} x^2 \lambda e^{-\lambda x} dx = \frac{2}{\lambda^2},$$

and so

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

Note that this means the expectation and standard deviation are both $1/\lambda$.

Notes

1. As λ increases, the probability of small values of X increases and the mean decreases.
2. The median m is given by

$$m = \frac{\log 2}{\lambda} = \log 2 E(X) < E(X).$$

3. The exponential distribution is often used to model times between random events (such as biochemical reactions). Some of the reasons are given below.

Proposition 3.16 (memoryless property) *If $X \sim \text{Exp}(\lambda)$, then for any $s, t \geq 0$ we have*

$$P(X > (s+t) | X > t) = P(X > s).$$

If we think of the exponential random quantity as representing the time to an event, we can regard that time as a “lifetime.” Then the proposition states that the probability of “surviving” a further time s , having survived time t , is the same as the original probability of surviving a time s . This is called the “memoryless” property of the distribution (as the distribution “forgets” that it has survived to time t). It is therefore the continuous analogue of the geometric distribution, which is the (unique) discrete distribution with such a property.

Proof:

$$\begin{aligned} P(X > (s+t) | X > t) &= \frac{P([X > (s+t)] \cap [X > t])}{P(X > t)} \\ &= \frac{P(X > (s+t))}{P(X > t)} \\ &= \frac{1 - P(X \leq (s+t))}{1 - P(X \leq t)} \\ &= \frac{1 - F_X(s+t)}{1 - F_X(t)} \\ &= \frac{1 - [1 - e^{-\lambda(s+t)}]}{1 - [1 - e^{-\lambda t}]} \\ &= e^{-\lambda s} \\ &= 1 - [1 - e^{-\lambda s}] \\ &= 1 - F_X(s) \\ &= 1 - P(X \leq s) \\ &= P(X > s). \end{aligned}$$

□

Proposition 3.17 *Consider a Poisson process with rate λ . Let T be the time to the first event (after zero). Then $T \sim \text{Exp}(\lambda)$.*

Proof. Let N_t be the number of events in the interval $(0, t]$ (for given fixed $t > 0$). We have seen previously that (by definition) $N_t \sim Po(\lambda t)$. Consider the CDF of T ,

$$\begin{aligned} F_T(t) &= P(T \leq t) \\ &= 1 - P(T > t) \\ &= 1 - P(N_t = 0) \\ &= 1 - \frac{(\lambda t)^0 e^{-\lambda t}}{0!} \\ &= 1 - e^{-\lambda t}. \end{aligned}$$

This is the distribution function of an $Exp(\lambda)$ random quantity, and so $T \sim Exp(\lambda)$.

□

So the time to the *first* event of a Poisson process is an exponential random variable. But then using the independence properties of the Poisson process, it should be reasonably clear that the time between any two such events has the same exponential distribution. Thus the times between events of the Poisson process are exponential.

There is another way of thinking about the Poisson process that this result makes clear. For an infinitesimally small time interval dt we have

$$P(T \leq dt) = 1 - e^{-\lambda dt} = 1 - (1 - \lambda dt) = \lambda dt,$$

and due to the independence property of the Poisson process, this is the probability for any time interval of length dt . The Poisson process can therefore be thought of as a process with constant event “hazard” λ , where the “hazard” is essentially a measure of “event density” on the time axis. The exponential distribution with parameter λ can therefore also be reinterpreted as the time to an event of constant hazard λ .

The two properties above are probably the most fundamental. However, there are several other properties that we will require of the exponential distribution when we come to use it to simulate discrete stochastic models of biochemical networks, and so they are mentioned here for future reference. The first describes the distribution of the minimum of a collection of independent exponential random quantities.

Proposition 3.18 *If $X_i \sim Exp(\lambda_i)$, $i = 1, 2, \dots, n$, are independent random variables, then*

$$X_0 \equiv \min_i \{X_i\} \sim Exp(\lambda_0), \text{ where } \lambda_0 = \sum_{i=1}^n \lambda_i.$$

Proof. First note that for $X \sim \text{Exp}(\lambda)$ we have $P(X > x) = e^{-\lambda x}$. Then

$$\begin{aligned} P(X_0 > x) &= P\left(\min_i\{X_i\} > x\right) \\ &= P([X_1 > x] \cap [X_2 > x] \cap \cdots \cap [X_n > x]) \\ &= \prod_{i=1}^n P(X_i > x) \\ &= \prod_{i=1}^n e^{-\lambda_i x} \\ &= e^{-x \sum_{i=1}^n \lambda_i} \\ &= e^{-\lambda_0 x}. \end{aligned}$$

So $P(X_0 \leq x) = 1 - e^{-\lambda_0 x}$ and hence $X_0 \sim \text{Exp}(\lambda_0)$. \square

The next lemma is for the following proposition.

Lemma 3.9 *Suppose that $X \sim \text{Exp}(\lambda)$ and $Y \sim \text{Exp}(\mu)$ are independent random variables. Then*

$$P(X < Y) = \frac{\lambda}{\lambda + \mu}.$$

Proof.

$$\begin{aligned} P(X < Y) &= \int_0^\infty P(X < Y|Y = y) f(y) dy && \text{(Proposition 3.11)} \\ &= \int_0^\infty P(X < y) f(y) dy \\ &= \int_0^\infty (1 - e^{-\lambda y}) \mu e^{-\mu y} dy \\ &= \frac{\lambda}{\lambda + \mu} \end{aligned}$$

\square

This next result gives the likelihood of a particular exponential random quantity of an independent collection being the smallest.

Proposition 3.19 *If $X_i \sim \text{Exp}(\lambda_i)$, $i = 1, 2, \dots, n$ are independent random variables, let j be the index of the smallest of the X_i . Then j is a discrete random variable with PMF*

$$\pi_i = \frac{\lambda_i}{\lambda_0}, \quad i = 1, 2, \dots, n, \text{ where } \lambda_0 = \sum_{i=1}^n \lambda_i.$$

Proof.

$$\begin{aligned} \pi_j &= P\left(X_j < \min_{i \neq j}\{X_i\}\right) \\ &= P(X_j < Y) \end{aligned}$$

where $Y = \min_{i \neq j} \{X_i\}$, so that $Y \sim \text{Exp}(\lambda_{-j})$, where $\lambda_{-j} = \sum_{i \neq j} \lambda_i$

$$\begin{aligned} &= \frac{\lambda_j}{\lambda_j + \lambda_{-j}} \quad (\text{by the lemma}) \\ &= \frac{\lambda_j}{\lambda_0} \end{aligned}$$

□

This final result is a trivial consequence of Proposition 3.12.

Proposition 3.20 Consider $X \sim \text{Exp}(\lambda)$. Then for $\alpha > 0$, $Y = \alpha X$ has distribution

$$Y \sim \text{Exp}(\lambda/\alpha).$$

3.10 The normal/Gaussian distribution

3.10.1 Definition and properties

Another fundamental distribution in probability theory is the normal or Gaussian distribution. It turns out that sums of random quantities often approximately follow a normal distribution. Since the change of state of a biochemical network can sometimes be represented as a sum of random quantities, it turns out that normal distributions are useful in this context.

Definition 3.15 A random quantity X has a normal distribution with parameters μ and σ^2 , written

$$X \sim N(\mu, \sigma^2)$$

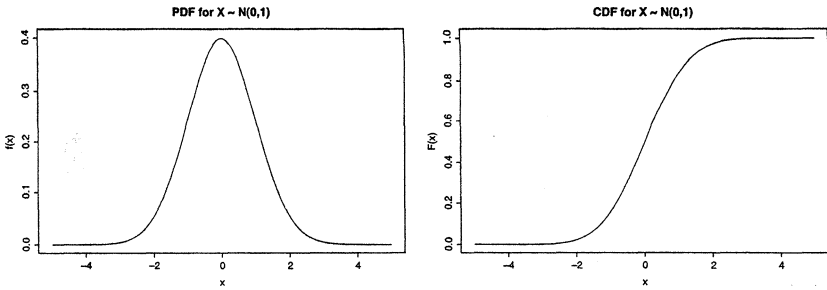
if it has probability density function

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right\}, \quad -\infty < x < \infty,$$

for $\sigma > 0$.

Note that $f_X(x)$ is symmetric about $x = \mu$, and so (provided the density integrates to 1), the median of the distribution will be μ . Checking that the density integrates to 1 requires the computation of a difficult integral. However, it follows directly from the known ‘‘Gaussian’’ integral

$$\int_{-\infty}^{\infty} e^{-\alpha x^2} dx = \sqrt{\frac{\pi}{\alpha}}, \quad \alpha > 0,$$

Figure 3.6 PDF and CDF for a $N(0, 1)$ distribution

since then

$$\begin{aligned}
 \int_{-\infty}^{\infty} f_X(x) dx &= \int_{-\infty}^{\infty} \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} dx \\
 &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp\left\{-\frac{1}{2\sigma^2}z^2\right\} dz \quad (\text{putting } z = x - \mu) \\
 &= \frac{1}{\sigma\sqrt{2\pi}} \sqrt{\frac{\pi}{1/2\sigma^2}} \\
 &= 1.
 \end{aligned}$$

Now that we know that the given PDF represents a valid density, we can calculate the expectation and variance of the normal distribution as follows:

$$\begin{aligned}
 E(X) &= \int_{-\infty}^{\infty} x f_X(x) dx \\
 &= \int_{-\infty}^{\infty} x \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} dx \\
 &= \mu.
 \end{aligned}$$

The last line follows after a little algebra and calculus. Similarly,

$$\begin{aligned}
 \text{Var}(X) &= \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx \\
 &= \int_{-\infty}^{\infty} (x - \mu)^2 \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\} dx \\
 &= \sigma^2.
 \end{aligned}$$

The PDF and CDF for a $N(0, 1)$ are shown in Figure 3.6.

3.10.2 The standard normal distribution

A standard normal random quantity is a normal random quantity with zero mean and variance equal to 1. It is usually denoted Z , so that

$$Z \sim N(0, 1).$$

Therefore, the density of Z , which is usually denoted $\phi(z)$, is given by

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z^2\right\}, \quad -\infty < z < \infty.$$

It is important to note that the PDF of the standard normal is symmetric about zero. The distribution function of a standard normal random quantity is denoted $\Phi(z)$, that is

$$\Phi(z) = \int_{-\infty}^z \phi(x) dx.$$

There is no neat analytic expression for $\Phi(z)$, so tables of the CDF are used. Of course, we do know that $\Phi(-\infty) = 0$ and $\Phi(\infty) = 1$, as it is a distribution function. Also, because of the symmetry of the PDF about zero, it is clear that we must also have $\Phi(0) = 1/2$, and this can prove useful in calculations. The standard normal distribution is important because it is easy to transform any normal random quantity to a standard normal random quantity by means of a simple linear scaling. Consider $Z \sim N(0, 1)$ and put

$$X = \mu + \sigma Z,$$

for $\sigma > 0$. Then $X \sim N(\mu, \sigma^2)$. To show this, we must show that the PDF of X is the PDF for a $N(\mu, \sigma^2)$ random quantity. Using Proposition 3.12 we have

$$\begin{aligned} f_X(x) &= \frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right) \\ &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right\}, \end{aligned}$$

which is the PDF of a $N(\mu, \sigma^2)$ distribution. Conversely, if

$$X \sim N(\mu, \sigma^2)$$

then

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1).$$

Even more importantly, the distribution function of X is given by

$$F_X(x) = \Phi\left(\frac{x - \mu}{\sigma}\right),$$

and so the cumulative probabilities for any normal random quantity can be calculated using tables for the standard normal distribution.

It is a straightforward generalisation of the above result to see that any linear rescaling of a normal random quantity is another normal random quantity. However,

the normal distribution also has an additive property similar to that of the Poisson distribution.

Proposition 3.21 *If $X_1 \sim N(\mu_1, \sigma_1^2)$ and $X_2 \sim N(\mu_2, \sigma_2^2)$ are independent normal random quantities, then their sum $Y = X_1 + X_2$ is also normal, and*

$$Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2).$$

The proof is straightforward, but omitted. Putting the above results together, we can see that any linear combination of independent normal random quantities will be normal. We can then use the results for the mean and variance of a linear combination to deduce the mean and variance of the resulting normal distribution.

3.10.3 The Central Limit Theorem

Now that we know about the normal distribution and its properties, we need to understand how and why it arises so frequently. The answer is contained in the Central Limit Theorem, which is stated without proof; see (for example) Miller & Miller (2004) for a proof and further details.

Theorem 3.3 (Central Limit Theorem) *If X_1, X_2, \dots, X_n are independent realizations of an arbitrary random quantity with mean μ and variance σ^2 , let \bar{X}_n be the sample mean and define*

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}.$$

Then the limiting distribution of Z_n as $n \rightarrow \infty$ is the standard normal distribution. In other words, $\forall z$,

$$P(Z_n \leq z) \xrightarrow[n \rightarrow \infty]{} \Phi(z).$$

Note that Z_n is just \bar{X}_n linearly scaled so that it has mean zero and variance 1. By rescaling it this way, it is then possible to compare it with the standard normal distribution. Typically, Z_n will not have a normal distribution for any finite value of n , but what the CLT tells us is that the distribution of Z_n becomes more like that of a standard normal as n increases. Then since \bar{X}_n is just a linear scaling of Z_n , \bar{X}_n must also become more normal as n increases. And since $S_n \equiv \sum_{i=1}^n X_i$ is just a linear scaling of \bar{X}_n , this must also become more normal as n increases. In practice, values of n as small as 20 are often quite adequate for a normal approximation to work quite well.

3.10.4 Normal approximation of binomial and Poisson

A Bernoulli random quantity with parameter p , written $Bern(p)$, is a $B(1, p)$ random quantity. It therefore can take only the values zero or 1, and has mean p and variance $p(1-p)$. It is of interest because it can be used to construct the binomial distribution. Let $I_k \sim Bern(p)$, $k = 1, 2, \dots, n$. Then put

$$X = \sum_{k=1}^n I_k.$$

It is clear that X is a count of the number of successes in n trials, and therefore $X \sim B(n, p)$. Then, because of the Central Limit Theorem, this will be well approximated by a normal distribution if n is large and p is not too extreme (if p is very small or very large, a Poisson approximation will be more appropriate). A useful guide is that if

$$0.1 \leq p \leq 0.9 \quad \text{and} \quad n > \max \left[\frac{9(1-p)}{p}, \frac{9p}{1-p} \right]$$

then the binomial distribution may be adequately approximated by a normal distribution. It is important to understand exactly what is meant by this statement. No matter how large n is, the binomial will always be a discrete random quantity with a PMF, whereas the normal is a continuous random quantity with a PDF. These two distributions will always be qualitatively different. The similarity is measured in terms of the CDF, which has a consistent definition for both discrete and continuous random quantities. It is the CDF of the binomial which can be well approximated by a normal CDF. Fortunately, it is the CDF which matters for typical computations involving cumulative probabilities.

When the n and p of a binomial distribution are appropriate for approximation by a normal distribution, the approximation is achieved by matching expectation and variance. That is

$$B(n, p) \simeq N(np, np[1-p]).$$

Since the Poisson is derived from the binomial, it is unsurprising that in certain circumstances, the Poisson distribution may also be approximated by the normal. It is generally considered appropriate to apply the approximation if the mean of the Poisson is bigger than 20. Again the approximation is done by matching mean and variance:

$$X \sim Po(\lambda) \simeq N(\lambda, \lambda) \text{ for } \lambda > 20.$$

3.11 The gamma distribution

The *gamma function*, $\Gamma(x)$, is defined by the integral

$$\Gamma(x) = \int_0^{\infty} y^{x-1} e^{-y} dy.$$

By integrating by parts, it is easy to see that $\Gamma(x+1) = x\Gamma(x)$, and it is similarly clear that $\Gamma(1) = 1$. Together these give $\Gamma(n+1) = n!$ for integer $n > 0$, and so the gamma function can be thought of as a generalisation of the factorial function to non-integer values, x .[§] A graph of the function for small positive values is given in Figure 3.7. It is also worth noting that $\Gamma(1/2) = \sqrt{\pi}$. The gamma function is used in the definition of the *gamma distribution*.

Definition 3.16 *The random variable X has a gamma distribution with parameters*

[§] Recall that $n!$ (pronounced “n-factorial”) is given by $n! = 1 \times 2 \times 3 \times \dots \times (n-1) \times n$.

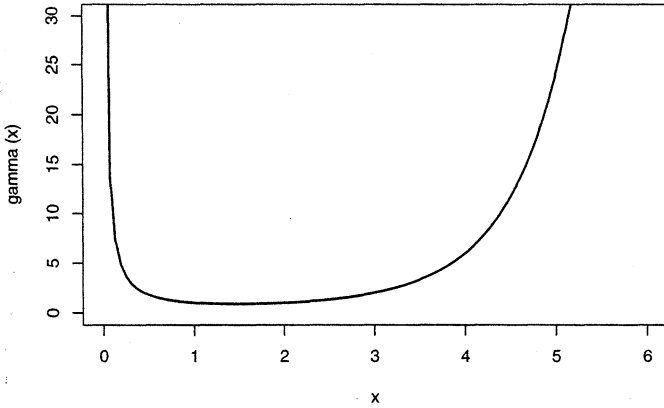


Figure 3.7 Graph of $\Gamma(x)$ for small positive values of x

$\alpha, \beta > 0$, written $\Gamma(\alpha, \beta)$ if it has PDF

$$f(x) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & x > 0 \\ 0 & x \leq 0. \end{cases}$$

Using the substitution $y = \beta x$ and the definition of the gamma function, it is straightforward to see that this density must integrate to 1. It is also clear that $\Gamma(1, \lambda) = \text{Exp}(\lambda)$, so the gamma distribution is a generalisation of the exponential distribution. It is also relatively straightforward to compute the mean and variance as

$$E(X) = \frac{\alpha}{\beta}, \quad \text{Var}(X) = \frac{\alpha}{\beta^2}.$$

We compute $E(X)$ as follows ($\text{Var}(X)$ is similar),

$$\begin{aligned} E(X) &= \int_0^\infty x f(x) dx \\ &= \int_0^\infty x \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} dx \\ &= \frac{\alpha}{\beta} \int_0^\infty \frac{\beta^{\alpha+1}}{\alpha \Gamma(\alpha)} x^\alpha e^{-\beta x} dx \\ &= \frac{\alpha}{\beta} \int_0^\infty \frac{\beta^{\alpha+1}}{\Gamma(\alpha+1)} x^\alpha e^{-\beta x} dx \\ &= \frac{\alpha}{\beta}. \end{aligned}$$

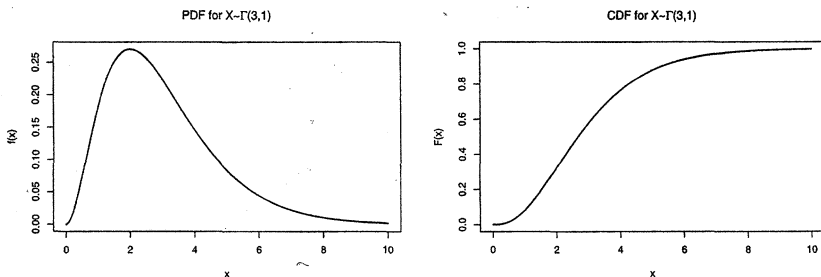


Figure 3.8 PDF and CDF for a $\Gamma(3, 1)$ distribution

The last line follows as the integral on the line above is the integral of a $\Gamma(\alpha + 1, \beta)$ density, and hence must be 1. A plot of the PDF and CDF of a gamma distribution is shown in Figure 3.8. The gamma distribution has a very important property.

Proposition 3.22 *If $X_1 \sim \Gamma(\alpha_1, \beta)$ and $X_2 \sim \Gamma(\alpha_2, \beta)$ are independent, and $Y = X_1 + X_2$, then*

$$Y \sim \Gamma(\alpha_1 + \alpha_2, \beta).$$

The proof is straightforward but omitted. This clearly generalises to more than two gamma random quantities, and hence implies that the sum of n independent $Exp(\lambda)$ random quantities is $\Gamma(n, \lambda)$. Since we know that the interevent times of a Poisson process with rate λ are $Exp(\lambda)$, we now know that the time to the n th event of a Poisson process is $\Gamma(n, \lambda)$. This close relationship between the gamma distribution, the exponential distribution, and the Poisson process is the reason why the gamma distribution turns out to be important for discrete stochastic modelling.

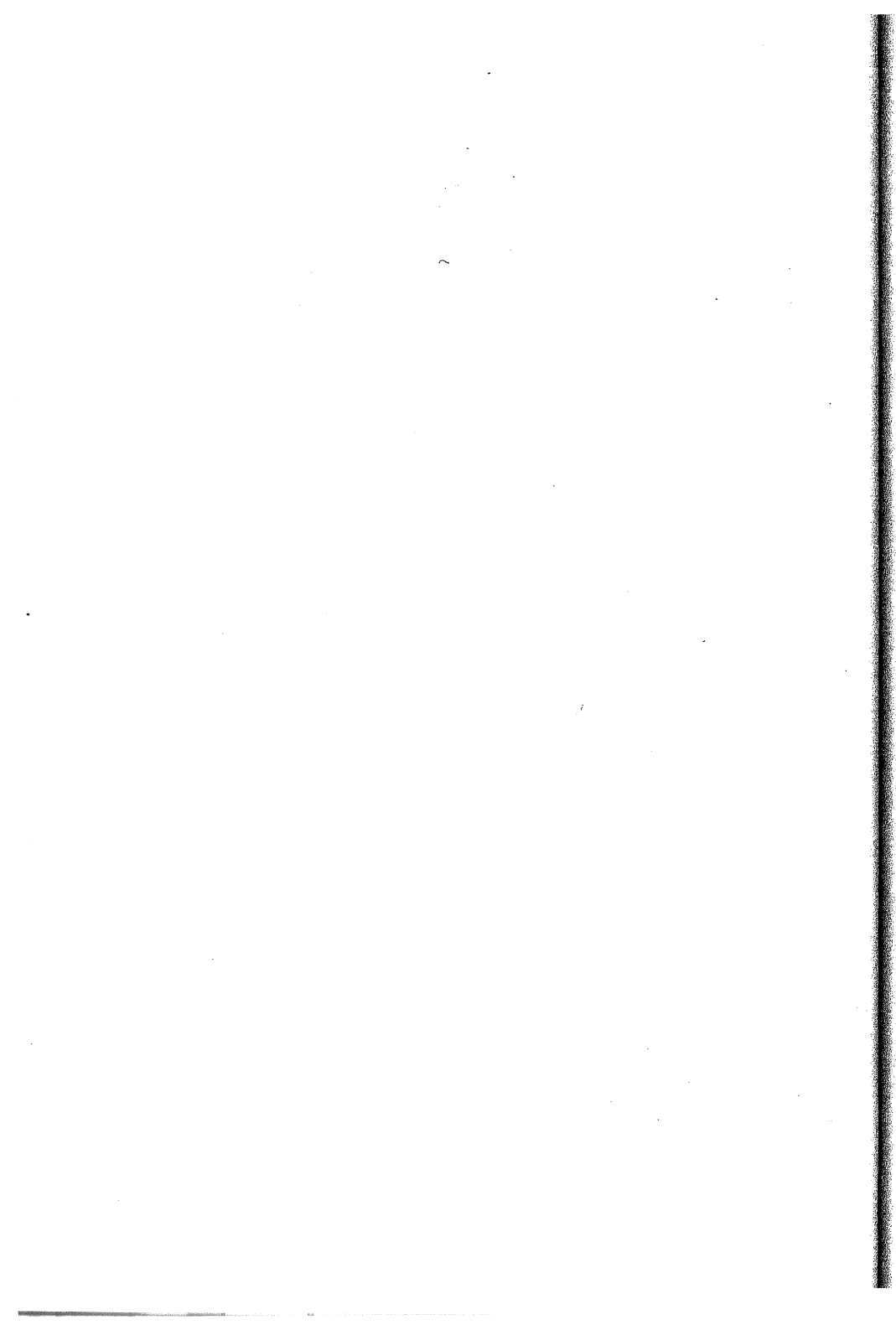
3.12 Exercises

1. Prove the addition law (item 4 of Proposition 3.2).
2. Show that if events E and F both have positive probability they cannot be both independent and exclusive.
3. Prove Proposition 3.13.
4. Suppose that a cell produces mRNA transcripts for a particular gene according to a Poisson process with a rate of 2 per second.
 - (a) What is the distribution of the number of transcripts produced in 5 seconds?
 - (b) What is the mean and variance of the number produced in 5 seconds?
 - (c) What is the probability that exactly the mean number will be produced?
 - (d) What is the distribution of the time until the first transcript?
 - (e) What is the mean and variance of the time until the first transcript?
 - (f) What is the probability that the time is more than 1 second?
 - (g) What is the distribution of the time until the third transcript?

- (h) What is the mean and variance of the time until the third transcript?
 - (i) What is the probability that the third transcript happens in the first second?
5. Derive from first principles the variance of the gamma distribution.

3.13 Further reading

For a more comprehensive and less terse introduction to probability models, classic texts such as Ross (2003) are ideal. However, most introductory statistics textbooks also cover much of the more basic material. More advanced statistics texts, such as Rice (1993) or Miller & Miller (2004), also cover this material in addition to statistical methodology that will be helpful in later chapters.



Stochastic simulation

4.1 Introduction

As we saw in the previous chapter, probability theory is a powerful framework for understanding random phenomena. For simple random systems, mathematical analysis alone can provide a complete description of all properties of interest. However, for the kinds of random systems that we are interested in (stochastic models of biochemical networks), mathematical analysis is not possible, and the systems are described as *analytically intractable*. However, this does not mean that it is not possible to understand such systems. With the aid of a computer, it is possible to *simulate* the time-evolution of the system dynamics. Stochastic simulation is concerned with the computer simulation of random (or stochastic) phenomena, and it is therefore essential to know something about this topic before proceeding further.

4.2 Monte-Carlo integration

The rationale for stochastic simulation can be summarised very easily: to understand a statistical model, simulate many realisations from it and study them. To make this more concrete, one way to understand stochastic simulation is to perceive it as a way of numerically solving the difficult integration problems that naturally arise in probability theory.

Suppose we have a (continuous) random variable X , with probability density function (PDF), $f(x)$, and we wish to evaluate $E(g(X))$ for some function $g(\cdot)$. We know that

$$E(g(X)) = \int_X g(x)f(x) dx,$$

and so the problem is one of integration. For simple $f(\cdot)$ and $g(\cdot)$ this integral might be straightforward to compute directly. On the other hand, in more complex scenarios, it is likely to be analytically intractable. However, if we can *simulate* realisations x_1, \dots, x_n of X , then we can form realisations of the random variable $g(X)$ as $g(x_1), \dots, g(x_n)$. Then, provided that the variance of $g(X)$ is finite, the laws of large numbers (Propositions 3.14 and 3.15) assure us that for large n we may approximate the integral by

$$E(g(X)) \simeq \frac{1}{n} \sum_{i=1}^n g(x_i).$$

In fact, even if we cannot simulate realisations of X , but can simulate realisations y_1, \dots, y_n of Y (a random variable with the same support as X), which has PDF

$h(\cdot)$, then

$$\begin{aligned} E(g(X)) &= \int_X g(x)f(x) dx \\ &= \int_X \frac{g(x)f(x)}{h(x)} h(x) dx \end{aligned}$$

and so $E(g(X))$ may be approximated by

$$E(g(X)) \simeq \frac{1}{n} \sum_{i=1}^n \frac{g(y_i)f(y_i)}{h(y_i)}.$$

This procedure is known as *importance sampling*, and it can be very useful when there is reasonable agreement between $f(\cdot)$ and $h(\cdot)$.

The above examples show that it is possible to compute summaries of interest such as expectations provided we have a mechanism for simulating realisations of random quantities. It turns out that doing this is a rather non-trivial problem. We begin first by thinking about the generation of uniform random quantities, and then move on to thinking about the generation of random quantities from other standard distributions.

4.3 Uniform random number generation

4.3.1 Discrete uniform random numbers

Most stochastic simulation begins with a uniform random number generator (Ripley 1987). Computers are essentially deterministic, so most random number generation strategies begin with the development of algorithms which produce a deterministic sequence of numbers that appears to be random. Consequently, such generators are often referred to as *pseudo-random number generators*.

Typically, a number theoretic method is used to generate an apparently random integer from 0 to $2^N - 1$ (often, $N = 16, 32$, or 64). The methods employed these days are often very sophisticated, as modern cryptographic algorithms rely heavily on the availability of "good" random numbers. A detailed discussion would be out of place in the context of this book, but a brief explanation of a very simple algorithm is perhaps instructive.

Linear congruential generators are the simplest class of algorithm used for pseudo-random number generation. The algorithm begins with a *seed* x_0 then generates new values according to the (deterministic) rule

$$x_{n+1} = (a x_n + b \pmod{2^N})$$

for carefully chosen a and b . Clearly such a procedure must return to a previous value at some point and then cycle indefinitely. However, if a "good" choice of a , b , and N are used, this deterministic sequence of numbers will have a very large cycle length (significantly larger than the number of random quantities one is likely to want to generate in a particular simulation study) and give every appearance of being uniformly random. One reasonable choice is

$$N = 59, b = 0, a = 13^{13}.$$

4.3.2 Standard uniform random numbers

Most computer programming languages have a built-in function for returning a pseudo-random integer. The technique used will be similar in principle to the linear congruential generator described above, but is likely to be more sophisticated (and therefore less straightforward for cryptographers to “crack”). Provided the maximum value is very large, the random number can be divided by the maximum value to give a pseudo-random $U(0, 1)$ number (and again, most languages will provide a function which does this). We will not be concerned with the precise mechanisms of how this is done, but rather take it as our starting point to examine how these uniform random numbers can be used to simulate more interesting distributions. Once we have a method for simulating $U \sim U(0, 1)$, we can use it in order to simulate random quantities from any distribution we like.

4.4 Transformation methods

Suppose that we wish to simulate realisations of a random variable X , with PDF $f(x)$, and that we are able to simulate $U \sim U(0, 1)$.

Proposition 4.1 (inverse distribution method) *If $U \sim U(0, 1)$ and $F(\cdot)$ is a valid invertible cumulative distribution function (CDF), then*

$$X = F^{-1}(U)$$

has CDF $F(\cdot)$.

Proof.

$$\begin{aligned} P(X \leq x) &= P(F^{-1}(U) \leq x) \\ &= P(U \leq F(x)) \\ &= F_U(F(x)) \\ &= F(x). \end{aligned} \quad (\text{as } F_U(u) = u)$$

□

So for a given $f(\cdot)$, assuming that we also compute the probability distribution function $F(\cdot)$ and its inverse $F^{-1}(\cdot)$, we can simulate a realisation of X using a single $U \sim U(0, 1)$ by applying the inverse CDF to the uniform random quantity.

4.4.1 Uniform random variates

Recall that the properties of the uniform distribution were discussed in Section 3.8. Given $U \sim U(0, 1)$, we can simulate $V \sim U(a, b)$ ($a < b$) in the obvious way, that is

$$V = a + (b - a)U.$$

We can justify this as V has CDF

$$F(v) = \frac{v - a}{b - a}, \quad a \leq v \leq b$$

and hence inverse

$$F^{-1}(u) = a + (b - a)u.$$

4.4.2 Exponential random variates

For discrete-event simulation of the time evolution of biochemical networks, it is necessary to have an efficient way to simulate exponential random quantities (Section 3.9). Fortunately, this is very straightforward, as the following result illustrates.

Proposition 4.2 *If $U \sim U(0, 1)$ and $\lambda > 0$, then*

$$X = -\frac{1}{\lambda} \log(U)$$

has an $Exp(\lambda)$ distribution.

Proof. Consider $X \sim Exp(\lambda)$. This has density $f(x)$ and distribution $F(x)$, where

$$f(x) = \lambda e^{-\lambda x}, \quad F(x) = 1 - e^{-\lambda x}, \quad x \geq 0$$

and so

$$F^{-1}(u) = -\frac{1}{\lambda} \log(1 - u), \quad 0 \leq u \leq 1.$$

The result follows as U and $1 - U$ clearly both have a $U(0, 1)$ distribution. \square

So, to simulate a realisation of X , simulate u from $U(0, 1)$, and then put

$$x = -\frac{1}{\lambda} \log(u).$$

4.4.3 Scaling

It is worth noting the scaling issues in the above example, as these become more important for distributions which are more difficult to simulate from. If $U \sim U(0, 1)$, then $Y = -\log U \sim Exp(1)$. The parameter, λ , of an exponential distribution is a *scale parameter* because we can obtain exponential variates with other parameters from a variate with a unit parameter by a simple linear scaling. That is, if

$$Y \sim Exp(1)$$

then

$$X = \frac{1}{\lambda} Y \sim Exp(\lambda).$$

In general, we can spot *location* and *scale* parameters in a distribution as follows. If Y has PDF $f(y)$ and CDF $F(y)$, and $X = aY + b$, then X has CDF

$$\begin{aligned} F_X(x) &= P(X \leq x) \\ &= P(aY + b \leq x) \\ &= P\left(Y \leq \frac{x - b}{a}\right) \\ &= F\left(\frac{x - b}{a}\right) \end{aligned}$$

and PDF

$$f_X(x) = \frac{1}{a} f\left(\frac{x-b}{a}\right).$$

Parameters like a are scale parameters, and parameters like b are location parameters.

4.4.4 Gamma random variates

The inverse transformation method is ideal for distributions with tractable CDF and inverse CDF, but the gamma distribution (Section 3.11) is a good example of a distribution where it is not practical to compute an analytic expression for the inverse of the CDF, and so it is not possible to use the inverse transformation method directly. In this case it is possible to use numerical techniques to compute the CDF and its inverse, but this is not very efficient. More sophisticated techniques are therefore required in general. One way to simulate $X \sim \Gamma(n, \lambda)$ random variates for integer n is to use the fact that if

$$Y_i \sim \text{Exp}(\lambda),$$

and the Y_i are independent, then

$$X = \sum_{i=1}^n Y_i \sim \Gamma(n, \lambda).$$

So, just simulate n exponential random variates and add them up.

The first parameter of the gamma distribution (here n) is known as the *shape* parameter, and the second (here λ) is known as the *scale* parameter. It is important to understand that the second parameter is a scale parameter, because many gamma generation algorithms work by first generating gamma variables with arbitrary shape but unit scale. Once this has been achieved, it can then easily be rescaled to give the precise gamma variate required. This can be done easily because the $\Gamma(\alpha, \beta)$ PDF is

$$f_X(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}, \quad x > 0$$

and so the $\Gamma(\alpha, 1)$ PDF is

$$f_Y(y) = \frac{1}{\Gamma(\alpha)} y^{\alpha-1} e^{-y}, \quad y > 0.$$

We can see that if $Y \sim \Gamma(\alpha, 1)$, then $X = Y/\beta \sim \Gamma(\alpha, \beta)$ because

$$f_X(x) = \beta f_Y(\beta x).$$

Consequently, the CDFs must be related by

$$F_X(x) = F_Y(\beta x).$$

Techniques for efficiently generating gamma variates with arbitrary shape parameter are usually based on rejection techniques (to be covered later). Note, however, that for shape parameters which are an integer multiple of 0.5, use can be made of the fact that $\chi_n^2 = \Gamma(n/2, 1/2)$.* So, given a technique for generating χ^2 quantities, gamma

* We did not discuss the χ^2 (or *chi-square*) distribution in Chapter 3, but it is usually defined as the

variates with a shape parameter that is an integer multiple of $1/2$ can be generated by a simple rescaling.

4.4.5 Normal random variates

The ability to efficiently simulate normal (or Gaussian) random quantities is of crucial importance for stochastic simulation. However, the inverse transformation method is not really appropriate due to the analytic intractability of the CDF and inverse CDF of the normal distribution, so another method needs to be found. Note that all we need is a technique for simulating $Z \sim N(0, 1)$ random variables, as then $X = \mu + \sigma Z \sim N(\mu, \sigma^2)$. Also note that standard normal random variables can be used to generate χ^2 random variables. If $Z_i \sim N(0, 1)$ and the Z_i are independent, then

$$C = \sum_{i=1}^n Z_i^2$$

has a χ_n^2 distribution.

CLT-based method

One simple way to generate normal random variables is to make use of the Central Limit Theorem (Theorem 3.3). Consider

$$Z = \sum_{i=1}^{12} U_i - 6$$

where $U_i \sim U(0, 1)$. Clearly $E(Z) = 0$ and $\text{Var}(Z) = 1$, and by the Central Limit Theorem (CLT), Z is approximately normal. However, this method is not exact. For example, Z only has support on $[-6, 6]$ and is poorly behaved in the extreme tails. However, for $Z \sim N(0, 1)$ we have $P(|Z| > 6) \simeq 2 \times 10^{-9}$, and so the truncation is not much of a problem in practice, and this method is good enough for many purposes. The main problem with this method is that for “good” random number generators, simulating uniform random numbers is quite slow, and so simulating 12 in order to get one normal random quantity is undesirable.

Box-Muller method

A more efficient (and “exact”) method for generating normal random variates is the following. Simulate

$$\Theta \sim U(0, 2\pi), \quad R^2 \sim \text{Exp}(1/2)$$

independently. Then

$$X = R \cos \Theta, \quad Y = R \sin \Theta$$

sum of squares of independent standard normal random quantities, which leads directly to the simulation method discussed in the next section. That the χ^2 distribution is a special case of the gamma distribution is demonstrated in Miller & Miller (2004).

are two independent standard normal random variables. It is easier to show this the other way around.

Proposition 4.3 *If X and Y are independent standard normal random quantities, and they are regarded as the Cartesian coordinates (X, Y) of a 2-d random variable, then the polar coordinates of the variable (R, Θ) are independent with $R^2 \sim \text{Exp}(1/2)$ and $\Theta \sim U(0, 2\pi)$.*

Proof. Suppose that $X, Y \sim N(0, 1)$ and that X and Y are independent. Then

$$f_{X,Y}(x, y) = \frac{1}{2\pi} \exp\{-(x^2 + y^2)/2\}.$$

Put

$$X = R \cos \Theta \quad \text{and} \quad Y = R \sin \Theta.$$

Then,

$$\begin{aligned} f_{R,\Theta}(r, \theta) &= f_{X,Y}(x, y) \left| \frac{\partial(x, y)}{\partial(r, \theta)} \right| \\ &= \frac{1}{2\pi} e^{-r^2/2} \begin{vmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{vmatrix} \\ &= \frac{1}{2\pi} \times r e^{-r^2/2}. \end{aligned}$$

So, Θ and R are independent, $\Theta \sim U(0, 2\pi)$, and $f_R(r) = r e^{-r^2/2}$. It is then easy to show (using Proposition 3.13) that $R^2 \sim \text{Exp}(1/2)$. \square

Note that here we have used a couple of results that were not covered directly in Chapter 3: namely, that independence of continuous random quantities corresponds to factorisation of the density, and also that the random variable transformation formula (Proposition 3.13), generalises to higher dimensions (where the ‘‘Jacobian’’ of the transformation is now the *determinant* of the matrix of partial derivatives). Note also that algorithms for efficiently simulating normal random quantities are built in to most scientific and numerical libraries, and it is preferable to use these where available.

4.5 Lookup methods

The so-called lookup method is just the discrete version of the inverse transformation method. However, as it looks a little different at first sight, it is worth examining in detail. Suppose one is interested in simulating a discrete random quantity X with outcome space $S = \{0, 1, 2, \dots\}$, or a subset thereof. Put $p_k = P(X = k)$, $k = 0, 1, 2, \dots$, the PMF of X . Some of the p_k may be zero. Indeed, if X is finite, then all but finitely many of the p_k will be zero. Next define

$$q_k = P(X \leq k) = \sum_{i=0}^k p_i.$$

Realisations of X can be generated by first simulating $U \sim U(0, 1)$ and then putting

$$X = \min\{k | q_k \geq U\}.$$

This works because then by definition we must have $q_{k-1} < U \leq q_k$, and so

$$P(X = k) = P(U \in (q_{k-1}, q_k]) = q_k - q_{k-1} = p_k,$$

as required. In practice, for a one-off simulation, first simulate a U and then compute a running sum of the p_i . At each stage check to see if the running sum is as big as U . If it is, return the index, k of the final p_k . If many realisations are required, more efficient strategies are possible, but in the context of discrete-event simulation of biochemical networks, a one-off simulation is usually all that is required. As we will see, typically one picks one of a finite set of reactions using a lookup method and chooses a time to wait until that reaction occurs by simulating an exponential random quantity.

4.6 Rejection samplers

We have now examined a range of useful techniques for random number generation, but we still have not seen a technique that can be used for general continuous random quantities where the inverse CDF cannot easily be computed. In this situation, rejection samplers are most often used.

Proposition 4.4 (uniform rejection method) *Suppose we want to simulate from $f(x)$ with (finite) support on (a, b) , and that $f(x) \leq m$, $\forall x \in (a, b)$. Then consider simulating*

$$X \sim U(a, b) \quad \text{and} \quad Y \sim U(0, m).$$

Accept X if $Y < f(X)$, otherwise reject and try again. Then the accepted X values have PDF $f(x)$.

Intuitively, we can see that this will work because it has the effect of scattering points uniformly over the region bounded by the PDF and the x -axis.

Proof. Call the acceptance region A , and the accepted value \tilde{X} . Then

$$\begin{aligned} F_{\tilde{X}}(x) &= P(\tilde{X} \leq x) \\ &= P(X \leq x | (X, Y) \in A) \\ &= \frac{P((X \leq x) \cap ((X, Y) \in A))}{P((X, Y) \in A)} \\ &= \frac{\int_a^b P((X \leq x) \cap ((X, Y) \in A) | X = z) \times \frac{1}{b-a} dz}{\int_a^b P((X, Y) \in A | X = z) \times \frac{1}{b-a} dz} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{b-a} \int_a^x P((X, Y) \in A | X = z) dz \\
&= \frac{1}{b-a} \int_a^b P((X, Y) \in A | X = z) dz \\
&= \frac{\int_a^x \frac{f(z)}{m} dz}{\int_a^b \frac{f(z)}{m} dz} \\
&= \int_a^x f(z) dz \\
&= F(x).
\end{aligned}$$

□

So, in summary, we simulate a value x uniformly from the support of X and accept this value with probability $f(x)/m$, otherwise we reject and try again. Obviously the efficiency of this method depends on the overall proportion of candidate points that are accepted. The actual acceptance probability for this method is

$$\begin{aligned}
P(\text{Accept}) &= P((X, Y) \in A) \\
&= \int_a^b P((X, Y) \in A | X = x) \times \frac{1}{b-a} dx \\
&= \int_a^b \frac{f(x)}{m} \times \frac{1}{b-a} dx \\
&= \frac{1}{m(b-a)} \int_a^b f(x) dx \\
&= \frac{1}{m(b-a)}.
\end{aligned}$$

If this acceptance probability is very low, the procedure will be very inefficient, and a better procedure should be sought — the *envelope method* is one possibility.

4.6.1 Envelope method

Once we have established that scattering points uniformly over the region bounded by the density and the x -axis generates x -values with the required distribution, we can extend it to distributions with infinite support and make it more efficient by choosing our *enveloping* region more carefully.

Suppose that we wish to simulate X with PDF $f(\cdot)$, but that we can already simulate values of Y (with the same support as X), which has PDF $g(\cdot)$. Suppose further that there exists some constant a such that

$$f(x) \leq a g(x), \quad \forall x.$$

That is, a is an upper bound for $f(x)/g(x)$. Note also that $a \geq 1$, as both $f(x)$ and $g(x)$ integrate to 1.

Consider the following algorithm. Draw $Y = y$ from $g(\cdot)$, and then $U = u \sim U(0, a g(y))$. Accept y as a simulated value of X if $u < f(y)$, otherwise reject and try again. This works because it distributes points uniformly over a region covering $f(x)$, and then only keeps points in the required region (under $f(x)$):

$$\begin{aligned}
 P(\tilde{X} \leq x) &= P(Y \leq x | U \leq f(Y)) \\
 &= \frac{P([Y \leq x] \cap [U \leq f(Y)])}{P(U \leq f(Y))} \\
 &= \frac{\int_{-\infty}^{\infty} P([Y \leq x] \cap [U \leq f(Y)] | Y = y) g(y) dy}{\int_{-\infty}^{\infty} P(U \leq f(Y) | Y = y) g(y) dy} \\
 &= \frac{\int_{-\infty}^x P(U \leq f(Y) | Y = y) g(y) dy}{\int_{-\infty}^{\infty} P(U \leq f(Y) | Y = y) g(y) dy} \\
 &= \frac{\int_{-\infty}^x \frac{f(y)}{a g(y)} g(y) dy}{\int_{-\infty}^{\infty} \frac{f(y)}{a g(y)} g(y) dy} \\
 &= \frac{\int_{-\infty}^x \frac{f(y)}{a} dy}{\int_{-\infty}^{\infty} \frac{f(y)}{a} dy} \\
 &= \int_{-\infty}^x f(y) dy \\
 &= F(x).
 \end{aligned}$$

To summarise, just simulate a proposed value from $g(\cdot)$ and accept this with probability $f(y)/[a g(y)]$, otherwise reject and try again. The accepted values will have PDF $f(\cdot)$.

Obviously, this method will work well if the overall acceptance rate is high, but not otherwise. The overall acceptance probability can be computed as

$$\begin{aligned}
 P(U < f(Y)) &= \int_{-\infty}^{\infty} P(U < f(Y) | Y = y) g(y) dy \\
 &= \int_{-\infty}^{\infty} \frac{f(y)}{a g(y)} g(y) dy \\
 &= \int_{-\infty}^{\infty} \frac{f(y)}{a} dy \\
 &= \frac{1}{a}.
 \end{aligned}$$

Consequently, we want a to be as small as possible (that is, as close as possible to 1). What “small enough” means is context-dependent, but generally speaking, if $a > 10$, the envelope is not adequate — too many points will be rejected, so a better envelope needs to be found. If this is not practical, then an entirely new approach is required.

4.7 The Poisson process

Consider a Poisson process with rate λ defined on the interval $[0, T]$. As we know, the inter-event times are $Exp(\lambda)$, and so we have a very simple algorithm for simulating realisations of the process. Initialise the process at time zero. Then simulate $t_1 \sim Exp(\lambda)$, the time to the first event, and put $X_1 = t_1$. Next simulate $t_2 \sim Exp(\lambda)$, the time from the first to the second event, and put $X_2 = X_1 + t_2$. At step k , simulate $t_k \sim Exp(\lambda)$, the time from the $k-1$ to k th event, and put $X_k = X_{k-1} + t_k$. Repeat until $X_k > T$, and keep the X_1, X_2, \dots as the realisation of the process.

4.8 Using the statistical programming language, R

4.8.1 Introduction

R is a programming language for data analysis and statistics. It is a completely free open-source software application, and very widely used by professional statisticians. It is also very popular in certain application areas, including bioinformatics. R is a dynamically typed interpreted language, and it is typically used interactively. It has many built-in functions and libraries and is extensible, allowing users to define their own functions and procedures using R, C, or Fortran. It also has a simple object system. R is a particularly convenient environment to use for stochastic simulation, visualisation, and analysis. It will therefore be used in the forthcoming chapters in order to provide concrete illustrations of the theory being described. It is strongly recommended that readers unfamiliar with R download and install it, then work through this mini-tutorial.

4.8.2 Vectors

Vectors are a fundamental concept in R, as many functions operate on and return vectors, so it is best to master these as soon as possible. For the technically inclined, in R, a (numeric) vector is an object consisting of a one-dimensional array of scalars.

```
> rep(1, 10)
[1] 1 1 1 1 1 1 1 1 1 1
>
```

Here `rep` is a function that returns a vector (here, 1 repeated 10 times). You can get documentation for `rep` by typing

```
> ?rep
```

You can assign any object (including vectors) using the assignment operator `<-` (or `=`), and combine vectors and scalars with the `c` function.

```

> a<-rep(1,10)
> b<-1:10
> c(a,b)
 [1] 1 1 1 1 1 1 1 1 1 1 1 2 3 4 5 6
[17] 7 8 9 10
> a+b
 [1] 2 3 4 5 6 7 8 9 10 11
> a+2*b
 [1] 3 5 7 9 11 13 15 17 19 21
> a/b
 [1] 1.0000000 0.5000000 0.3333333 0.2500000 0.2000000
 [6] 0.1666667 0.1428571 0.1250000 0.1111111 0.1000000
> c(1,2,3)
 [1] 1 2 3
>

```

Note that arithmetic operations act element-wise on vectors. To look at any object (function or data), just type its name.

```

> b
 [1] 1 2 3 4 5 6 7 8 9 10

```

To list all of your objects, use `ls()`. Note that because of the existence of a function called `c` (and another called `t`) it is best to avoid using these as variable names in user-defined functions (this is a common source of bugs).

Vectors can be “sliced” very simply:

```

> d<-c(3,5,3,6,8,5,4,6,7)
> d
 [1] 3 5 3 6 8 5 4 6 7
> d[4]
 [1] 6
> d[2:4]
 [1] 5 3 6
> d<7
 [1] TRUE TRUE TRUE TRUE FALSE TRUE TRUE TRUE FALSE
> d[d<7]
 [1] 3 5 3 6 5 4 6
>

```

Vectors can be sorted and randomly sampled. The following command generates some lottery numbers.

```

> sort(sample(1:49,6))
 [1] 2 17 23 24 25 35

```

Get help (using `?`) on `sort` and `sample` to see how they work. R is also good at stochastic simulation of vectors of quantities from standard probability distributions, so

```

> rpois(20,2)
[1] 1 4 0 2 2 3 3 3 3 4 2 2 4 2 1 3 4 1 2 1
generates 20 Poisson random variates with mean 2, and

> rnorm(5,1,2)
[1] -0.01251322 -0.03181018 0.30426031 3.24302197 -2.04370284
generates 5 normal random variates with mean 1 and standard deviation (not variance) 2. There are lots of functions that act on vectors.

> x<-rnorm(50,2,3)
> x
[1] 2.04360635 5.01113289 -1.52215979 -0.19789766 1.41945311 -0.08850784
[7] -0.91161025 3.47199019 6.13447194 4.62796165 0.07600234 -2.99687943
[13] 1.75153104 8.55000833 3.11921624 3.38411717 3.86860456 0.29103619
[19] 1.25823419 3.88427191 -0.77722215 -0.57774833 2.99937058 4.29042603
[25] 6.10597239 2.83832381 3.73618138 4.12999252 6.23009274 1.07251421
[31] -0.19645150 1.77581296 2.08783542 1.62948606 2.74911850 0.44028844
[37] 1.80996899 1.86436309 0.29372974 2.37077354 1.54285955 4.40098545
[43] -3.01913118 -0.23174209 3.58252631 5.18954147 3.61988373 4.08815220
[49] 6.30878505 4.56744882
> mean(x)
[1] 2.361934
> length(x)
[1] 50
> var(x)
[1] 6.142175
> summary(x)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-3.0190 0.3304  2.2290  2.3620  4.0370  8.5500
> cumsum(x)
 [1] 2.043606  7.054739  5.532579  5.334682  6.754135  6.665627
 [7] 5.754017  9.226007 15.360479 19.988441 20.064443 17.067563
[13] 18.819095 27.369103 30.488319 33.872436 37.741041 38.032077
[19] 39.290311 43.174583 42.397361 41.819613 44.818983 49.109409
[25] 55.215382 58.053705 61.789887 65.919879 72.149972 73.222486
[31] 73.026035 74.801848 76.889683 78.519169 81.268288 81.708576
[37] 83.518545 85.382908 85.676638 88.047412 89.590271 93.991257
[43] 90.972125 90.740383 94.322910 99.512451 103.132335 107.220487
[49] 113.529272 118.096721
> sum(x)
[1] 118.0967
>

```

Note that the function `var` computes and returns the *sample variance* of a data vector (and the function `sd` returns the square root of this). The sample variance, usually denoted by s^2 , is an estimator of the population variance in the same sense that the sample mean, \bar{X} , is an estimator of the population mean, μ . It is defined by

$$s^2 = \frac{1}{n-1} \left(\left[\sum_{i=1}^n X_i^2 \right] - n\bar{X}^2 \right).$$

4.8.3 Plotting

R has many functions for data visualisation and for producing publication quality plots — running `demo(graphics)` will give an idea of some of the possibilities. Some more basic commands are given below. Try them in turn and see what they do.

```

> plot(1:50, cumsum(x))
> plot(1:50, cumsum(x), type="l", col="red")

```



```

> plot(x, 0.6*x+rnorm(50, 0.3))
> curve(0.6*x, -5, 10, add=TRUE)
> hist(x)
> hist(x, 20)
> hist(x, freq=FALSE)
> curve(dnorm(x, 2, 3), -5, 10, add=TRUE)
> boxplot(x, 2*x)
> barplot(d)
>

```

Study the Help file for each of these commands to get a feel for the way each can be customised.

4.8.4 User-defined functions

R is a full programming language, and before long, you are likely to want to add your own functions. Consider the following declaration.

```

rchi<-function(n, p=2) {
  x<-matrix(rnorm(n*p), nrow=n, ncol=p)
  y<-x*x
  as.vector(y %*% rep(1, p))
}

```

The first line declares the object `rchi` to be a function with two arguments, `n` and `p`, the second of which will default to a value of 2 if not specified. Then everything between `{` and `}` is the function body, which can use the variables `n` and `p` as well as any globally defined objects. The second line declares a local variable `x` to be a matrix with `n` rows and `p` columns, whose elements are standard normal random variables. The next line forms a new matrix `y` whose elements are the squares of the elements in `x`. The last line computes the matrix-vector product of `y` and a vector of `p` ones, then coerces the resulting `n` by 1 matrix into a vector. The result of the last line of the function body is the return result of the function. In fact, this function provides a fairly efficient way of simulating Chi-squared random quantities with `p` degrees of freedom, but that is not particularly important. The function is just another R object, and hence can be viewed by entering `rchi` on a line by itself. It can be edited by doing `fix(rchi)`. The function can be called just like any other, so

```

> rchi(10, 3)
[1] 1.847349 5.590369 3.994036 4.243734 2.104224
[6] 1.027634 1.119508 6.653095 5.660968 5.384954
> rchi(10)
[1] 0.09356735 3.63633129 1.34073206 1.79412360
[5] 1.46038656 2.67362870 0.50413958 6.04307710
[9] 1.03116671 1.39662895
>

```

generates 10 chi-squared random variates with 3 and 2 degrees of freedom, respectively.

4.8.5 Reading and writing data

Of course, in order to use R for data analysis, it is necessary to be able to read data into R from other sources. It is often also desirable to be able to output data from R in a format that can be read by other applications. Unsurprisingly, R has a range of functions for accomplishing these tasks, but we shall just look here at the two simplest.

The simplest way to get data into R is to read a list of numbers from a text file using the `scan` command. This is most easily illustrated by first writing some data out to a text file and then reading it back into R. A vector of numbers can be output with a command like

```
> write(x, "scandata.txt")
>
```

Then, to load data from the file `scandata.txt`, use a command like

```
> x<-scan("scandata.txt")
Read 50 items
>
```

In general, you may need to use the `getwd` and `setwd` commands to inspect and change the working directory that R is using.

More often, we will be concerned with loading tabular data output from a spreadsheet or database or even another statistics package. R copes best with whitespace-separated data, but can be persuaded to read other formats with some effort. The key command here is `read.table` (and the corresponding output command `write.table`). So, suppose that `mytable.txt` is a plain text file containing the following lines.

```
"Name" "Shoe size" "Height"
"Fred" 9 170
"Jim" 10 180
"Bill" 9 185
"Jane" 7 175
"Jill" 6 170
"Janet" 8 180
```

To read this data into R, do

```
> mytab<-read.table("mytable.txt", header=TRUE)
> mytab
  Name Shoe.size Height
1  Fred         9    170
2   Jim        10    180
3  Bill         9    185
4  Jane         7    175
5  Jill         6    170
6 Janet         8    180
>
```

Note that R does contain some primitive functions for editing data frames like this (and other objects), so

```
> mytabnew<-edit(mytab)
```

will pop up a simple editor for `mytab`, and on quitting, the edited version will be stored in `mytabnew`. Data frames like `mytab` are a key object type in R, and tend to be used often. Here are some ways to interact with data frames.

```
> mytab$Height
[1] 170 180 185 175 170 180
> mytab[,2]
[1] 9 10 9 7 6 8
> plot(mytab[,2],mytab[,3])
> mytab[4,]
  Name Shoe.size Height
4 Jane          7    175
> mytab[5,3]
[1] 170
> mytab[mytab$Name=="Jane",]
  Name Shoe.size Height
4 Jane          7    175
> mytab$Height[mytab$Shoe.size > 8]
[1] 170 180 185
>
```

Also see the Help on `source` and `dump` for input and output of R objects of other sorts.

4.8.6 Further reading for R

One of the great things about R is that it comes with a great deal of excellent documentation (from the Comprehensive R Archive Network — CRAN). The next thing to work through is the official *Introduction to R*, which covers more material in more depth than this very quick introduction. Further pointers are given on this book's website.

4.9 Analysis of simulation output

This chapter finishes with the analysis of a random quantity using stochastic simulation. Suppose interest lies in $Y = \exp(X)$, where $X \sim N(2, 1)$. In fact, Y is a standard distribution (it is log-normally distributed), and all of the interesting properties of Y can be derived directly, analytically. However, we will suppose that we are not able to do this, and instead study Y using stochastic simulation, using only the ability to simulate normal random quantities. Using R, samples from Y can be generated as follows:

```
> x<-rnorm(10000,2,1)
> y<-exp(x)
```

The variable Y has a long-tailed distribution, which can be visualised with

```
> hist(y,breaks=50,freq=FALSE)
```

and a version of this plot is shown in Figure 4.1. The samples can also be used for computing summary statistics. Basic sample statistics can be obtained with

```
> summary(y)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 0.134  3.733   7.327  12.220  14.490  728.700
> sd(y)
[1] 17.21968
>
```

and the sample mean, median, and quartiles provide estimates of the true population quantities. Focussing on the sample mean, \bar{x} , the value obtained here (12.220) is an estimate of the population mean. Of course, should the experiment be repeated, the estimate will be different. However, we can use the fact that $\bar{X} \sim N(\mu, \sigma^2/n)$ (by the CLT), to obtain $Z \sim N(0, 1)$, where $Z \sim \sqrt{n}(\bar{X} - \mu)/\sigma$. Then since $P(|Z| < 2) \simeq 0.95$, we have

$$P(|\bar{X} - \mu| < 2\sigma/\sqrt{n}) \simeq 0.95,$$

and substituting in n and the estimated value of σ (17.21968), we get

$$P(|\bar{X} - \mu| < 0.344) \simeq 0.95.$$

We therefore expect that \bar{X} is likely to be within 0.344 of μ (though careful readers will have noted that the conditioning here is really the wrong way around). In fact, in this particular case, the true mean of this distribution can be calculated analytically as $\exp(2.5) = 12.18249$, which is seen to be consistent with the simulation estimate. Thus, in more complex examples where the true population properties are not available, estimated sample quantities can be used as a substitute, provided that enough samples can be generated to keep the “Monte-Carlo error” to a minimum.

4.10 Exercises

1. The random variable X has PDF

$$f(x) = \begin{cases} \sin(x), & 0 \leq x \leq \pi/2, \\ 0, & \text{otherwise.} \end{cases}$$

- Derive a transformation method for simulating values of X based on $U(0, 1)$ random variates.
- Derive a uniform rejection method for simulating values from X . What is the acceptance probability?
- Derive an envelope rejection method for simulating values of X based on a proposal with density

$$g(x) = \begin{cases} kx, & 0 \leq x \leq \pi/2, \\ 0, & \text{otherwise,} \end{cases}$$

for some fixed k . You should use the fact that $\sin(x) \leq x, \forall x \geq 0$. What is the acceptance probability?

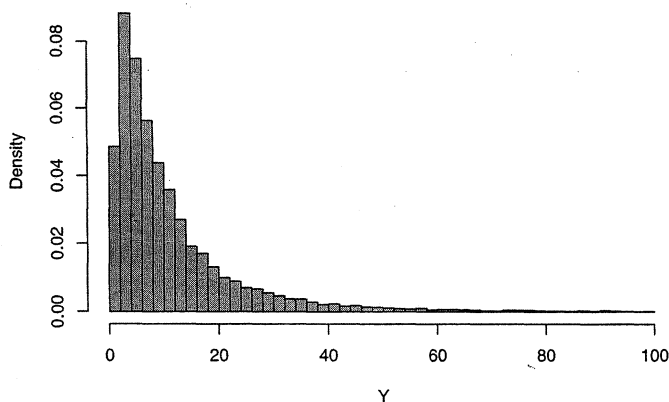


Figure 4.1 Density of $Y = \exp(X)$, where $X \sim N(2, 1)$.

2. If you have not already done so, follow the links from this book's website and download and install R. Work through the mini-tutorial from this chapter.
3. Download the official *Introduction to R* from CRAN (linked from this book's website) and work through the first half at least.
4. Write your own function `myrexp`, which does the same as `rexp`, but does not rely on the built-in version.
5. Write a function to simulate normal random quantities using the CLT method. Use plots and summary statistics to compare the distribution you obtain with those of the built-in `rnorm` function (which is exact).
6. Write your own function to simulate $\Gamma(3, 2)$ random quantities (and again, compare with the built-in version). See if you can also write your own function to simulate $\Gamma(3.5, 5)$ random quantities.
7. Obtain Monte-Carlo solutions to the problems posed in Exercise 4 from Chapter 3.

4.11 Further reading

There are several good introductory texts on stochastic simulation, including Morgan (1984) and Ripley (1987). Devroye (1986) is an excellent reference work on the subject. The standard reference for R is the R Development Core Team (2005).

Markov processes

5.1 Introduction

We now have a grounding in elementary probability theory and an understanding of stochastic simulation. The only remaining theory required before studying the dynamics of genetic and biochemical networks (and chemical kinetics more generally) is an introduction to the theory of stochastic processes. A stochastic process is a random variable (say, the state of a biochemical network) which evolves through time. The state may be continuous or discrete, and it can evolve through time in a discrete or continuous way. A Markov process is a stochastic process which possesses the property that the future behaviour depends only on the current state of the system. Put another way, given information about the current state of the system, information about the past behaviour of the system is of no help in predicting the time-evolution of the process. It turns out that Markov processes are particularly amenable to both theoretical and computational analyses. Fortunately, the dynamic behaviour of biochemical networks can be effectively modelled by a Markov process, so familiarity with Markov processes is sufficient for studying many problems that arise naturally in systems biology.

5.2 Finite discrete time Markov chains

5.2.1 Introduction

The set $\{\theta^{(t)} | t = 0, 1, 2, \dots\}$ is a *discrete time stochastic process*. The *state space* S is such that $\theta^{(t)} \in S, \forall t$ and may be discrete or continuous.

A (first order) *Markov chain* is a stochastic process with the property that the future states are independent of the past states given the present state. Formally, for $A \subseteq S, n = 0, 1, 2, \dots$, we have

$$\begin{aligned} P(\theta^{(n+1)} \in A | \theta^{(n)} = x, \theta^{(n-1)} = x_{n-1}, \dots, \theta^{(0)} = x_0) \\ = P(\theta^{(n+1)} \in A | \theta^{(n)} = x), \quad \forall x, x_{n-1}, \dots, x_0 \in S. \end{aligned}$$

The past states provide no information about the future state if the present state is known. The behaviour of the chain is therefore determined by $P(\theta^{(n+1)} \in A | \theta^{(n)} = x)$. In general this depends on n, A and x . However, if there is no n dependence, so that

$$P(\theta^{(n+1)} \in A | \theta^{(n)} = x) = P(x, A), \quad \forall n,$$

then the Markov chain is said to be (time) *homogeneous*, and the *transition kernel*,

$P(x, A)$ determines the behaviour of the chain. Note that $\forall x \in S$, $P(x, \cdot)$ is a probability measure over S .

5.2.2 Notation

When dealing with discrete state spaces, it is easier to write

$$P(x, \{y\}) = P(x, y) = P(\theta^{(n+1)} = y | \theta^{(n)} = x).$$

In the case of a finite discrete state space, $S = \{x_1, \dots, x_r\}$, we can write $P(\cdot, \cdot)$ as a matrix

$$P = \begin{pmatrix} P(x_1, x_1) & \cdots & P(x_1, x_r) \\ \vdots & \ddots & \vdots \\ P(x_r, x_1) & \cdots & P(x_r, x_r) \end{pmatrix}.$$

The matrix P is a *stochastic matrix*.

Definition 5.1 A real $r \times r$ matrix P is said to be a stochastic matrix if its elements are all non-negative and its rows sum to 1.

Proposition 5.1 The product of two stochastic matrices is another stochastic matrix. Every eigenvalue λ of a stochastic matrix satisfies $|\lambda| \leq 1$.* Also, every stochastic matrix has at least one eigenvalue equal to 1.

The proof of this proposition is straightforward and left to the end-of-chapter exercises.

Suppose that at time n , we have

$$P(\theta^{(n)} = x_1) = \pi^{(n)}(x_1)$$

$$P(\theta^{(n)} = x_2) = \pi^{(n)}(x_2)$$

$$\vdots \quad \quad \quad \vdots$$

$$P(\theta^{(n)} = x_r) = \pi^{(n)}(x_r).$$

We can write this as an r -dimensional row vector

$$\pi^{(n)} = (\pi^{(n)}(x_1), \pi^{(n)}(x_2), \dots, \pi^{(n)}(x_r)).$$

The probability distribution at time $n + 1$ can be computed using Theorem 3.1, as

$$\begin{aligned} P(\theta^{(n+1)} = x_1) &= P(x_1, x_1) \pi^{(n)}(x_1) + P(x_2, x_1) \pi^{(n)}(x_2) + \\ &\quad \cdots + P(x_r, x_1) \pi^{(n)}(x_r), \end{aligned}$$

* When considering a matrix A , the vector v is called a (column) *eigenvector* of A if $Av = \lambda v$ for some real number λ , which is known as an *eigenvalue* of A , corresponding to the eigenvector v . The row eigenvectors of A are the column eigenvectors of A' . Although row and column eigenvectors are different, the corresponding eigenvalues are the same. That is, A and A' have the same (column) eigenvalues.

and similarly for $P(\theta^{(n+1)} = x_2)$, $P(\theta^{(n+1)} = x_3)$, etc. We can write this in matrix form as

$$(\pi^{(n+1)}(x_1), \pi^{(n+1)}(x_2), \dots, \pi^{(n+1)}(x_r)) = (\pi^{(n)}(x_1), \pi^{(n)}(x_2), \dots, \pi^{(n)}(x_r)) \\ \times \begin{pmatrix} P(x_1, x_1) & \cdots & P(x_1, x_r) \\ \vdots & \ddots & \vdots \\ P(x_r, x_1) & \cdots & P(x_r, x_r) \end{pmatrix}$$

or equivalently

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

So,

$$\begin{aligned} \pi^{(1)} &= \pi^{(0)}P \\ \pi^{(2)} &= \pi^{(1)}P = \pi^{(0)}PP = \pi^{(0)}P^2 \\ \pi^{(3)} &= \pi^{(2)}P = \pi^{(0)}P^2P = \pi^{(0)}P^3 \\ &\vdots = \vdots \\ \pi^{(n)} &= \pi^{(0)}P^n. \end{aligned}$$

That is, the initial distribution $\pi^{(0)}$, together with the transition matrix P , determine the probability distribution for the state at all future times. Further, if the one-step transition matrix is P , then the n -step transition matrix is P^n . Similarly, if the m -step transition matrix is P^m and the n -step transition matrix is P^n , then the $(m+n)$ -step transition matrix is $P^mP^n = P^{m+n}$. The set of linear equations corresponding to this last statement are known as the *Chapman-Kolmogorov equations*.

5.2.3 Stationary distributions

A distribution π is said to be a *stationary distribution* of the homogeneous Markov chain governed by the transition matrix P if

$$\pi = \pi P. \quad (5.1)$$

Note that π is a row eigenvector of the transition matrix, with corresponding eigenvalue equal to 1. It is also a fixed point of the linear map induced by P . The stationary distribution is so-called because if at some time n , we have $\pi^{(n)} = \pi$, then $\pi^{(n+1)} = \pi^{(n)}P = \pi P = \pi$, and similarly $\pi^{(n+k)} = \pi$, $\forall k \geq 0$. That is, if a chain has a stationary distribution, it retains that distribution for all future time. Note that

$$\begin{aligned} \pi = \pi P &\iff \pi - \pi P = 0 \\ &\iff \pi(I - P) = 0 \end{aligned}$$

where I is the $r \times r$ identity matrix. Hence the stationary distribution of the chain may be found by solving

$$\pi(I - P) = 0. \quad (5.2)$$

Note that the trivial solution $\pi = 0$ is not of interest here, as it does not correspond to a probability distribution (its elements do not sum to 1). However, there are always infinitely many solutions to (5.2), so proper solutions can be found by finding a positive solution and then imposing the unit-sum constraint. In the case of a unique stationary distribution (just one eigenvalue of P equal to 1), then there will be a one-dimensional set of solutions to (5.2), and the unique stationary distribution will be the single solution with positive elements summing to 1.

5.2.4 Convergence

Convergence of Markov chains is a rather technical topic, which we do not have time to examine in detail here. This short section presents a very informal explanation of why Markov chains often do converge to their stationary distribution and how the rate of convergence can be understood.

By convergence to stationary distribution, we mean that irrespective of the starting distribution, $\pi^{(0)}$, the distribution at time n , $\pi^{(n)}$, will converge to the stationary distribution, π as n tends to infinity. If the limit of $\pi^{(n)}$ exists, it is referred to as the *equilibrium* distribution of the chain (sometimes referred to as the *limiting* distribution). Clearly the equilibrium distribution will be a stationary distribution, but a stationary distribution is not guaranteed to be an equilibrium distribution.[†] The relationship between stationary and equilibrium distributions is therefore rather subtle.

Let π be a (row) eigenvector of P with corresponding eigenvalue λ . Then

$$\pi P = \lambda \pi.$$

Also $\pi P^n = \lambda^n \pi$. It is easy to show that for stochastic P we must have $|\lambda| \leq 1$ (see Exercises). We also know that at least one eigenvector is equal to 1 (the corresponding eigenvector is a stationary distribution). Let

$$(\pi_1, \lambda_1), (\pi_2, \lambda_2), \dots, (\pi_r, \lambda_r)$$

be the full eigen-decomposition of P , with $|\lambda_i|$ in decreasing order, so that $\lambda_1 = 1$, and π_1 is a (rescaled) stationary distribution. To keep things simple, let us now make the assumption that the initial distribution $\pi^{(0)}$ can be written in the form

$$\pi^{(0)} = a_1 \pi_1 + a_2 \pi_2 + \dots + a_r \pi_r$$

for appropriate choice of a_i (this might not always be possible, but making the assumption keeps the mathematics simple). Then

$$\begin{aligned} \pi^{(n)} &= \pi^{(0)} P^n \\ &= (a_1 \pi_1 + a_2 \pi_2 + \dots + a_r \pi_r) P^n \\ &= a_1 \pi_1 P^n + a_2 \pi_2 P^n + \dots + a_r \pi_r P^n \\ &= a_1 \lambda_1^n \pi_1 + a_2 \lambda_2^n \pi_2 + \dots + a_r \lambda_r^n \pi_r \\ &\rightarrow a_1 \pi_1, \quad \text{as } n \rightarrow \infty, \end{aligned}$$

[†] This is clear if there is more than one stationary distribution, but in fact, even in the case of a unique stationary distribution, there might not exist an equilibrium distribution at all.

provided that $|\lambda_2| < 1$. The rate of convergence is governed by the second eigenvalue, λ_2 . Therefore, provided that $|\lambda_2| < 1$, the chain eventually converges to an equilibrium distribution, which corresponds to the unique stationary distribution, irrespective of the initial distribution. If there is more than one unit eigenvalue, then there is an infinite family of stationary distributions, and convergence to any particular distribution is not guaranteed. For more details on the theory of Markov chains and their convergence, it is probably a good idea to start with texts such as Ross (1996) and Cox & Miller (1977), and then consult the references therein as required. For the rest of this chapter we will assume that an equilibrium distribution exists and that it corresponds to a unique stationary distribution.

5.2.5 Reversible chains

If $\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(N)}$ is a Markov chain, then the reversed sequence of states, $\theta^{(N)}, \theta^{(N-1)}, \dots, \theta^{(0)}$ is also a Markov chain. To see this, consider the conditional distribution of the current state given the future:

$$\begin{aligned} & P\left(\theta^{(n)} = y \mid \theta^{(n+1)} = x_{n+1}, \dots, \theta^{(N)} = x_N\right) \\ &= \frac{P\left(\theta^{(n+1)} = x_{n+1}, \dots, \theta^{(N)} = x_N \mid \theta^{(n)} = y\right) P\left(\theta^{(n)} = y\right)}{P\left(\theta^{(n+1)} = x_{n+1}, \dots, \theta^{(N)} = x_N\right)} \\ &= \frac{P\left(\theta^{(n+1)} = x_{n+1} \mid \theta^{(n)} = y\right) \cdots \cdots P\left(\theta^{(N)} = x_N \mid \theta^{(N-1)} = x_{N-1}\right) P\left(\theta^{(n)} = y\right)}{P\left(\theta^{(n+1)} = x_{n+1}\right) P\left(\theta^{(n+2)} = x_{n+2} \mid \theta^{(n+1)} = x_{n+1}\right) \cdots \cdots P\left(\theta^{(N)} = x_N \mid \theta^{(N-1)} = x_{N-1}\right)} \\ &= \frac{P\left(\theta^{(n+1)} = x_{n+1} \mid \theta^{(n)} = y\right) P\left(\theta^{(n)} = y\right)}{P\left(\theta^{(n+1)} = x_{n+1}\right)} \\ &= P\left(\theta^{(n)} = y \mid \theta^{(n+1)} = x_{n+1}\right). \end{aligned}$$

This is exactly the condition required for the reversed sequence of states to be Markovian.

Now let $P_n^*(x, y)$ be the transition kernel for the reversed chain. Then

$$\begin{aligned} P_n^*(x, y) &= P\left(\theta^{(n)} = y \mid \theta^{(n+1)} = x\right) \\ &= \frac{P\left(\theta^{(n+1)} = x \mid \theta^{(n)} = y\right) P\left(\theta^{(n)} = y\right)}{P\left(\theta^{(n+1)} = x\right)} && \text{(by Theorem 3.2)} \\ &= \frac{P(y, x) \pi^{(n)}(y)}{\pi^{(n+1)}(x)}. \end{aligned}$$

Therefore, in general, the reversed chain is not homogeneous. However, if the chain has reached its stationary distribution, then

$$P^*(x, y) = \frac{P(y, x) \pi(y)}{\pi(x)},$$

and so the reversed chain is homogeneous, and has a transition matrix which may

be determined from the transition matrix for the forward chain (and its stationary distribution).

If

$$P^*(x, y) = P(x, y), \quad \forall x, y$$

then the chain is said to be (time) *reversible* (as then the sequence of states appear the same if time is reversed), and we have the *detailed balance equations*:

$$\pi(x) P(x, y) = \pi(y) P(y, x), \quad \forall x, y. \quad (5.3)$$

The detailed balance equations capture symmetry in the flow of probability between pairs of states. The left hand side is the probability that a particular transition will be a move from x to y , and the right hand side is the probability that a particular transition will be a move from y to x .

If we have a chain with transition kernel $P(x, y)$ and a distribution $\pi(\cdot)$ satisfying (5.3), then it follows that the chain is reversible with stationary distribution $\pi(\cdot)$. The chain also has other nice properties (such as positive recurrence) which make it comparatively well behaved.

Proposition 5.2 *Consider a Markov chain with transition matrix P , satisfying (5.3) for some probability vector π . Then the chain has π as a stationary distribution and is time reversible.*

Proof. Summing both sides of (5.3) with respect to x gives

$$\begin{aligned} \sum_x \pi(x) P(x, y) &= \sum_x \pi(y) P(y, x) \\ &= \pi(y) \sum_x P(y, x) \\ &= \pi(y), \quad \forall y. \end{aligned}$$

In other words, $\pi P = \pi$. Hence π is a stationary distribution. Assuming that this stationary distribution is the equilibrium distribution of the chain, we can immediately conclude that the chain is reversible by dividing (5.3) through by $\pi(x)$. \square

5.2.6 Stochastic simulation and analysis

We next turn our attention to the problem of stochastic simulation of Markov chains on a computer and analysis of the simulation results. The key requirement is the ability to simulate a new state randomly, with probabilities given by an arbitrary probability vector p . The solution is given in Section 4.5. For simulation of a Markov chain to take place, a transition matrix P and an initial distribution $\pi^{(0)}$ is required. Simulation begins with sampling an initial state $\theta^{(0)}$ from $\pi^{(0)}$ using a lookup method. Once the initial state has been obtained, a value for $\theta^{(1)}$ can be sampled using the set of probabilities from the $\theta^{(0)}$ th row of P . Indeed, at time t , the state $\theta^{(t)}$ can be sampled using the probabilities from the $\theta^{(t-1)}$ th row of P . A simple R function for simulating the path of a Markov chain is given in Figure 5.1, and an example R session illustrating its use is given in Figure 5.2. The example session shows how to plot and summarise the simulated sample path and how to calculate the proportion of time

```

rfmc <- function(n,P,pi0)
{
  v=vector("numeric",n)
  r=length(pi0)
  v[1]=sample(r,1,prob=pi0)
  for (i in 2:n) {
    v[i]=sample(r,1,prob=P[v[i-1],])
  }
  ts(v)
}

```

Figure 5.1 An R function to simulate a sample path of length n from a Markov chain with transition matrix P and initial distribution π_0

spent in each state. These latter proportions approximate the equilibrium distribution of the chain.

5.3 Markov chains with continuous state space

Mostly we will regard biochemical networks as having a discrete state, but sometimes it is helpful to regard the state of certain quantities as continuous. In this case, we have to understand how the concept of a discrete state Markov chain extends to the continuous state case. In fact, this extension is exactly analogous to the generalisation of discrete random quantities to that of continuous random quantities.

Here we are still working with discrete time, but we are allowing the state space S of the Markov chain to be continuous (e.g. $S \subseteq \mathbb{R}$).

Example — First-order auto-regressive model

Consider the AR(1) model, which arises in elementary time-series analysis. Here, AR stands for auto-regressive, and the (1) means that the order of the auto-regression is 1. The essential structure of the model for an AR(1) process $\{Z_t | t = 1, 2, \dots\}$ can be summarised as follows:

$$Z_t = \alpha Z_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim N(0, \sigma^2).$$

This model captures the idea that the value of the stochastic process at time t , Z_t , depends on the value of the stochastic process at time $t - 1$, and that the relationship is non-deterministic, with the non-determinism is captured in the “noise” term, ε_t . It is assumed that the noise process, $\{\varepsilon_t | t = 1, 2, \dots\}$ is independent (that is, the pair of random quantities ε_i and ε_j will be independent of one another $\forall i \neq j$). However, it is clear that the stochastic process of interest $\{Z_t | t = 1, 2, \dots\}$ is not independent due to the dependence introduced by auto-regressing each value on the value at the previous time.

It is clear that the conditional distribution of Z_t given $Z_{t-1} = z_{t-1}$ is just

$$Z_t | (Z_{t-1} = z_{t-1}) \sim N(\alpha z_{t-1}, \sigma^2),$$

For continuous state spaces we always have $P(x, \{y\}) = 0$, so in this case we define $P(x, y)$ by

$$\begin{aligned} P(x, y) &= P\left(\theta^{(n+1)} \leq y | \theta^{(n)} = x\right) \\ &= P\left(\theta^{(1)} \leq y | \theta^{(0)} = x\right), \quad \forall x, y \in S, \end{aligned}$$

the conditional cumulative distribution function (CDF). This is the distributional form of the transition kernel for continuous state space Markov chains, but we can also define the corresponding conditional density

$$p(x, y) = \frac{\partial}{\partial y} P(x, y), \quad x, y \in S.$$

We can use this to define the density form of the *transition kernel* of the chain. Note that $p(x, y)$ is just the conditional density for the next state (with variable y) given that the current state is x , so it could also be written $p(y|x)$. The density form of the kernel can be used more conveniently than the CDF form for vector Markov chains, where the state space is multidimensional (say $S \subseteq \mathbb{R}^n$).

Example

If we write our AR(1) in the form

$$\theta^{(t+1)} = \alpha\theta^{(t)} + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2)$$

then

$$(\theta^{(t+1)} | \theta^{(t)} = x) \sim N(\alpha x, \sigma^2),$$

and so the density form of the transition kernel is just the normal density

$$p(x, y) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{y - \alpha x}{\sigma}\right)^2\right\},$$

and this is to be interpreted as a density for y for a given fixed value of x .

5.3.2 Stationarity and reversibility

Let the state at time n , $\theta^{(n)}$ be represented by a probability density function, $\pi^{(n)}(x)$, $x \in S$. By the continuous version of the theorem of total probability (Proposition 3.11), we have

$$\pi^{(n+1)}(y) = \int_S p(x, y)\pi^{(n)}(x) dx. \quad (5.4)$$

We see from (5.4) that a stationary distribution must satisfy

$$\pi(y) = \int_S p(x, y)\pi(x) dx, \quad (5.5)$$

which is clearly just the continuous version of the discrete matrix equation $\pi = \pi P$.

Again, we can use Bayes Theorem to get the transition density for the reversed

chain

$$p_n^*(x, y) = \frac{p(y, x)\pi^{(n)}(y)}{\pi^{(n+1)}(x)},$$

which homogenises in the stationary limit to give

$$p^*(x, y) = \frac{p(y, x)\pi(y)}{\pi(x)}.$$

So, if the chain is (time) reversible, we have the continuous form of the detailed balance equations

$$\pi(x)p(x, y) = \pi(y)p(y, x), \quad \forall x, y \in S. \quad (5.6)$$

Proposition 5.3 *Any homogeneous Markov chain satisfying (5.6) is reversible with stationary distribution $\pi(\cdot)$.*

Proof. We can see that detailed balance implies stationarity of $\pi(\cdot)$ by integrating both sides of (5.6) with respect to x and comparing with (5.5). Once we know that $\pi(\cdot)$ is the stationary distribution, reversibility follows immediately. \square

This result is of fundamental importance in the study of Markov chains with continuous state space, as it allows us to verify the stationary distribution $\pi(\cdot)$ of a (reversible) Markov chain with transition kernel $p(\cdot, \cdot)$ without having to directly verify the integral equation (5.5). Note, however, that a given density $\pi(\cdot)$ will fail detailed balance (5.6) regardless of whether or not it is a stationary distribution of the chain if the chain itself is not time reversible.

Example

We know that linear combinations of normal random variables are normal (Section 3.10), so we expect the stationary distribution of our example AR(1) to be normal. The normal distribution is characterised by its mean and variance, so if we can deduce the mean and variance of the stationary distribution, we can check to see if it really is normal. At convergence, successive distributions are the same. In particular, the first and second moments at successive time points remain constant. First, $E(\theta^{(n+1)}) = E(\theta^{(n)})$, and so

$$\begin{aligned} E(\theta^{(n)}) &= E(\theta^{(n+1)}) \\ &= E(\alpha\theta^{(n)} + \epsilon_n) \\ &= \alpha E(\theta^{(n)}) \\ \Rightarrow E(\theta^{(n)}) &= 0. \end{aligned}$$

Similarly,

$$\begin{aligned}\text{Var}(\theta^{(n)}) &= \text{Var}(\theta^{(n+1)}) \\ &= \text{Var}(\alpha\theta^{(n)} + \epsilon_n) \\ &= \alpha^2 \text{Var}(\theta^{(n)}) + \sigma^2 \\ \Rightarrow \text{Var}(\theta^{(n)}) &= \frac{\sigma^2}{1 - \alpha^2}.\end{aligned}$$

So, we think the stationary distribution is normal with mean zero and variance $\sigma^2/(1 - \alpha^2)$. That is, we think the stationary density is

$$\pi(x) = \frac{1}{\sqrt{\frac{2\pi\sigma^2}{1-\alpha^2}}} \exp\left\{-\frac{1}{2} \frac{x^2}{\frac{\sigma^2}{1-\alpha^2}}\right\} = \sqrt{\frac{1-\alpha^2}{2\pi\sigma^2}} \exp\left\{-\frac{x^2(1-\alpha^2)}{2\sigma^2}\right\}.$$

Since we know the transition density for this chain, we can see if this density satisfies detailed balance:

$$\begin{aligned}\pi(x)p(x, y) &= \sqrt{\frac{1-\alpha^2}{2\pi\sigma^2}} \exp\left\{-\frac{x^2(1-\alpha^2)}{2\sigma^2}\right\} \times \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \left(\frac{y-\alpha x}{\sigma}\right)^2\right\} \\ &= \frac{\sqrt{1-\alpha^2}}{2\pi\sigma^2} \exp\left\{-\frac{1}{2\sigma^2}[x^2 - 2\alpha xy + y^2]\right\}\end{aligned}$$

after a little algebra. But this expression is exactly symmetric in x and y , and so

$$\pi(x)p(x, y) = \pi(y)p(y, x).$$

So we see that $\pi(\cdot)$ does satisfy detailed balance, and so the AR(1) is a *reversible* Markov chain and *does* have stationary distribution $\pi(\cdot)$.

5.3.3 Stochastic simulation and analysis

Simulation of Markov chains with a continuous state in discrete time is easy provided that methods are available for simulating from the initial distribution, $\pi^{(0)}(x)$, and from the conditional distribution represented by the transition kernel, $p(x, y)$.

1. First $\theta^{(0)}$ is sampled from $\pi^{(0)}(\cdot)$, using one of the techniques discussed in Chapter 4.
2. We can then simulate $\theta^{(1)}$ from $p(\theta^{(0)}, \cdot)$, as this is just a density.
3. In general, once we have simulated a realisation of $\theta^{(n)}$, we can simulate $\theta^{(n+1)}$ from $p(\theta^{(n)}, \cdot)$, using one of the standard techniques from Chapter 4.

Example

Let us start our AR(1) off at $\theta^{(0)} = 0$, so we do not need to simulate anything for the initial value. Next we want to simulate $\theta^{(1)}$ from $p(\theta^{(0)}, \cdot) = p(0, \cdot)$, that is, we simulate $\theta^{(1)}$ from $N(0, \sigma^2)$. Next we simulate $\theta^{(2)}$ from $p(\theta^{(1)}, \cdot)$, that is,

we simulate $\theta^{(2)}$ from $N(\alpha\theta^{(1)}, \sigma^2)$. In general, having simulated $\theta^{(n)}$, we simulate $\theta^{(n+1)}$ from $N(\alpha\theta^{(n)}, \sigma^2)$.

As n gets large, the distribution of $\theta^{(n)}$ tends to the distribution with density $\pi(\cdot)$, the equilibrium distribution of the chain. All values sampled after convergence has been reached are draws from $\pi(\cdot)$. There is a "burn-in" period before convergence is reached, so if interest is in $\pi(\cdot)$, these values should be discarded before analysis takes place.

If we are interested in an integral

$$\int_S g(x)\pi(x)dx = E_\pi(g(\Theta)),$$

then if $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n)}$ are draws from $\pi(\cdot)$, this integral may be approximated by

$$E_\pi(g(\Theta)) \simeq \frac{1}{n} \sum_{i=1}^n g(\theta^{(i)}).$$

However, draws from a Markov chain are not independent, so the variance of the sample mean cannot be computed in the usual way.

Suppose $\theta^{(i)} \sim \pi(\cdot)$, $i = 1, 2, \dots$. Then

$$E_\pi(\Theta) \simeq \frac{1}{n} \sum_{i=1}^n \theta^{(i)} = \bar{\theta}_n.$$

Let $\text{Var}(\Theta) = \text{Var}(\theta^{(i)}) = \nu^2$. Then if the θ_i were independent, we would have

$$\text{Var}(\bar{\theta}_n) = \frac{\nu^2}{n}.$$

However, if the $\theta^{(i)}$ are not independent (e.g. sampled from a non-trivial Markov chain), then

$$\text{Var}(\bar{\theta}_n) \neq \frac{\nu^2}{n}.$$

Example

This example relies on some results from multivariate probability theory not covered in the text. It can be skipped without significant loss.

$$\text{Var}(\theta^{(i)}) = \frac{\sigma^2}{1 - \alpha^2} = \nu^2$$

and

$$\gamma(k) = \text{Cov}(\theta^{(i)}, \theta^{(i+k)}) = \nu^2 \alpha^k,$$

so

$$\begin{aligned}\text{Var}(\bar{\theta}_n) &= \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n \theta^{(i)}\right) \\ &= \frac{\nu^2}{n^2} \left(n + \sum_{i=1}^{n-1} 2(n-i)\alpha^i\right),\end{aligned}$$

which sums to give

$$\begin{aligned}\text{Var}(\bar{\theta}_n) &= \frac{1}{n^2} \frac{\sigma^2}{1-\alpha^2} \frac{n + 2\alpha^{n+1} - 2\alpha - n\alpha^2}{(1-\alpha)^2} \\ &= \frac{1}{n} \frac{\sigma^2}{1-\alpha^2} \left[\frac{1+\alpha}{1-\alpha} - \frac{2\alpha(1-\alpha^n)}{n(1-\alpha)^2}\right].\end{aligned}$$

To a first approximation, we have

$$\text{Var}(\bar{\theta}_n) \simeq \frac{1}{n} \frac{\sigma^2}{1-\alpha^2} \frac{1+\alpha}{1-\alpha},$$

and so the “correction factor” for the naive calculation based on an assumption of independence is $(1+\alpha)/(1-\alpha)$. For α close to one, this can be very large, for example, for $\alpha = 0.95$, $(1+\alpha)/(1-\alpha) = 39$, and so the variance of the sample mean is actually around 40 times bigger than calculations based on assumptions of independence would suggest. Similarly, confidence intervals should be around six times wider than calculations based on independence would suggest.

We can actually use this analysis in order to analyse other Markov chains. If the Markov chain is reasonably well approximated by an AR(1) (and many are), then we can estimate the variance of our sample estimates by the AR(1) variance estimate. For an AR(1), α is just the lag 1 auto-correlation of the chain ($\text{Corr}(\theta^{(i)}, \theta^{(i+1)}) = \alpha$), and so we can estimate the α of any simulated chain by the sample auto-correlation at lag 1. We can then use this to compute or correct sample variance estimates based on sample means of chain values.

5.4 Markov chains in continuous time

We have now looked at Markov chains in discrete time with both discrete and continuous state spaces. However, biochemical processes evolve continuously in time, and so we now turn our attention to the continuous time case. We begin by studying chains with a finite number of states, but relax this assumption in due course.

Before we begin we should try to be explicit about what exactly a Markov process is in the continuous time case. Intuitively, it is a straightforward extension of the discrete time definition. In continuous time, we could write this as

$$\begin{aligned}P(X(t+dt) = x | \{X(t) = x(t) | t \in [0, t]\}) \\ = P(X(t+dt) = x | X(t) = x(t)), \quad \forall t \in [0, \infty), x \in S.\end{aligned}$$

Again, this expresses the idea that the future behaviour of the process depends on the past behaviour of the process only via the current state.

5.4.1 Finite state-space

Consider first a process which can take on one of r states, which we label $S = \{1, 2, \dots, r\}$. If at time t the process is in state $x \in S$, its future behaviour can be characterised by the transition kernel

$$p(x, t, x', t') \equiv P(X(t+t') = x' | X(t) = x).$$

If this function does not depend explicitly on t , the process is said to be *homogeneous*, and the kernel can be written $p(x, x', t')$. For each value of t' , this kernel can be expressed as an $r \times r$ transition matrix, $P(t')$. It is clear that $P(0) = I$, the $r \times r$ identity matrix, as no transitions will take place in a time interval of length zero. Also note that since $P(\cdot)$ is a transition matrix for each value of t , we can multiply these matrices together to give combined transition matrices in the usual way. In particular, we have $P(t+t') = P(t)P(t') = P(t')P(t)$, just as in the discrete time case.[§] Now define the *transition rate matrix*, Q , to be the derivative of $P(t')$ at $t' = 0$. Then

$$\begin{aligned} Q &= \left. \frac{d}{dt'} P(t') \right|_{t'=0} \\ &= \lim_{\delta t \rightarrow 0} \frac{P(\delta t) - P(0)}{\delta t} \\ &= \lim_{\delta t \rightarrow 0} \frac{P(\delta t) - I}{\delta t}. \end{aligned}$$

The elements of the Q matrix give the “hazards” of moving to different states. Rearranging gives the infinitesimal transition matrix

$$P(dt) = I + Q dt.$$

Note that for $P(dt)$ to be a stochastic matrix (with non-negative elements and rows summing to 1), the above implies several constraints which must be satisfied by the rate matrix Q . Since the off-diagonal elements of I are zero, the off-diagonal elements of $P(dt)$ and $Q dt$ must be the same, and so the off-diagonal elements of Q must be non-negative. Also, since the diagonal elements of $P(dt)$ are bounded above by one, the diagonal elements of Q must be non-positive. Finally, since the rows of $P(dt)$ and I both sum to 1, the rows of Q must each sum to zero. These properties must be satisfied by *any* rate matrix Q .

The above rearrangement gives us a way of computing the stationary distribution of the Markov chain, as a probability row vector π will be stationary only if

$$\begin{aligned} \pi P(dt) &= \pi \\ \Rightarrow \pi(I + Q dt) &= \pi \\ \Rightarrow \pi Q &= 0. \end{aligned}$$

[§] When these equations are written out in full, as $p(i, j, t+t') = \sum_{k=1}^r p(i, k, t)p(k, j, t')$, $i, j = 1, \dots, r$, they are known as the *Chapman-Kolmogorov equations*.

Solving this last equation (subject to the constraint that the elements of π sum to 1), will give a stationary distribution for the system.

If $P(t)$ is required for finite t , it may be computed by solving a matrix differential equation. This can be derived by considering the derivative of $P(t)$ for arbitrary times t .

$$\begin{aligned} \frac{d}{dt}P(t) &= \frac{P(t + dt) - P(t)}{dt} \\ &= \frac{P(dt)P(t) - P(t)}{dt} \\ &= \frac{P(dt) - I}{dt}P(t) \\ &= Q P(t). \end{aligned}$$

Therefore, $P(t)$ is a solution to the matrix differential equation

$$\frac{d}{dt}P(t) = QP(t),$$

subject to the initial condition $P(0) = I$. This differential equation has solution

$$P(t) = \exp\{Q t\},$$

where $\exp\{\cdot\}$ denotes the matrix exponential function (Golub & Van Loan 1996). Working with the matrix exponential function is straightforward, but beyond the scope of this book.[¶] Note that if we prefer, we do not have to work with the differential equation in matrix form. If we write it out in component form we have

$$\frac{d}{dt}p(i, j, t) = \sum_{k=1}^r q_{ik}p(k, j, t), \quad i, j = 1, 2, \dots, r. \tag{5.7}$$

These are known as *Kolmogorov's backward equations* (Allen 2003). It is also worth noting that had we expanded $P(t + dt)$ as $P(t)P(dt)$ rather than $P(dt)P(t)$, we would have ended up with the matrix differential equation

$$\frac{d}{dt}P(t) = P(t)Q,$$

which we can write out in component form as

$$\frac{d}{dt}p(i, j, t) = \sum_{k=1}^r q_{kj}p(i, k, t), \quad i, j = 1, 2, \dots, r. \tag{5.8}$$

These are known as *Kolmogorov's forward equations*.

[¶] Note that the exponential of a matrix is not the matrix obtained by applying the exponential function to each element of the matrix. In fact, it is usually defined by its series expansion, $\exp\{A\} = \sum_{k=0}^{\infty} A^k/k! = I + A + A^2/2 + \dots$, but this expansion does not lead to an efficient way of computing the function. Note also that many scientific libraries for numerical linear algebra provide a function for computing the matrix exponential.

```

@model:2.1.1=Activation
@units
  substance=item
@compartments
  Cell
@species
  Cell:A=0 s
  Cell:I=1 s
@reactions
@r=Activation
  I -> A
  alpha : alpha=0.5
@r=Inactivation
  A -> I
  beta : beta=1

```

Figure 5.3 *SBML-shorthand for the simple gene activation process with $\alpha = 0.5$ and $\beta = 1$*

Example

Consider a very simple model for the activation of a single prokaryotic gene. In this model, the gene will be activated unless a repressor protein is bound to its regulatory region. We will consider just two states in our system: state 0 (inactive), and state 1 (active). In the inactive state (0), we will assume a constant hazard of $\alpha > 0$ for activation. In the active state, we will assume a constant hazard of $\beta > 0$ for inactivation. Given that the rows of Q must sum to zero, it is now completely specified as

$$Q = \begin{pmatrix} -\alpha & \alpha \\ \beta & -\beta \end{pmatrix}.$$

Solving $\pi Q = 0$ gives the stationary distribution

$$\pi = \left(\frac{\beta}{\alpha + \beta}, \frac{\alpha}{\alpha + \beta} \right).$$

We can also compute the infinitesimal transition matrix

$$P(dt) = I + Q dt = \begin{pmatrix} 1 - \alpha dt & \alpha dt \\ \beta dt & 1 - \beta dt \end{pmatrix}.$$

It is straightforward to encode this model in SBML. The SBML-shorthand for it is given in Figure 5.3 (and the full SBML can be downloaded from this book's website). A simulated realisation of this process is shown in Figure 5.4.

5.4.2 Stochastic simulation

There are three straightforward approaches one can take to simulating this process on a computer. The first is based on a fine time-discretisation of the process, similar in spirit to the first-order Euler method for integrating ordinary differential equations.

Given the definition of the infinitesimal transition matrix

$$P(dt) = I + Q dt,$$

for small time steps Δt we will have

$$P(\Delta t) \simeq I + Q \Delta t.$$

$P(\Delta t)$ can then be regarded as the transition matrix of a discrete time Markov chain, and a simulated sequence of states at times $0, \Delta t, 2 \Delta t, 3 \Delta t, \dots$ may be generated in the usual way.

The above method can be easily improved by replacing the above approximation for $P(\Delta t)$ by its *exact* value

$$P(\Delta t) = \exp\{Q \Delta t\},$$

provided a method for computing the matrix exponential is available. Then it does not matter how small Δt is chosen to be, provided it is small enough to clearly show the behaviour of the process and not so large that interesting transitions are "missed."

A third approach to simulation may be taken by simulating each transition event and its corresponding time sequentially, rather than simply looking at the processes only at times on a given lattice. Like the previous method, this gives an exact realisation of the process and offers the additional advantage that recording every reaction event ensures none will be "missed." Such an approach is known as *discrete event simulation*. If the process is currently in state x , then the x th row of Q gives the hazards for the transition to other states. As the row sums to zero, $-q_{xx}$ gives the combined hazard for moving away from the current state — a discrete transition event (note that q_{xx} is non-positive). So, the time to a transition event is exponential with rate $-q_{xx}$. When that transition event occurs, the new state will be random with probabilities proportional to the x th row of Q (with q_{xx} omitted). The above intuitive explanation can be formalised as follows.

To understand how to simulate the process we must consider being in state i at time t , and think about the probability that the next event will be in the time interval $(t + t', t + t' + dt]$, and will consist of a move to state j . Let this probability divided by dt be denoted by $f(t', j|t, i)$, so that the probability is $f(t', j|t, i)dt$. It is clear that as the Markov process is homogeneous, there will be no explicit dependence on t in this probability, but we will include it to be clear about exactly what we mean. Then

$$f(t', j|t, i)dt = P(\text{Next event in } (t + t', t + t' + dt]|t, i) \\ \times P(j|\text{Next event in } (t + t', t + t' + dt], t, i).$$

Thinking about the first term, we know that the hazards for the individual transitions are given by the off-diagonal elements of the i th row of Q . The combined hazard is the sum of these off-diagonal elements, which is $-q_{ii}$ (as the row sums to zero). Combined hazards can always be computed as the sum of hazards in this way because the probability that two events occur in the interval $(t, t + dt]$ is of order dt^2 and can therefore be neglected. Now we know from our consideration of the exponential distribution as the time to an event of constant hazard that the time to the next event

is $Exp(-q_{ii})$, and so the first term must be $-q_{ii}e^{q_{ii}t'} dt$. The second term is

$$\begin{aligned} P(X(t+t'+dt) = j | [X(t+t') = i] \cap [X(t+t'+dt) \neq i]) \\ = \frac{P(X(t+t'+dt) = j | X(t+t') = i)}{P(X(t+t'+dt) \neq i | X(t+t') = i)} = \frac{q_{ij} dt}{\sum_{k \neq i} q_{ik} dt} = \frac{q_{ij}}{-q_{ii}}. \end{aligned}$$

Taking the two terms together we have

$$f(t', j | t, i) = -q_{ii}e^{q_{ii}t'} \times \frac{q_{ij}}{-q_{ii}}.$$

The fact that this function factorises into the form of a probability density for the time to the next event and a probability mass function for the type of that event means that we can simulate the next event with the generation of two random variables. Note also that there is no j dependence in the PDF for t' and no t' dependence in the PMF for j , so the two random variables are independent of one another and hence can be simulated independently.

It is the consideration of $f(t', j | t, i)$ that leads to the standard discrete event simulation algorithm which could be stated as follows:

1. Initialise the process at $t = 0$ with initial state i ;
2. Call the current state i . Simulate the time to the next event, t' , as an $Exp(-q_{ii})$ random quantity;
3. Put $t := t + t'$;
4. Simulate new state j as a discrete random quantity with PMF $-q_{ik}/q_{ii}$, $k \neq i$;
5. Output the time t and state j ;
6. If $t < T_{max}$, return to step 2.

This particular discrete event simulation technique is known as the *direct method*. A simple R function to implement this algorithm is given in Figure 5.5. The function returns a step-function object, which is easy to plot. Using this function, a plot similar to Figure 5.4 can be obtained with the following command:

```
plot(rcfmc(20, matrix(c(-0.5, 0.5, 1, -1), ncol=2,
                      byrow=TRUE), c(1, 0)))
```

All of these simulation methods give a single realisation of the Markov process. Now obviously, just as one would not study a normal distribution by looking at a single simulated value, the same is true with Markov processes. Many realisations must be simulated in order to get a feel for the *distribution* of values at different times. In the case of a finite number of states, this distribution is relatively straightforward to compute directly without any simulation at all, but for the stochastic kinetic models we will consider later, simulation is likely to be the only tool we have available to us for gaining insight into the behaviour of the process.

5.4.3 Countable state-space

Before moving on to thinking about continuous state spaces, it is worth spending a little time looking at the case of a countably infinite state space. Rather than at-

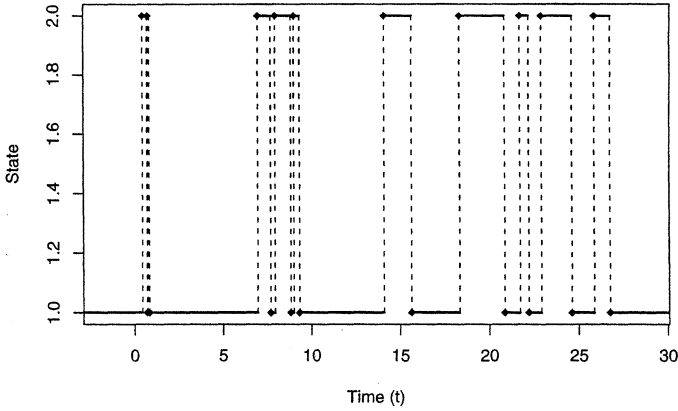


Figure 5.4 A simulated realisation of the simple gene activation process with $\alpha = 0.5$ and $\beta = 1$

tempting to present the theory in generality, we will concentrate on a simple example, which illustrates many of the interesting features. The model is known as the *immigration-death process*. In this model, individuals arrive into the population with constant hazard λ , and each individual dies independently with constant hazard μ . Consequently, the population of individuals increases by one when an immigration event occurs and decreases by one when a death event occurs. There is no reproduction in this model. Figure 5.6 gives the SBML-shorthand corresponding to this model. The key transition equations are:

$$\begin{aligned}
 P(X(t + dt) = x + 1 | X(t) = x) &= \lambda dt \\
 P(X(t + dt) = x - 1 | X(t) = x) &= x\mu dt \\
 P(X(t + dt) = x | X(t) = x) &= 1 - (\lambda + x\mu)dt \\
 P(X(t + dt) = y | X(t) = x) &= 0, \forall y \notin \{x - 1, x, x + 1\}.
 \end{aligned}$$

These equations clearly define a homogeneous Markov process, but with infinite state space $S = 0, 1, 2, \dots$. We therefore cannot easily write down a set of matrix equations for the process, as the matrices are infinite dimensional, but this does not prevent us from working with the process or from simulating it on a computer.

First let's think about understanding this process theoretically. Although the Q

```
rcfmc <- function(n,Q,pi0)
{
  xvec=vector("numeric",n+1)
  tvec=vector("numeric",n)
  r=length(pi0)
  x=sample(r,1,prob=pi0)
  t=0
  xvec[1]=x
  for (i in 1:n) {
    t=t+rexp(1,-Q[x,x])
    weights=Q[x,]
    weights[x]=0
    x=sample(r,1,prob=weights)
    xvec[i+1]=x
    tvec[i]=t
  }
  stepfun(tvec,xvec)
}
```

Figure 5.5 An R function to simulate a sample path with n events from a continuous time Markov chain with transition rate matrix Q and initial distribution π_0

```
@model:2.1.1=ImmigrationDeath
@units
  substance=item
@compartments
  Cell
@species
  Cell:X=0 s
@reactions
@r=Immigration
  -> X
  lambda : lambda=1
@r=Death
  X ->
  mu*X : mu=0.1
```

Figure 5.6 SBML-shorthand for the immigration-death process with $\lambda = 1$ and $\mu = 0.1$

matrix is infinite in extent, we can write its general form as follows:

$$Q = \begin{pmatrix} -\lambda & \lambda & 0 & 0 & 0 & \dots \\ \mu & -\lambda - \mu & \lambda & 0 & 0 & \dots \\ 0 & 2\mu & -\lambda - 2\mu & \lambda & 0 & \dots \\ 0 & 0 & 3\mu & -\lambda - 3\mu & \lambda & \ddots \\ \vdots & & \ddots & \ddots & \ddots & \ddots \end{pmatrix}.$$

Then for an infinite dimensional $\pi = (\pi_0, \pi_1, \pi_2, \dots)$ we can solve $\pi Q = 0$ to get the stationary distribution one equation at a time, expressing each π_k in terms of π_0 to find the general form

$$\pi_k = \frac{\lambda^k}{k! \mu^k} \pi_0, \quad k = 1, 2, \dots$$

But these are terms in the expansion of $\pi_0 e^{\lambda/\mu}$, and so imposing the unit-sum constraint we get $\pi_0 = e^{-\lambda/\mu}$ giving the general solution

$$\pi_k = \frac{(\lambda/\mu)^k e^{-\lambda/\mu}}{k!}, \quad k = 0, 1, 2, \dots$$

This is easily recognised as the PMF of a Poisson random quantity with mean λ/μ (Section 3.6). Hence, the stationary distribution of this process is Poisson with mean λ/μ .

We can also simulate realisations of this process on a computer. Here it is easiest to use the technique of discrete event simulation. If the current state of the process is x , the combined hazard for moving away from the current state is $\lambda + x\mu$, and so the time to the next event is an exponentially distributed random quantity with rate $\lambda + x\mu$. When that event occurs, the process will move up or down with probabilities proportional to their respective hazards, λ and $x\mu$. That is, the state will increase by 1 with probability $\lambda/(\lambda + x\mu)$ and decrease by 1 with probability $x\mu/(\lambda + x\mu)$. This sequence can be easily simulated on a computer to give a set of states and event times which can be plotted, summarised, etc. A simulated realisation of this immigration-death process is shown in Figure 5.7. An R function to simulate the process is given in Figure 5.8.

5.4.4 Inhomogeneous Poisson process

Our treatment of the Poisson process has so far been fairly low level and intuitive. The relationship between the (homogeneous) Poisson process, the Poisson distribution, and the exponential distribution was made explicit in Proposition 3.17. The Poisson process is described as *homogeneous* because the event hazard λ is constant throughout time. In this section we will see that the Poisson process can be regarded as a Markov process and understand how it may be generalised to the *inhomogeneous* case. Understanding the inhomogeneous Poisson process will be necessary for some of the fast, accurate hybrid stochastic simulation algorithms considered in Chapter 8.

Recall that for a (homogeneous) Poisson process with rate λ , we previously defined N_t to be the number of events of the Poisson process in the interval $(0, t]$, noting that we therefore have $N_t \sim Po(\lambda t)$. The process N_t is the *counting process* associated with the *point process* that is the Poisson process itself (represented by a collection of random event times). It turns out that the counting process N_t is a

© Cambridge University Press 2008

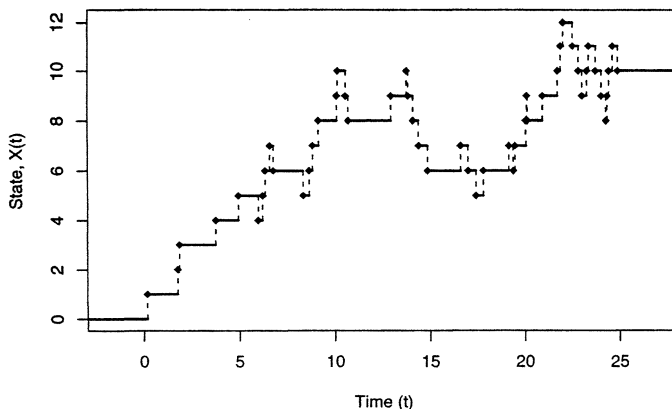


Figure 5.7 A single realisation of the immigration-death process with parameters $\lambda = 1$ and $\mu = 0.1$, initialised at $X(0) = 0$. Note that the stationary distribution of this process is Poisson with mean 10.

Markov process, governed by the homogeneous Markovian transition equations

$$\begin{aligned} P(N_{t+dt} = x + 1 | N_t = x) &= \lambda dt \\ P(N_{t+dt} = x | N_t = x) &= 1 - \lambda dt \\ P(N_{t+dt} = y | N_t = x) &= 0, \quad \forall y \notin \{x, x + 1\}. \end{aligned}$$

The Poisson process is special because its transition rates do not depend on the current state. It is this lack of dependence on the current state that leads to the analytical tractability of the Poisson process. Associating the Poisson (point) process with the Markov (counting) process N_t gives a new and powerful way of thinking about the relation between the two. Indeed, many standard texts choose to *define* the Poisson process in terms of the associated counting process, as the latter is much more convenient mathematically (it is a state-independent time-homogeneous Markov process). However, in my opinion, the distinction between a point process and its associated counting process is conceptually quite important.

All of the Markov processes that we have considered so far have been homogeneous in the sense that the transition equations do not depend explicitly on the current time, t . Consider now a generalisation of the Poisson process where the event hazard is not a constant λ , but a function, $\lambda(t)$. That is, the probability of an event in the interval $(t, t + dt]$ is $\lambda(t)dt$, so the inhomogeneous Markovian transition equations

```

imdeath <- function(n=20,x0=0,lambda=1,mu=0.1)
{
  xvec=vector("numeric",n+1)
  tvec=vector("numeric",n)
  t=0
  x=x0
  xvec[1]<-x
  for (i in 1:n) {
    t=t+rexp(1,lambda+x*mu)
    if ( runif(1,0,1) < lambda/(lambda+x*mu) )
      x <- x+1
    else
      x <- x-1
    xvec[i+1]<-x
    tvec[i]<-t
  }
  stepfun(tvec,xvec)
}

```

Figure 5.8 *R* function for discrete-event simulation of the immigration-death process

for the associated counting process N_t are

$$\begin{aligned}
 P(N_{t+dt} = x + 1 | N_t = x) &= \lambda(t)dt \\
 P(N_{t+dt} = x | N_t = x) &= 1 - \lambda(t)dt \\
 P(N_{t+dt} = y | N_t = x) &= 0, \quad \forall y \notin \{x, x + 1\}.
 \end{aligned}$$

A formal analysis of this process is fairly straightforward, and the reader is referred to a standard text such as Ross (1996) for the technical details. Intuitively, as the hazard is approximately constant in a sufficiently small interval, the number of events in that interval will be approximately Poisson. In the limit, the number of events in the interval $(t, t + dt]$ will be $Po(\lambda(t)dt)$, independent of all other intervals. The number of events in the interval $(0, t]$, N_t will then be the sum (integral) over all such intervals. As the sum of independent Poissons is Poisson, we get

$$N_t \sim Po\left(\int_0^t \lambda(s)ds\right).$$

It is helpful to define the *cumulative hazard*

$$\Lambda(t) = \int_0^t \lambda(s)ds$$

which then gives $N_t \sim Po(\Lambda(t))$. Similarly, the number of events in the interval $(s, t]$, $0 < s < t$ is $Po(\Lambda(t) - \Lambda(s))$.

Proposition 5.4 *For the inhomogeneous Poisson process with rate function $\lambda(t)$, the time, T , to the first event has distribution function*

$$F(t) = 1 - \exp\{-\Lambda(t)\},$$

and hence has density function

$$f(t) = \lambda(t) \exp\{-\Lambda(t)\}, \quad t > 0,$$

where $\Lambda(t)$ is the cumulative hazard defined above.

Proof.

$$\begin{aligned} F(t) &= P(T \leq t) \\ &= 1 - P(T > t) \\ &= 1 - P(N_t = 0) \\ &= 1 - \frac{\Lambda(t)^0 \exp\{-\Lambda(t)\}}{0!} \\ &= 1 - \exp\{-\Lambda(t)\}. \end{aligned}$$

□

In stochastic simulation one will often want to simulate the time to the first (or next) event of such a process. Using the inverse distribution method (Proposition 4.1) we can simulate $u \sim U(0, 1)$ and then solve $F(t) = u$ for t . Rearranging gives $\Lambda(t) = -\log(1 - u)$. However, as observed previously, $1 - u$ has the same distribution as u , so we just want to solve

$$\Lambda(t) = -\log u \tag{5.9}$$

for t . For simple hazards it will often be possible to solve this analytically, but in general a numerical procedure will be required.

Note that by construction the function $\Lambda(t)$ is monotonically increasing. So, for a given $u \in (0, 1)$, (5.9) will have at most one solution. It is also clear that unless the function $\Lambda(t)$ has the property that it tends to infinity as t tends to infinity, there may not be a solution to (5.9) at all. That is, the first event may not happen at all ever. Consequently, when dealing with the inhomogeneous Poisson process, attention is usually restricted to the case where this is true. In practice this means that the event hazard $\lambda(t)$ is not allowed to decay faster than $1/t$. Assuming this to be the case, there will always be exactly one solution to (5.9). This turns out to be useful if rather than knowing the exact event time, one simply needs to know whether or not the event has occurred before a given time t . For then it is clear that if $\Lambda(t) + \log u$ is negative, the event has not yet occurred, and hence the event time is greater than t , and if it is positive, the event time is less than t . The event time itself is clearly the unique root of the expression $\Lambda(t) + \log u$ (regarded as a function of t). Then if the event time really is required, it can either be found analytically, or failing this, an interval bisection method can be used to find it extremely quickly. Although this discussion might currently seem a little theoretical, it turns out to be a very important practical part of several hybrid stochastic simulation algorithms for biochemical network simulation, so is important to understand.^{||}

^{||} In the context of biochemical network simulation, the cumulative hazard, $\Lambda(t)$, is often known analytically, which makes the procedure just discussed fairly efficient. However, if only the hazard function, $\lambda(t)$ is known, and the cumulative hazard cannot be evaluated without numerical integration, then the procedure is not at all efficient. It turns out that as long as one can establish an upper bound for $\lambda(t)$

Example

Consider the inhomogeneous Poisson process with rate function $\lambda(t) = \lambda t$ for some constant $\lambda > 0$. This process has an event hazard that linearly increases with time. The cumulative hazard is clearly given by $\Lambda(t) = \lambda t^2/2$. From this we can immediately deduce that $N_t \sim Po(\lambda t^2/2)$ and that the number of events in the interval $(s, t]$ is $Po(\lambda(t^2 - s^2)/2)$. The time to the first event has PDF

$$f(t) = \lambda t \exp \left\{ -\frac{\lambda t^2}{2} \right\}, \quad t > 0$$

and this time can be simulated by sampling $u \sim U(0, 1)$ and solving $\lambda t^2/2 = -\log u$ for t to get

$$t = \sqrt{-\frac{2 \log u}{\lambda}}.$$

5.5 Diffusion processes

Markov processes with continuous states that evolve continuously in time are often termed *diffusion processes*. While these processes are extremely important, a formal discussion of the theory of such processes is beyond the scope of a book such as this. Nevertheless, it is useful to provide a brief non-technical introduction at this point, as these processes provide an excellent approximation to biochemical network models in certain situations.

A d -dimensional Itô diffusion process Y is governed by a stochastic differential equation (SDE) of the form

$$dY_t = \mu(Y_t)dt + \Lambda(Y_t)dW_t, \tag{5.10}$$

where $\mu : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a d -dimensional drift vector and $\Lambda : \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d$ is a $(d \times d)$ -dimensional diffusion matrix. The SDE can be thought of as a recipe for constructing a realisation of Y from a realisation of a d -dimensional Brownian motion (or Wiener process), W . A d -dimensional Brownian motion has d independent components, each of which is a univariate Brownian motion, B . A univariate Brownian motion B is a process defined for $t \geq 0$ in the following way.

1. $B_0 = 0$,
2. $B_t - B_s \sim N(0, t - s)$, $\forall t > s$,
3. The increment $B_t - B_s$ is independent of the increment $B_{t'} - B_{s'}$, $\forall t > s \geq t' > s'$.

It is clear from property 2 that $B_t \sim N(0, t)$ (and so $E(B_t) = 0$ and $\text{Var}(B_t) = t$). It is also clear that if for some small time increment Δt we define the process increment $\Delta B_t = B_{t+\Delta t} - B_t$, we then have $\Delta B_t \sim N(0, \Delta t)$, $\forall t$, and since we know that the

over the time interval of interest (usually trivial), it is possible to use *exact sampling* techniques to definitively decide if an event has occurred and if so at what time, based only on the ability to evaluate the hazard at a small number of time points. Such techniques are used frequently in applied probability and computational statistics, but do not yet seem to be widely known in the stochastic kinetics literature; see Wilkinson (2006) for further details.

```

rdiff <- function(afun, bfun, x0 = 0, t = 50, dt = 0.01, ...)
{
  n <- t/dt
  xvec <- vector("numeric", n)
  x <- x0
  sdt <- sqrt(dt)
  for (i in 1:n) {
    t <- i*dt
    x <- x + afun(x, ...) * dt +
              bfun(x, ...) * rnorm(1, 0, sdt)
    xvec[i] <- x
  }
  ts(xvec, deltat = dt)
}

```

Figure 5.9 R function for simulation of a diffusion process using the Euler method

increments are independent of one another, this provides a mechanism for simulating the process on a regular grid of time points.

If we define the increment in the diffusion process Y (and the multivariate Brownian motion W) similarly, then we can interpret the SDE (5.10) as the limit of the difference equation

$$\Delta Y_t = \mu(Y_t)\Delta t + \Lambda(Y_t)\Delta W_t, \quad (5.11)$$

as Δt gets infinitesimally small. For finite Δt , (5.11) is known as the *Euler approximation* (or, more correctly, as the *Euler-Maruyama approximation*) of the SDE, and it provides a simple mechanism for approximate simulation of the process Y on a regular grid of time points.**

In the case $d = 1$ we have a univariate diffusion process, and it is clear that then the increments of the process are approximately distributed as

$$\Delta Y_t \sim N(\mu(Y_t)\Delta t, \Lambda(Y_t)^2 \Delta t).$$

An R function to simulate a univariate diffusion using an Euler approximation is given in Figure 5.9. Note that more efficient simulation strategies are possible; see Kloeden & Platen (1992) for further details.

We can approximate a discrete Markov process using a diffusion by choosing the functions $\mu(\cdot)$ and $\Lambda(\cdot)$ so that the mean and variance of the increments match. This is best illustrated by example.

Example — diffusion approximation of the immigration-death process

Suppose we have an immigration-death process with immigration rate λ and death rate μ , and that at time t the current state of the system is x . Then at time $t + dt$, the

** “Euler” is pronounced, “oil-er.”

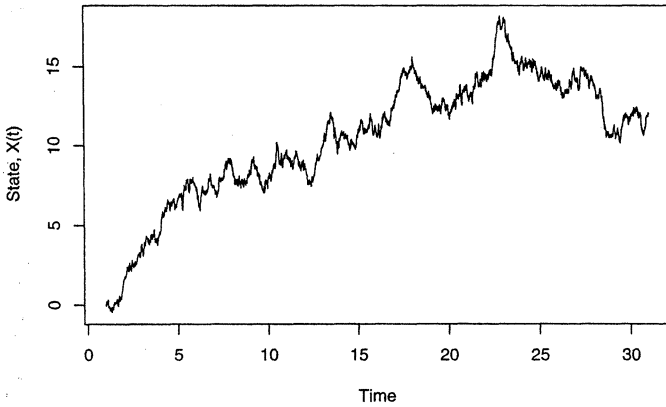


Figure 5.10 A single realisation of the diffusion approximation to the immigration-death process with parameters $\lambda = 1$ and $\mu = 0.1$, initialised at $X(0) = 0$

state of the system is a discrete random quantity with PMF,

$$\begin{aligned} P(X_{t+dt} = x - 1) &= x\mu dt, \\ P(X_{t+dt} = x) &= 1 - (\lambda + x\mu)dt, \\ P(X_{t+dt} = x + 1) &= \lambda dt. \end{aligned}$$

So the increment of the process has PMF

$$P(dX_t = -1) = x\mu dt, \quad P(dX_t = 0) = 1 - (\lambda + x\mu)dt, \quad P(dX_t = 1) = \lambda dt.$$

From this PMF we can calculate the expectation and variance as

$$E(dX_t) = (\lambda - \mu x)dt, \quad \text{Var}(dX_t) = (\lambda + \mu x)dt.$$

We therefore set $\mu(x) = \lambda - \mu x$ and $\Lambda(x)^2 = \lambda + \mu x$ to get the diffusion approximation

$$dX_t = (\lambda - \mu x)dt + \sqrt{\lambda + \mu x} dW_t.$$

We can use our code for simulating diffusion processes to get sample paths like that shown in Figure 5.10 using the R code shown in Figure 5.11.

5.6 Exercises

1. (a) Show that the product of two stochastic matrices is stochastic.
- (b) Show that for stochastic P , and any row vector π , we have $\|\pi P\|_1 \leq \|\pi\|_1$, where $\|v\|_1 = \sum_i |v_i|$. Deduce that all eigenvalues, λ , of P must satisfy


```

afun <- function(x, lambda, mu)
{
  lambda-mu*x
}
bfun <- function(x, lambda, mu)
{
  sqrt(lambda+mu*x)
}
plot(rdifff(afun,bfun,lambda=1,mu=0.1,t=30))

```

Figure 5.11 R code for simulating the diffusion approximation to the immigration-death process

$|\lambda| \leq 1$. As an aside, note that as a consequence of (a), it is clear that for all probability row vectors, π , we have $\|\pi P\|_1 = \|\pi\|_1$.

- (c) Show that a stochastic matrix always has a row eigenvector with eigenvalue 1. Show that this row eigenvector can be chosen to correspond to a stationary distribution of the induced Markov chain. *Hint: It is an easy consequence of the Brouwer Fixed Point Theorem, which says, inter alia, that any continuous map from a compact convex set to itself must have a fixed point.*
2. For the AR(1) process

$$Z_t = \alpha Z_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim \text{Exp}(\lambda), \quad \alpha \in (0, 1), \lambda > 0,$$

- (a) what is the transition kernel, $p(x, y)$, of the chain?
- (b) Hence, deduce an integral equation satisfied by the stationary distribution, $\pi(x)$.
- (c) If Z has density $\pi(x)$, calculate the expected value of Z and show that the variance of Z is given by

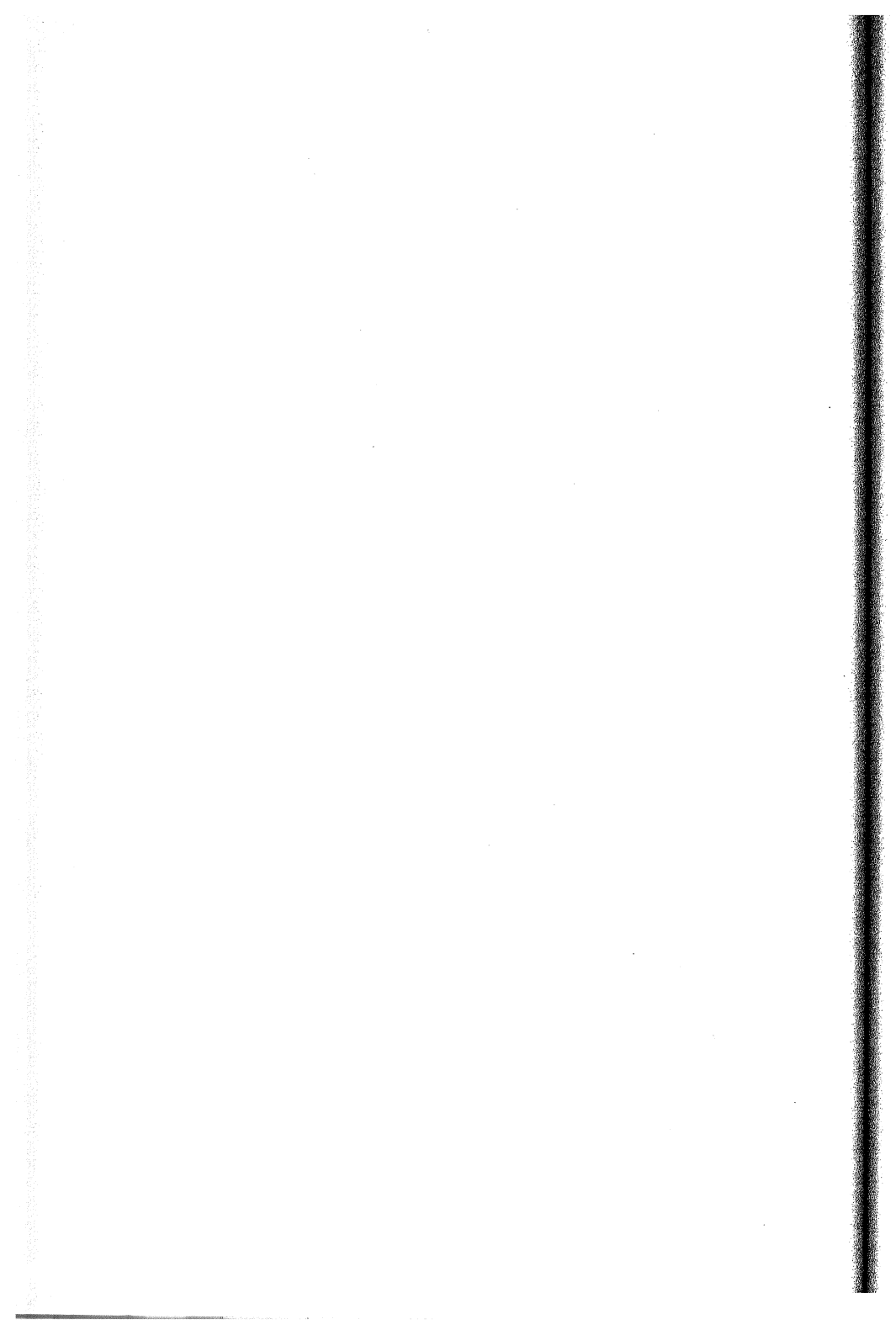
$$\text{Var}(Z) = \frac{1}{\lambda^2(1 - \alpha^2)}.$$

- (d) Is this chain reversible? *Hint: You do not need to work out the stationary density!*
- (e) Simulate the time evolution of the chain by writing appropriate functions for R. Discard the initial burn-in phase from a long run and then use the rest to get empirical estimates of the mean and variance of the stationary distribution. Use this to verify your calculations from (c) for a couple of different combinations of λ and α .
3. Starting with the R function for the simulation of the immigration-death process, modify it to include “births” in addition to immigrations and deaths. Use your modified function to simulate a long realisation from the birth-immigration-death process where the birth, death, and immigration rates are all one, starting from an initial condition of zero.

4. Use stochastic simulation to investigate the stationary distribution of the diffusion approximation to the immigration-death process. How well does it approximate the Poisson distribution of the original discrete process? By solving $E(dX_t) = 0$, show that the stationary mean of the diffusion approximation is λ/μ . Then (this one is slightly tricky) by solving $\text{Var}(X_t) = \text{Var}(X_t + dX_t)$, show that the stationary variance is also λ/μ .

5.7 Further reading

The literature on the theory of Markov processes is vast. For further information, it is sensible to start with classic texts such as Cox & Miller (1977) and Ross (1996), before moving on to applied texts such as Allen (2003), Gillespie (1992*a*), Van Kampen (1992), and references therein. The theory of diffusion processes and stochastic differential equations is particularly technical. An excellent starting point in this area, which is more accessible than most, is Øksendal (2003), but it is not appropriate for novices. For numerical methods related to SDEs, the standard text is Kloeden & Platen (1992).

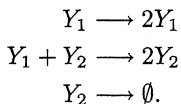


Chemical and biochemical kinetics

6.1 Classical continuous deterministic chemical kinetics

6.1.1 Introduction

Chemical kinetics is concerned with the time-evolution of a reaction system specified by a given set of coupled chemical reactions. In particular, it is concerned with system behaviour away from equilibrium. In order to introduce the concepts it is helpful to use a very simple model system. Consider the “Lotka-Volterra” (LV) system introduced in Section 1.6,



Although the reaction equations capture the key interactions between the competing species, on their own they are not enough to determine the full dynamic behaviour of the system. For that, we need to know the *rates* at which each of the reactions occurs (together with some suitable initial conditions).

6.1.2 Mass-action kinetics

The above model encourages us to think about the number of prey (Y_1) and predators (Y_2) as integers, which can change only by discrete (integer) amounts when a reaction *event* occurs. This picture is entirely correct, and we will study the implications of such an interpretation later in this chapter. However, we will introduce the study of kinetics by thinking about a more classical chemical reaction setting of macroscopic amounts of chemicals reacting in a “beaker of water.” There, the amount of each chemical is generally regarded as a *concentration*, measured in (say) moles per litre, M , which can vary continuously as the reaction progresses. Conventionally, the concentration of a chemical species X is denoted $[X]$.

It is generally the case that the instantaneous rate of a reaction is directly proportional to the concentration (in turn directly proportional to mass) of each reactant raised to the power of its stoichiometry. We will see the reason behind this when we study stochastic kinetics later, but for now we will accept it as an empirical law. This kinetic “law” is known as *mass-action kinetics*. So, for the LV system, the second reaction will proceed at a rate proportional to $[Y_1][Y_2]$. Consequently, due to the effect of this reaction, $[Y_1]$ will decrease at instantaneous rate $k_2[Y_1][Y_2]$ (where k_2 is the constant of proportionality for this reaction), and $[Y_2]$ will increase at the same

rate (since the overall effect of the reaction is to decrease $[Y_1]$ at the same rate $[Y_2]$ increases). $k_2[Y_1][Y_2]$ is known as the *rate law* of the reaction, and k_2 is the *rate constant*. Considering all three reactions, we can write down a set of ordinary differential equations (ODEs) for the system:

$$\begin{aligned}\frac{d[Y_1]}{dt} &= k_1[Y_1] - k_2[Y_1][Y_2] \\ \frac{d[Y_2]}{dt} &= k_2[Y_1][Y_2] - k_3[Y_2].\end{aligned}$$

The three rate constants, k_1 , k_2 , and k_3 (measured in appropriate units) must be specified, as well as the initial concentrations of each species. Once this has been done, the entire dynamics of the system are completely determined and can be revealed by “solving” the set of ODEs, either analytically (in the rare cases where this is possible) or numerically using a computer.

It is instructive to rewrite the above ODE system in matrix form as

$$\frac{d}{dt} \begin{pmatrix} [Y_1] \\ [Y_2] \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & 1 \end{pmatrix} \begin{pmatrix} k_1[Y_1] \\ k_2[Y_1][Y_2] \\ k_3[Y_2] \end{pmatrix},$$

where the 2×3 matrix is just the stoichiometry matrix, S , of the reaction system (Definition 2.4). This leads to a general strategy for constructing ODE models from the Petri net reaction network representation discussed in Chapter 2.

6.1.3 Equilibrium

Even when the set of ODEs is analytically intractable, it may be possible to discover an “equilibrium” solution of the system by analytic (or simple numerical) means. An equilibrium solution is a set of concentrations which will not change over time, and hence can be found by solving the set of simultaneous equations formed by setting the RHS of the ODEs to zero. In the context of the LV example, this is:

$$\begin{aligned}k_1[Y_1] - k_2[Y_1][Y_2] &= 0 \\ k_2[Y_1][Y_2] - k_3[Y_2] &= 0.\end{aligned}$$

Solving these for $[Y_1]$ and $[Y_2]$ in terms of k_1 , k_2 , and k_3 gives two solutions. The first is the rather uninteresting

$$[Y_1] = 0, \quad [Y_2] = 0,$$

and the second is

$$[Y_1] = \frac{k_3}{k_2}, \quad [Y_2] = \frac{k_1}{k_2}.$$

Further analysis (beyond the scope of this text and rather tangential to it), reveals that this second solution is not unstable, and hence corresponds to a realistic stable state of the system. However, it is not an “attractive” stable state, and so there is no reason to suppose that the system will tend to this state irrespective of the starting conditions.

Despite knowing the existence of an equilibrium solution to this system, we have

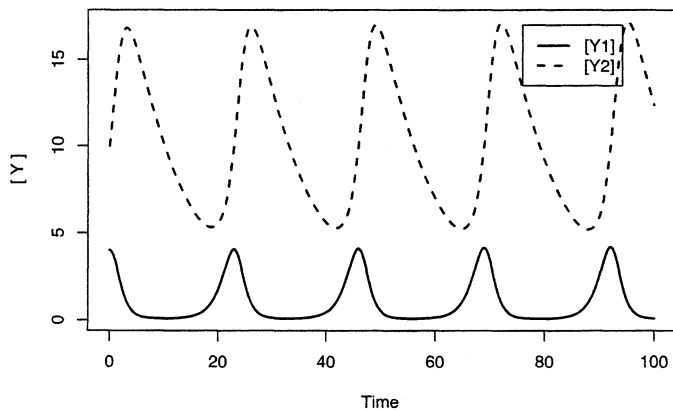
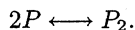


Figure 6.1 Lotka-Volterra dynamics for $[Y_1](0) = 4$, $[Y_2](0) = 10$, $k_1 = 1$, $k_2 = 0.1$, $k_3 = 0.1$. Note that the equilibrium solution for this combination of rate parameters is $[Y_1] = 1$, $[Y_2] = 10$.

no reason to suppose that any particular set of initial conditions will lead to this equilibrium, and even if we did, it would say nothing (or little) about how the system reaches it. To answer this question we need to specify the initial conditions of the system and integrate the ODEs to uncover the full dynamics. The dynamics for a particular combination of rate parameters and initial conditions are shown in Figure 6.1. An alternative way of displaying these dynamics is as an “orbit” in “phase-space” (where the value of one variable is plotted against the others, and time is not shown directly). Figure 6.2 shows the dynamics in this way.

6.1.4 Reversibility

Before going on to examine numerical integration of ODEs (which will explain how plots such as Figure 6.1 are produced), it is worth considering an important special class of reactions, namely reversible reactions. These are reactions that can proceed in both directions. For example, consider a dimerisation reaction,



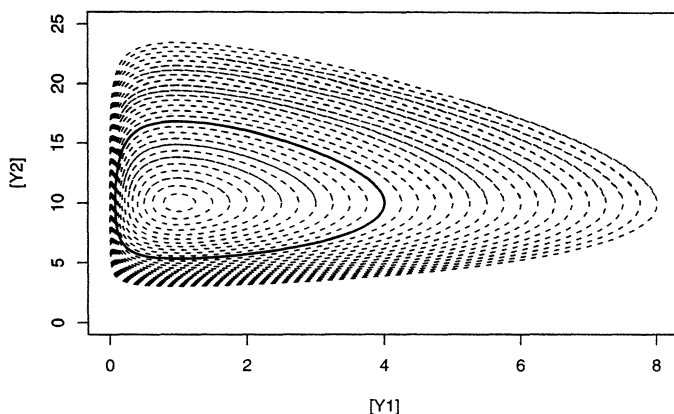


Figure 6.2 Lotka-Volterra dynamics in phase-space for rate parameters $k_1 = 1$, $k_2 = 0.1$, $k_3 = 0.1$. The dynamics for the initial condition $[Y_1](0) = 4$, $[Y_2](0) = 10$ are shown as the bold orbit. Note that the system moves around this orbit in an anti-clockwise direction. Orbits for other initial conditions are shown as dotted curves. Note that the equilibrium solution for this combination of rate parameters is $[Y_1] = 1$, $[Y_2] = 10$.

If we make the very strong assumption that *neither of these species are involved in any other reactions*, then we get the ODEs

$$\begin{aligned}\frac{d[P]}{dt} &= 2k_2[P_2] - 2k_1[P]^2 \\ \frac{d[P_2]}{dt} &= k_1[P]^2 - k_2[P_2],\end{aligned}\quad (6.1)$$

where k_1 and k_2 are the forward and backward rate constants, respectively. We clearly have equilibrium whenever

$$k_2[P_2] = k_1[P]^2.$$

Another way of writing this is

$$\frac{[P_2]}{[P]^2} = \frac{k_1}{k_2} \equiv K_{eq}\quad (6.2)$$

where K_{eq} is the *equilibrium constant* of the system. It turns out that this equilibrium is stable and attractive, as can be seen from the sample simulated dynamics in Figure 6.3.

For this system it is possible to make further progress by noting that $[P]$ and $[P_2]$ are deterministically related in this system. One way to see this is to add twice the

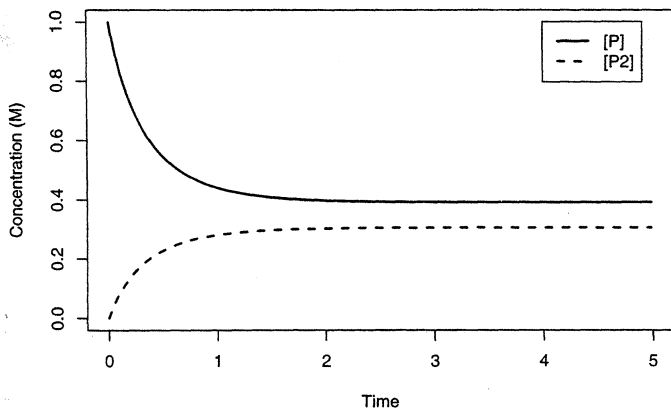


Figure 6.3 *Dimerisation kinetics for* $[P](0) = 1$, $[P_2](0) = 0$, $k_1 = 1$, $k_2 = 0.5$. This combination of parameters gives $K_{eq} = 2$, $c = 1$, and hence equilibrium concentrations of $[P] = 0.39$, $[P_2] = 0.30$.

second ODE to the first to get

$$\begin{aligned} \frac{d[P]}{dt} + 2\frac{d[P_2]}{dt} &= 0 \\ \Rightarrow \frac{d}{dt}([P] + 2[P_2]) &= 0 \\ \Rightarrow [P] + 2[P_2] &= c \end{aligned} \quad (6.3)$$

where c is the concentration of $[P]$ we would have if the dimers were fully disassociated. Equation (6.3) is known as a *conservation equation*, as the value of the LHS is conserved by the reaction system. A more direct approach to finding conservation equations is to use the Petri net theory from Chapter 2. Here, the reaction matrix (Definition 2.4) for the system is given by

$$A = \begin{pmatrix} -2 & 1 \\ 2 & -1 \end{pmatrix}.$$

The conservation equation corresponds to the P -invariant (Definition 2.5) $y = (1, 2)'$ which can be found directly by looking for non-trivial solutions of the linear system $Ay = 0$.

Conservation equations are useful for reducing the dimension of the system under consideration. Here, for example, we can rearrange (6.3) for $[P_2]$ and substitute back into the equilibrium relation (6.2) in order to find the equilibrium concentration of

$[P]$ to be a solution of the quadratic equation

$$2K_{eq}[P]^2 + [P] - c = 0.$$

This clearly has a single positive real root given by

$$[P] = \frac{\sqrt{8cK_{eq} + 1} - 1}{4K_{eq}}.$$

We can therefore find exactly the equilibrium concentrations of $[P]$ and $[P_2]$ in the absence of other reactions.

Alternatively, we can use the conservation equation (6.3) to reduce the pair of ODEs to a single first-order ODE

$$\frac{d[P]}{dt} = k_2(c - [P]) - 2k_1[P]^2, \quad (6.4)$$

which can be solved for given initial conditions to give the full dynamical behaviour of the system. It turns out that (6.4) can be solved analytically to give an explicit expression for $[P]$ as a function of t , but the solution is not elegant, and solving ODEs is not the main focus of this book. Note, however, that in order to be able to do this (analytically or otherwise), we must know both the forward and backward rate constants k_1 and k_2 , and not just their ratio, K_{eq} .

6.1.5 Numerical integration of ODEs

We will finish the section on deterministic kinetics by examining numerically integrate a system of ODEs on a computer. We will just look at this simplest possible technique, known as the first-order Euler method.* It is worth bearing in mind, however, that there are more sophisticated techniques that can be used which give much more accurate dynamics for an equivalent amount of computation time. A good example of this is the fourth-order Runge-Kutta method, which is implemented by many biochemical simulators.

We can write any of the ODE systems that we have considered so far very simply as a vector ODE

$$\frac{dX}{dt} = f(X)$$

where X is a p -dimensional vector and $f(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is an arbitrary (non-linear) p -dimensional function of X . Recall that the derivative is defined by

$$\frac{dX}{dt}(t) = \lim_{\Delta t \rightarrow 0} \frac{X(t + \Delta t) - X(t)}{\Delta t}.$$

So for small Δt we have

$$\frac{X(t + \Delta t) - X(t)}{\Delta t} \simeq f(X(t)),$$

* In fact, this technique is a special case of the Euler method for numerically solving SDEs that was examined in Section 5.5.

```
euler <- function(t=50, dt=0.001, fun=f, ic=c(1,1), ...)
{
  p=length(ic)
  n=t/dt
  xmat=matrix(0,ncol=p,nrow=n)
  x=ic
  xmat[1,]=x
  for (i in 2:n) {
    x = x + fun(x, ...) * dt
    xmat[i,]=x
  }
  ts(xmat,start=0,deltat=dt)
}
```

Figure 6.4 An R function to numerically integrate a system of coupled ODEs using a simple first-order Euler method

and rearranging gives

$$X(t + \Delta t) \simeq X(t) + \Delta t f(X(t)). \quad (6.5)$$

Equation (6.5) gives us a simple method for computing $X(t + \Delta t)$ from $X(t)$. If we start off at known $X(0)$, we can compute $X(\Delta t)$, $X(2\Delta t)$, $X(3\Delta t)$, ... to get the full dynamics of the system. This is the so-called Euler method.

A very simple R function to implement this algorithm is given in Figure 6.4. The function returns an R time series object, which can easily be plotted. For example, to simulate the dynamics of the LV system, first define the function

```
lv <- function(x, k1=1, k2=0.1, k3=0.1)
{
  c( k1*x[1] - k2*x[1]*x[2] ,
     k2*x[1]*x[2] - k3*x[2] )
}
```

Then typing `plot(euler(t=100, fun=lv, ic=c(4, 10)))` gives a plot similar to the one shown in Figure 6.1. For more sophisticated numerical integration strategies, consult a standard numerical analysis text such as Burden & Faires (2000). Note that the R package `odesolve` (available from CRAN) provides an interface to a fairly sophisticated ODE solving library.

6.2 Molecular approach to kinetics

The deterministic approach to kinetics fails to capture the discrete and stochastic nature of chemical kinetics at low concentrations. As many intra-cellular processes involve reactions at extremely low concentrations, such discrete stochastic effects are often relevant for systems biology models. We are now in a position to see how chemical kinetics can be modelled in this way.

Consider a bi-molecular reaction of the form



(the RHS is not important). What this reaction really means is that a molecule of X is able to react with a molecule of Y if the pair happen to collide with one another (with sufficient energy), while moving around randomly, driven by Brownian motion. Considering a single pair of such molecules in a container of volume V , it is possible to use statistical mechanical arguments to understand the hazard of molecules colliding. Under fairly weak assumptions regarding the container and its contents (essentially that it is small or well stirred, and in thermal equilibrium), it can be rigorously demonstrated that the collision hazard is *constant*, provided the volume is fixed and the temperature is constant. A comprehensive treatment of this issue is given in Gillespie (1992*b*), to which the reader is referred for further details. However, the essence of the argument is that as the molecules are uniformly distributed throughout the volume and this distribution does not depend on time, then the probability that the molecules are within reaction distance is also independent of time. In the case of time varying V (which can be quite relevant in the biological context), the hazard is inversely proportional to V . Again, for the careful statistical mechanical argument see Gillespie (1992*b*), but an intuitive explanation can be given as follows. Let the molecules position in space be denoted by P_1 and P_2 , respectively. Then P_1 and P_2 are uniformly and independently distributed over the volume V . This means that for a region of space d with volume v' we have

$$P(P_i \in d) = \frac{v'}{V}, \quad i = 1, 2.$$

Now if we are interested in the probability that X and Y are within a reacting distance (r) of one another at any given instant of time (assuming that r is very small relative to the dimensions of the container, so that boundary effects can be ignored), this probability can be computed as

$$P(|P_1 - P_2| < r) = E(P(|P_1 - P_2| < r | P_2)) \quad (\text{by Proposition 3.11})$$

but the conditional probability will be the same for any P_2 away from the boundary, rendering the expectation redundant, and reducing the expression to

$$\begin{aligned} &= P(|P_1 - p| < r) \quad (\text{for any } p \text{ away from the boundary}) \\ &= P(P_1 \in d) \quad (\text{where } d \text{ is a sphere of radius } r) \\ &= \frac{4\pi r^3}{3V}. \end{aligned}$$

This probability is inversely proportional to V . Then conditional on the molecules being within reaction distance, they will not necessarily react, but will do so with a probability independent of V (as other important variables, such as the velocity distributions, are independent of V), thus preserving the inverse dependence on V in the combined probability of being within reaction distance and reacting. The case of time varying V is a little messy to deal with, so a detailed discussion will be deferred until Section 8.2.3. A fixed volume V will be assumed throughout the rest

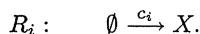
of this chapter. The case of non-constant temperature or other environmental factors are likely not to have such a straightforward relationship to the reaction hazard, and so again, throughout this chapter, it will be assumed that temperature, pressure, and all other environmental factors not explicitly described in the reaction network are held constant.

6.3 Mass-action stochastic kinetics

We will consider a system of reactions involving u species X_1, X_2, \dots, X_u and v reactions, R_1, R_2, \dots, R_v . Typically there will be more reactions than species, $v > u$. We will assume that the qualitative structure of the reaction network can be encoded in the form of a Petri net $N = (P, T, Pre, Post, M)$, where $P = (X_1, X_2, \dots, X_u)'$ and $T = (R_1, R_2, \dots, R_v)'$, as described in Section 2.3. In addition, each reaction, R_i , will have a *stochastic rate constant*, c_i , and an associated *rate law* (or *hazard function*), $h_i(x, c_i)$, where $x = (x_1, x_2, \dots, x_u)$ is the current state (or marking) of the system. The form of $h_i(x, c_i)$ (and the interpretation of the rate constant c_i) is determined by the order of reaction R_i . In all cases the hazard function has the same interpretation, namely that conditional on the state being x at time t , the probability that an R_i reaction (or transition) will occur in the time interval $(t, t + dt]$ is given by $h_i(x, c_i) dt$. Thus, *in the absence of any other reactions taking place*, the time to such a reaction event would be an $Exp(h_i(x, c_i))$ random quantity. Note, however, that since the hazard depends on the state x , and other reactions could change the state, the actually time until an R_i reaction most likely will not be exponential at all.

6.3.1 Zeroth-order reactions

First consider a reaction of the form

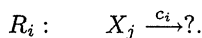


Although in practice things are not created from nothing, it can sometimes be useful to model a constant rate of production of a chemical species (or influx from another compartment) via a zeroth-order reaction. In this case, c_i is the hazard of a reaction of this type occurring, and so

$$h_i(x, c_i) = c_i.$$

6.3.2 First-order reactions

Consider the first-order reaction



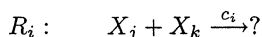
Here, c_i represents the hazard that a particular molecule of X_j will undergo the reaction. However, there are x_j molecules of X_j , each of which having a hazard of c_i of reacting. This gives a combined hazard of

$$h_i(x, c_i) = c_i x_j$$

for a reaction of this type. Note that first-order reactions of this nature are intended to capture the spontaneous change of a molecule into one or more other molecules, such as radioactive decay, or the spontaneous dissociation of a complex molecule into simpler molecules. It is not intended to model the conversion of one molecule into another in the presence of a catalyst, as this is really a second-order reaction. However, in the presence of a large pool of catalyst that can be considered not to vary in level during the time evolution of the reaction network, a first-order reaction may provide a reasonable approximation. Michaelis-Menten enzyme kinetics will be examined in detail in the next chapter.

6.3.3 Second-order reactions

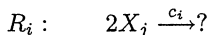
For second-order reactions of the form



c_i represents the hazard that a particular pair of molecules of type X_j and X_k will react. But since there are x_j molecules of X_j and x_k molecules of X_k , there are $x_j x_k$ different pairs of molecules of this type, and so this gives a combined hazard of

$$h_i(x, c_i) = c_i x_j x_k$$

for this type of reaction. There is another type of second-order reaction which needs to be considered:



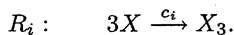
Again, c_i represents the hazard of a particular pair of molecules reacting. But here there are only $x_j(x_j - 1)/2$ pairs of molecules of type X_j , and so

$$h_i(x, c_i) = c_i \frac{x_j(x_j - 1)}{2}.$$

Note that this does not match exactly the form of the corresponding deterministic mass-action rate law — this will be further discussed in Section 6.6.

6.3.4 Higher-order reactions

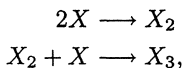
It is straightforward to extend this theory to higher-order reactions, but in reality, most (if not all) reactions that are normally written as a single reaction of order higher than two, in fact represent the combined effect of two or more reactions of order one or two. In these cases it is usually better to model the reactions in detail rather than via high-order stochastic kinetics. Consider, for example, a trimerisation reaction



Taken at face value, the rate constant c_i should represent the hazard of triples of molecules of X coming together simultaneously and reacting, leading to a rate law of the form

$$h(x, c_i) = c_i \binom{x}{3} = c_i \frac{x!}{(x-3)!3!} = c_i \frac{x(x-1)(x-2)}{6}.$$

However, in most cases it is likely to be more realistic to model the process as the pair of second-order reactions



and this pair of second-order reactions will have quite different dynamics to the corresponding third-order system.

6.4 The Gillespie algorithm

The discussion in the previous sections shows that the time-evolution of a reaction system can be regarded as a stochastic process. Further, due to the fact that the reaction hazards depend only on the current state of the system (the number of molecules of each type), it is clear that the time-evolution of the state of the reaction system can be regarded as a continuous time Markov process with a discrete state space. Detailed mathematical analysis of such systems is usually intractable, but stochastic simulation of the time-evolution of the system is quite straightforward.

In a given reaction system with v reactions, we know that the hazard for a type i reaction is $h_i(x, c_i)$, so the hazard for a reaction of some type occurring is

$$h_0(x, c) \equiv \sum_{i=1}^v h_i(x, c_i).$$

We now follow the discrete event stochastic simulation procedure from Section 5.4.2 to update the state of the process. It is clear that the time to the next reaction is $Exp(h_0(x, c))$, and also that this reaction will be a random type, picked with probabilities proportional to the $h_i(x, c_i)$, independent of the time to the next event. That is, the reaction type will be i with probability $h_i(x, c_i)/h_0(x, c)$. Using the time to the next event and the event type, the state of the system can be updated, and simulation can continue. In the context of chemical kinetics, this standard discrete event simulation procedure is known as “the Gillespie algorithm” (or “Gillespie’s direct method”), after Gillespie (1977).[†] The algorithm can be summarised as follows:

The Gillespie algorithm

1. Initialise the system at $t = 0$ with rate constants c_1, c_2, \dots, c_v and initial numbers of molecules for each species, x_1, x_2, \dots, x_u .
2. For each $i = 1, 2, \dots, v$, calculate $h_i(x, c_i)$ based on the current state, x .
3. Calculate $h_0(x, c) \equiv \sum_{i=1}^v h_i(x, c_i)$, the combined reaction hazard.
4. Simulate time to next event, t' , as an $Exp(h_0(x, c))$ random quantity.
5. Put $t := t + t'$.
6. Simulate the reaction index, j , as a discrete random quantity with probabilities $h_i(x, c_i) / h_0(x, c)$, $i = 1, 2, \dots, v$.

[†] Despite its publication date, this paper is still well worth reading.

```

gillespie <- function(N, n, ...)
{
  tt=0
  x=N$M
  S=t(N$Post-N$Pre)
  u=nrow(S)
  v=ncol(S)
  tvec=vector("numeric",n)
  xmat=matrix(0,ncol=u,nrow=n+1)
  xmat[1,]=x
  for (i in 1:n) {
    h=N$h(x, ...)
    tt=tt+rexp(1,sum(h))
    j=sample(v,1,prob=h)
    x=x+S[,j]
    tvec[i]=tt
    xmat[i+1,]=x
  }
  return(list(t=tvec,x=xmat))
}

```

Figure 6.5 An R function to implement the Gillespie algorithm for a stochastic Petri net representation of a coupled chemical reaction system

7. Update x according to reaction j . That is, put $x := x + S^{(j)}$, where $S^{(j)}$ denotes the j th column of the stoichiometry matrix S .
8. Output x and t .
9. If $t < T_{max}$, return to step 2.

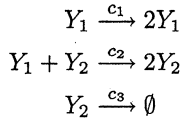
6.5 Stochastic Petri nets (SPNs)

A *stochastic Petri net* (SPN) is a Petri net where the state (represented by the number of tokens at each node) changes dynamically and randomly by choosing event firings in a carefully prescribed random manner (Goss & Peccoud 1998).

If the Petri net is used to describe a chemical reaction network (Section 2.3), where the number of tokens at a node represents the number of molecules of a given type, then the rates of event firings are given by stochastic rate laws, and the Gillespie algorithm can be used to determine the time to the next event firing and which event to fire. So a SPN is simply a convenient mathematical and graphical representation of a stochastic kinetic process.

A very simple R function for simulating the time-evolution of a SPN using the Gillespie algorithm is given in Figure 6.5. To illustrate its use, consider the stochastic kinetic formulation of the Lotka-Volterra system.

Here we will use the usual equations



leading to stochastic rate laws

$$\begin{aligned} h_1(y, c_1) &= c_1 y_1 \\ h_2(y, c_2) &= c_2 y_1 y_2 \\ h_3(y, c_3) &= c_3 y_2. \end{aligned}$$

The SPN corresponding to this system could be written

$$\begin{aligned} N &= (P, T, Pre, Post, M, h, c), \quad P = (\text{Prey}, \text{Predator})', \\ T &= (\text{Prey reproduction}, \text{Predator-prey interaction}, \text{Predator death})', \\ Pre &= \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{pmatrix}, \quad Post = \begin{pmatrix} 2 & 0 \\ 0 & 2 \\ 0 & 0 \end{pmatrix}, \quad h(y, c) = (c_1 y_1, c_2 y_1 y_2, c_3 y_2)'. \end{aligned}$$

It remains only to specify the initial state of the system, M , and the vector of rate constants, c . Some R code that formulates this problem as a SPN and then simulates it using the Gillespie algorithm assuming initial state $M = (50, 100)'$ and stochastic rate constants $c = (1, 0.005, 0.6)'$ is given in Figure 6.6. Note that the vector of rate constants is called τh (for θ) rather than c , as c has special meaning in R. The output from the Gillespie algorithm consists of a list containing two items. The first item, τ , is a vector of event times, and the second item, x , is a matrix whose rows represent the state of the system immediately *before* the corresponding event time. The matrix x therefore has an additional row, corresponding to the state of the system immediately *after* the final event simulated. As shown in the Figure 6.6, this output can be used to construct R “step function” objects which can be plotted. A single realisation of this process is shown in Figure 6.7 and Figure 6.8. These should be contrasted with the corresponding deterministic kinetics (Figure 6.1 and Figure 6.2). It is clear that although the stochastic solution approximately follows the path of the deterministic phase-space orbits, it is not constrained to follow them slavishly, but rather free to wander to nearby orbits in a stochastic manner.

The SBML-shorthand for this model is given in Figure 6.9, and the full SBML is listed in Appendix A.2 (as well as on this book’s website). Software suitable for taking the SBML as input and simulating the system dynamics will be discussed in Section 6.8.

Although it is sometimes useful to have output corresponding to each event that occurs in the simulation of the reaction network, often this is not desirable. This is because for systems of realistic size and complexity, there will be a very large number of events, and all that is likely to be of interest is the state of the system on a relatively fine regular grid of time points. One approach to this problem is to time-discretise the discrete-event output *post-hoc*. An R function to ac-


```

N=list()
N$M=c(50,100)
N$Pre=matrix(c(1,0,1,1,0,1),ncol=2,byrow=TRUE)
N$Post=matrix(c(2,0,0,2,0,0),ncol=2,byrow=TRUE)
N$h=function(y,th=c(1,0.005,0.6))
{ return(c(th[1]*y[1], th[2]*y[1]*y[2], th[3]*y[2] )) }

out=gillespie(N,10000)

op=par(mfrow=c(2,2))
plot(stepfun(out$t,out$x[,1]),pch="")
plot(stepfun(out$t,out$x[,2]),pch="")
plot(out$x,type="l")
par(op)

```

Figure 6.6 Some R code to set up the LV system as a SPN and then simulate it using the Gillespie algorithm. The state of the system is initialised to 50 prey and 100 predators, and the stochastic rate constants are $c = (1, 0.005, 0.6)'$.

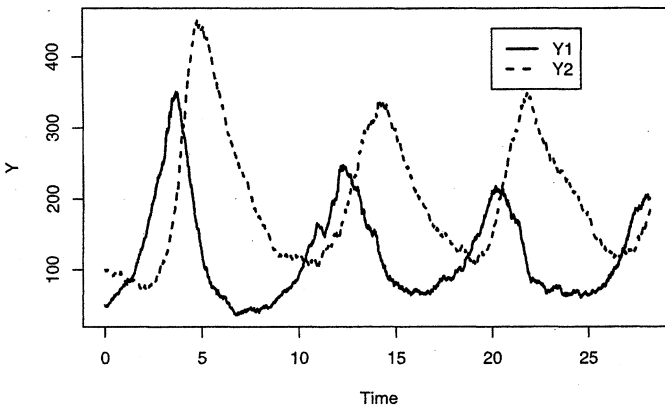


Figure 6.7 A single realisation of a stochastic LV process. The state of the system is initialised to 50 prey and 100 predators, and the stochastic rate constants are $c = (1, 0.005, 0.6)'$.

comply with this is given in Figure 6.10. This could be used with a command such as `plot(discretise(out, dt=0.01))`, where `out` is the result of a call to the `gillespie` function. However, even this solution turns out to be slightly unwieldy in practice. In particular, one cannot know in advance how many reaction events correspond to a particular length of simulation time. For this reason (and others, relating to code efficiency), it is usually better to build the discretisation

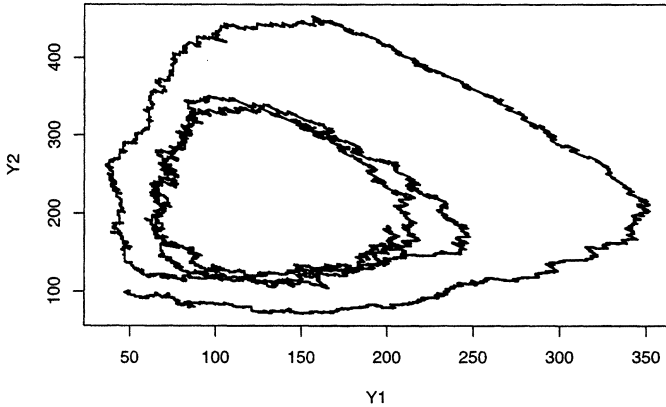


Figure 6.8 A single realisation of a stochastic LV process in phase-space. The state of the system is initialised to 50 prey and 100 predators, and the stochastic rate constants are $c = (1, 0.005, 0.6)'$.

process into the Gillespie algorithm itself. This turns out to be relatively straightforward, and a discretised version of `gillespie`, called `gillespied`, is given in Figure 6.11. It can be called with a pre-defined SPN with a command such as `plot(gillespied(N, T=100, dt=0.01))`. Note that the function `gillespied` also contains a small amount of code to gracefully handle the case of the reaction network going “extinct,” that is, reaching a point where no more reactions will occur. The previous function (Figure 6.5) would fail in that case, but can be easily modified to return something sensible in this eventuality.

Interestingly, although Petri nets are often used to model autonomous systems acting locally and in parallel (which is entirely appropriate for a reaction network), the Gillespie algorithm is “global,” in the sense that it acts on the network as a whole, and not “locally” at the level of reaction nodes. So the Gillespie algorithm does not really feel quite right in this context. However, the Gillespie algorithm is just one possible way of carrying out exact stochastic simulation of the underlying Markov process, and some of the other methods are more local in nature and more naturally lead to parallel implementations, which in turn fit better with a Petri net formulation of the problem. We will explore some of these alternative simulation strategies in Chapter 8.

6.6 Rate constant conversion

Much of the literature on biochemical reactions is dominated by a continuous deterministic view of kinetics. Consequently, where rate constants are documented, they

```

@model:2.1.1=LotkaVolterra
@units
  substance=item
@compartments
  Cell
@species
  Cell:Prey=50 s
  Cell:Predator=100 s
@reactions
@r=PreyReproduction
  Prey -> 2Prey
  c1*Prey : c1=1
@r=PredatorPreyInteraction
  Prey+Predator -> 2Predator
  c2*Prey*Predator : c2=0.005
@r=PredatorDeath
  Predator ->
  c3*Predator : c3=0.6

```

Figure 6.9 *SBML-shorthand for the stochastic Lotka-Volterra system*

```

discretise <- function(out, dt=1, start=0)
{
  events=length(out$t)
  end=out$t[events]
  len=(end-start)/%dt+1
  x=matrix(0,nrow=len,ncol=ncol(out$x))
  target=0
  j=1
  for (i in 1:events) {
    while (out$t[i]>=target) {
      x[j,]=out$x[i,]
      j=j+1
      target=target+dt
    }
  }
  ts(x,start=0,deltat=dt)
}

```

Figure 6.10 *An R function to discretise the output of gillespie onto a regular grid of time points. The result is returned as an R multivariate time series object.*

are usually deterministic rate constants, k . In order to carry out a stochastic simulation, these constants must be converted in an appropriate way to stochastic rate constants, c , representing molecular reaction hazards. In order to make this conversion, we need to understand the relationship between the deterministic and stochastic kinetic models.

```

gillespie <- function(N, T=100, dt=1, ...)
{
  tt=0
  n=T%/%dt
  x=N$M
  S=t(N$Post-N$Pre)
  u=nrow(S)
  v=ncol(S)
  xmat=matrix(0,ncol=u,nrow=n)
  i=1
  target=0
  repeat {
    h=N$h(x, ...)
    h0=sum(h)
    if (h0<1e-10)
      tt=1e99
    else
      tt=tt+rexp(1,h0)
    while (tt>=target) {
      xmat[i,]=x
      i=i+1
      target=target+dt
      if (i>n)
        return(ts(xmat,start=0,deltat=dt))
    }
    j=sample(v,1,prob=h)
    x=x+S[,j]
  }
}

```

Figure 6.11 An R function to implement the Gillespie algorithm for a SPN, recording the state on a regular grid of time points. The result is returned as an R multivariate time series object.

6.6.1 Concentrations to molecule numbers

The first issue that needs to be addressed is the difference in the representation of the *amount* of any species. In the stochastic model, this is an integer representing the number of molecules of the species, but in the deterministic model, it is usually a concentration, measured in M (moles per litre). In order to carry out the conversion from concentration to numbers of molecules, we also need to know the volume of the container, V , measured in litres.

Then for a concentration of X of $[X]$ M in a volume of V litres, there are clearly $[X]V$ moles of X and hence $n_A[X]V$ molecules, where $n_A \simeq 6.023 \times 10^{23}$ is Avogadro's constant (the number of molecules in a mole).

Example

Consider the following example, based loosely on an example from Bower & Bolouri (2000). An *E. coli* cell is a rod-shaped bacterium $2\mu\text{m}$ long with a diameter of $1\mu\text{m}$. Its volume is therefore

$$\begin{aligned} V &= \pi r^2 l \\ &= \pi(0.5 \times 10^{-6})^2 (2 \times 10^{-6}) \\ &= \frac{\pi}{2} \times 10^{-18} \text{ m}^3 \\ &= \frac{\pi}{2} \times 10^{-15} \text{ L.} \end{aligned}$$

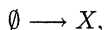
Now if a chemical species X has a concentration $[X] = 10^{-5} \text{ M}$ within an *E. coli* cell, the number of molecules is

$$n_A[X]V = 6.023 \times 10^{23} \times 10^{-5} \times \frac{\pi}{2} \times 10^{-15} = 9,461.$$

Once we are happy with converting amounts, we can think about converting rate constants.

6.6.2 Zeroth order

For the reaction



the deterministic rate law is $k \text{ Ms}^{-1}$, and so for a volume V , X is produced at a rate of $n_A k V$ molecules per second. As the stochastic rate law is just c molecules per second, we have

$$c = n_A V k.$$

6.6.3 First order

For the reaction



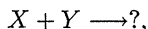
the deterministic rate law is $k[X] \text{ Ms}^{-1}$. As this involves $[X]$, we need to know that for a volume V , a concentration of $[X]$ corresponds to $x = n_A[X]V$ molecules. Now since X decreases at rate $n_A k [X]V = kx$ molecules per second, and the stochastic rate law is cx molecules per second we have

$$c = k.$$

That is, for first-order reactions, the stochastic and deterministic rate constants are always equal.

6.6.4 Second order

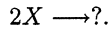
For the reaction



the deterministic rate law is $k[X][Y] \text{ Ms}^{-1}$. Here, for a volume V , the reaction proceeds at a rate of $n_A k[X][Y]V = kxy/(n_A V)$ molecules per second. Since the stochastic rate law is cxy molecules per second, we have

$$c = \frac{k}{n_A V}.$$

We also need to consider dimerisation-style reactions, of the form



Here the deterministic rate law is $k[X]^2$, so the concentration of X decreases at rate $n_A 2k[X]^2 V = 2kx^2/(n_A V)$ molecules per second. Now the stochastic rate law is $cx(x-1)/2$ so that molecules of X are consumed at a rate of $cx(x-1)$ molecules per second. Now these two laws do not match, but for large x , $x(x-1)$ can be approximated by x^2 , and so to the extent that the kinetics match, we have

$$c = \frac{2k}{n_A V}.$$

Note the additional factor of two in this case.

6.6.5 Higher order

It should be fairly clear how to extend this analysis to higher-order reactions, but such reactions are not often used in stochastic kinetic models.

6.7 The Master equation

In the stochastic kinetics literature there is often reference to “the (chemical) master equation.” Unfortunately this seems to be an overused (and misused) term, and seems now to apply to any set of differential equations whose solution gives the full transition probability kernel for the system dynamics. We have already seen sets of differential equations that determine the time evolution of the transition kernel of a Markov process — the Kolmogorov differential equations (5.7, 5.8). It turns out that the set of differential equations most often labelled as the “master equation” is just Kolmogorov’s forward equation for a stochastic kinetic process.

Proposition 6.1 *Kolmogorov’s forward equations (5.8) for a SPN can be written in the form*

$$\frac{d}{dt} p(x_0, t_0, x, t) = \sum_{i=1}^v \left[h_i(x - S^{(i)}, c_i) p(x_0, t_0, x - S^{(i)}, t) - h_i(x, c_i) p(x_0, t_0, x, t) \right], \quad \forall t_0 \in \mathbb{R}, x_0, x \in \mathcal{M},$$

where \mathcal{M} is the countable state space of the process (the set of all possible markings of the SPN). This set of differential equations is often referred to as the chemical master equation.

Proof. We start with the forward equation for a move from (x_0, t_0) to (x, t) , and then expand $q_{x,x}$ as follows,

$$\begin{aligned} \frac{d}{dt}p(x_0, t_0, x, t) &= \sum_{\{x' \in \mathcal{M}\}} q_{x',x}p(x_0, t_0, x', t) \\ &= \left[\sum_{\{x' \in \mathcal{M} | x' \neq x\}} q_{x',x}p(x_0, t_0, x', t) \right] + q_{x,x}p(x_0, t_0, x, t) \\ &= \left[\sum_{\{x' \in \mathcal{M} | x' \neq x\}} q_{x',x}p(x_0, t_0, x', t) \right] - p(x_0, t_0, x, t) \sum_{\{x' \in \mathcal{M} | x' \neq x\}} q_{x,x'} \\ &= \sum_{\{x' \in \mathcal{M} | x' \neq x\}} \left[q_{x',x}p(x_0, t_0, x', t) - q_{x,x'}p(x_0, t_0, x, t) \right]. \end{aligned}$$

This is just Kolmogorov's forward equation rewritten more appropriately for a Markov process with general countable state space \mathcal{M} . The equation involves a sum over all possible transitions, but for a SPN, only a finite number of transition events are possible, corresponding to the v different reaction channels. Considering first the hazard $q_{x,x'}$, we note that starting from x , it is only possible to move to $x + S^{(i)}$, $i = 1, 2, \dots, v$, and then by definition we have $q_{x, x+S^{(i)}} = h_i(x, c_i)$. Similarly, in order to get to x , the process must have come from one of $x - S^{(i)}$, $i = 1, 2, \dots, v$, and in this case we have $q_{x-S^{(i)}, x} = h_i(x - S^{(i)}, c_i)$. Substituting these into the above equation gives the result. \square

A more extensive discussion of the chemical master equation can be found in Van Kampen (1992). We will not have much more to say about it, as the cases where it can be solved exactly and explicitly are very few in number. These special cases are examined in McQuarrie (1967). The cases that can be solved exactly are interesting for a variety of reasons, including the testing of stochastic simulation algorithms. The derivation of the stationary distribution of the immigration-death process in Section 5.4.3 could be described as a "master equation approach," and we will do something similar with the analysis of stochastic dimerisation kinetics in Chapter 7. In general, however, a master equation approach to the analysis of stochastic kinetic models of realistic size and complexity will not be possible, and then stochastic simulation will be the only practical approach to gaining insight into the system dynamics.

Before leaving the master equation, it is instructive to see that it sheds light on the relationship between the continuous deterministic formulation and the expected value of the stochastic kinetic model. In certain special cases these are the same, and we can see this by using the master equation to derive a set of differential equations

for the expected value of the stochastic kinetic model as

$$\begin{aligned}
 \frac{\partial}{\partial t} \mathbb{E}(X_t) &= \frac{\partial}{\partial t} \sum_{x \in \mathcal{M}} x p(x, t) \\
 &= \sum_{x \in \mathcal{M}} x \frac{\partial}{\partial t} p(x, t) \\
 &= \sum_{x \in \mathcal{M}} x \sum_{i=1}^v \left[h_i(x - S^{(i)}, c_i) p(x - S^{(i)}, t) - h_i(x, c_i) p(x, t) \right] \\
 &= \sum_{i=1}^v \left[\sum_{x \in \mathcal{M}} x h_i(x - S^{(i)}, c_i) p(x - S^{(i)}, t) - \sum_{x \in \mathcal{M}} x h_i(x, c_i) p(x, t) \right] \\
 &= \sum_{i=1}^v \left[\sum_{x \in \mathcal{M}} (x + S^{(i)}) h_i(x, c_i) p(x, t) - \sum_{x \in \mathcal{M}} x h_i(x, c_i) p(x, t) \right] \\
 &= \sum_{i=1}^v \left[\mathbb{E} \left((X_t + S^{(i)}) h_i(X_t, c_i) \right) - \mathbb{E} \left(X_t h_i(X_t, c_i) \right) \right] \\
 &= \sum_{i=1}^v \mathbb{E} \left(S^{(i)} h_i(X_t, c_i) \right) \\
 &= \sum_{i=1}^v S^{(i)} \mathbb{E} \left(h_i(X_t, c_i) \right).
 \end{aligned}$$

Now in general it is not possible to solve this set of differential equations directly, but in the case where all reactions have zero- or first-order mass action rate laws we can use the linearity of expectation to get $\mathbb{E}(h_i(X_t, c_i)) = h_i(\mathbb{E}(X_t), c_i)$ giving

$$\frac{\partial}{\partial t} \mathbb{E}(X_t) = \sum_{i=1}^v S^{(i)} h_i(\mathbb{E}(X_t), c_i).$$

Putting $y(t) = \mathbb{E}(X_t)$ we get

$$\frac{d}{dt} y(t) = \sum_{i=1}^v S^{(i)} h_i(y(t), c_i) = Sh(y(t), c),$$

which is just the ODE system for the deterministic model.[‡] So, when all reactions are zero- and first-order mass-action kinetics, the deterministic solution will correctly describe the expected value of the stochastic kinetic model. However, it will not give any insight into variability, and in any case, cannot be used to describe the expectation of any system containing second-order reactions.

[‡] Note that this ODE model uses mass rather than concentration units (and the mass is measured in molecules rather than moles), and uses stochastic rate constants.

6.8 Software for simulating stochastic kinetic networks

Although it is very instructive to develop simple algorithms for simulating the dynamic evolution of biochemical networks using a high-level language such as R, this approach will not scale well to large, complex networks with many species and many reaction channels. For such models, it will be desirable to encode them in SBML, and then import them into simulation software designed with such models in mind. Such “industrial-strength” simulators are often developed in fast compiled languages such as C/C++, and written carefully to be memory efficient, accurate, and fast.

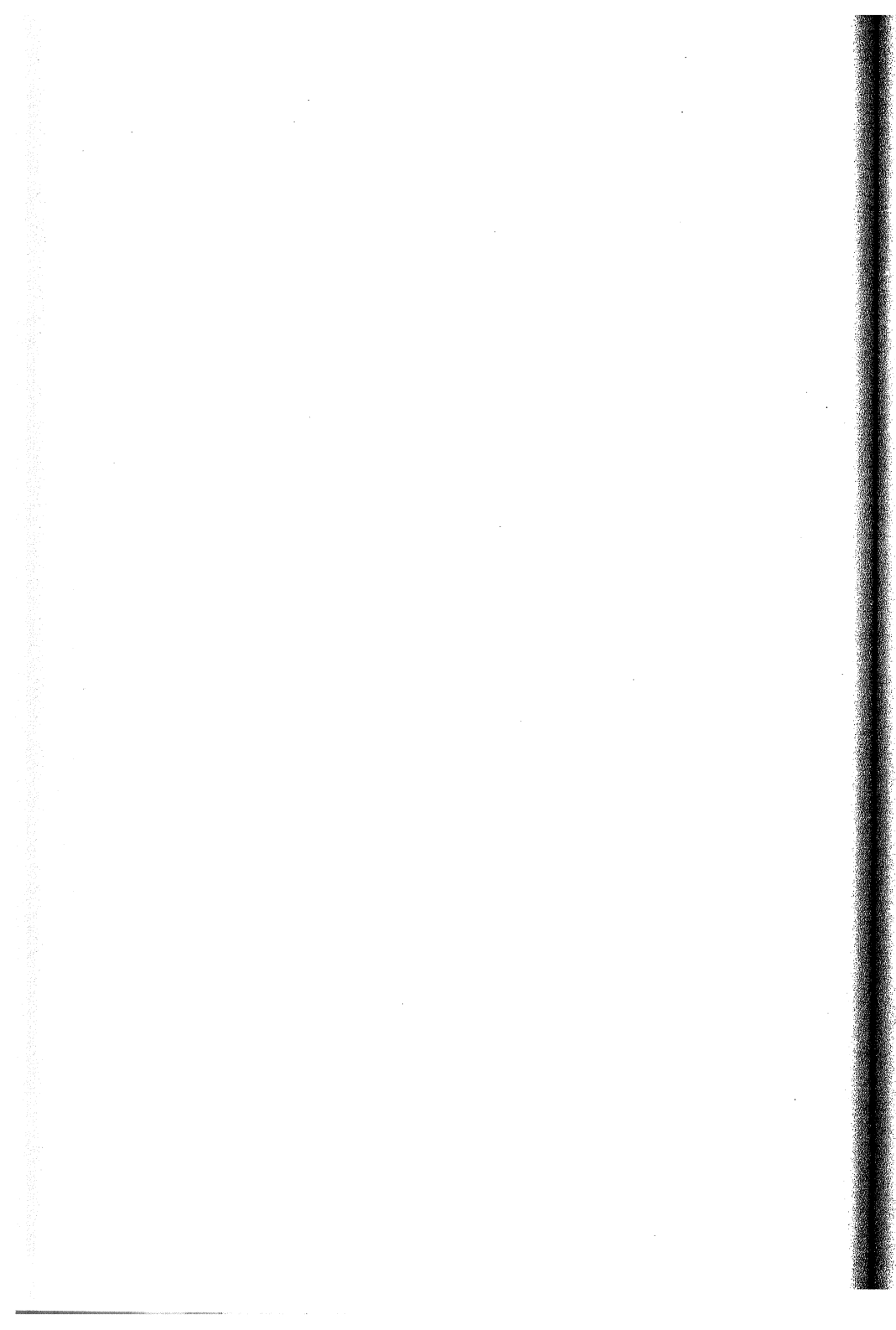
There are many software systems available for simulating the continuous deterministic kinetics corresponding to an SBML model (many such packages are listed on the SBML.org web page). However (at the time of writing), when it comes to discrete stochastic simulation there is less choice, and there are serious problems of lack of rigorous testing and inconsistencies in interpretation of SBML as one moves from one software product to another. There are several reasons for this. The deterministic framework has been around longer and is still more widely used than the stochastic framework, and so SBML was originally designed with continuous deterministic modelling in mind and involved the authors of several of the standard deterministic simulators. There is also a fairly sophisticated SBML test suite for deterministic simulators, which can be used to ensure the correctness of SBML interpretation and simulator correctness in the deterministic case. On the other hand, SBML Level 1 was not entirely appropriate for encoding discrete stochastic models (as discussed in Chapter 2), and so developers of stochastic simulators had to either ignore SBML or adopt their own particular conventions for the interpretation of SBML. Although it is possible to correctly encode a discrete stochastic model in SBML Level 2 in an unambiguous way, there is little in the way of guidance on this issue in version 1 of the specification, so the “Level 1 effect” has not yet gone away. This problem is further compounded by the difficulty of testing stochastic simulators, and the lack of a good test suite for stochastic simulators. In response to this, I have co-developed a test suite for discrete stochastic simulators that support SBML Level 2 (there is a link to it from this book’s website). It consists of a range of simple models for which direct analytic analysis of the implied stochastic process is possible (using a master equation approach), together with time-course data on the mean and standard deviation of the output that would be expected from many runs of the simulator on the same model. One simulator that passes most of the tests in the suite is a simulator known as `gillespie2`, developed for the BASIS project, described in Kirkwood et al. (2003), and based on a core stochastic simulation engine I wrote in C. Another well-known and highly regarded stochastic simulator is known as “Dizzy,” which is written in Java. However, at the time of writing, it only accepts models encoded in SBML Level 1 (in addition to its own native format), but this may well change in the future. There is also a simple Gillespie simulation service bundled with the Systems Biology Workbench (SBW).

6.9 Exercises

1. Define the auto-regulatory network from Chapter 2 (given as SBML in Appendix A.1) as a SPN in R, and then use the function `gillespie2` to simulate its time-course behaviour.
2. Write an R function to simulate the time-course behaviour of the LV system 1,000 times and compute the sample mean of the prey and predator numbers at times 1, 2, ..., 20. Plot them.
3. Install some SBML-compliant stochastic simulation software. Use it to simulate the auto-regulatory network from Appendix A.1. Check that the results seem consistent with those obtained from Exercise 1.
4. Simulate the auto-regulatory model using a deterministic simulator and interpret the results.
5. Consider the dimerisation kinetics example with parameters and initial conditions as given in Figure 6.3. Assume that the reaction is taking place in a compartment with volume V . By converting to a stochastic kinetic model and simulating, discover how small the volume V needs to be in order for stochastic effects to become important and prevalent.

6.10 Further reading

Cornish-Bowden (2004) is a classic text on modelling biochemical reactions from a continuous deterministic perspective. For numerical techniques for integrating ODEs, start with a basic numerical analysis text such as Burden & Faires (2000). Kitano (2001) and Bower & Bolouri (2000) give a good overview of systems biology and the role that biochemical network modelling and simulation has to play in it. The latter also includes a couple of chapters on stochastic simulation. To understand the role of Dan Gillespie in making physical scientists aware of the utility of stochastic modelling and simulation techniques, it is worth reading Gillespie (1977, 1992*a*, and 1992*b*). Another important book for physical scientists is Van Kampen (1992). Finally, seminal papers on the modelling of stochastic kinetic effects in the regulation of gene expression include McAdams & Arkin (1997) and Arkin et al. (1998).



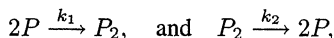
Case studies

7.1 Introduction

This chapter will focus on how the theory developed so far can be used in practice by applying it to a range of illustrative examples. The examples are relatively “small” compared to the kinds of models that systems biologists are mainly interested in (see McAdams & Arkin (1997) and Arkin et al. (1998) for a couple of good early examples), but smaller models tend to be more effective for elucidating key principles. Although the examples themselves are interesting in their own right, each one will be used to address particular modelling issues that arise in practice. So, Section 7.2 will illustrate the conversion of a deterministic model to a stochastic model, and the visualisation, summarisation, and analysis of the output of stochastic simulators. Section 7.3 will discuss conservation laws and dimensionality reduction, Section 7.4 will illustrate sensitivity and uncertainty analysis for stochastic models, and Section 7.5 will examine the analysis of external interventions (known as *events* in the SBML world).

7.2 Dimerisation kinetics

Let us consider in greater detail the problem of dimerisation kinetics briefly examined in Section 6.1.4. For this problem we will consider the dimerisation kinetics of a protein, P , at very low concentrations in a bacterial cell. We will begin by considering the usual continuous deterministic kinetics and then go on to examine the corresponding stochastic kinetic behaviour of the system. The forward and backward reactions respectively are



where k_1 and k_2 denote the usual deterministic mass-action kinetic rate constants, leading to the ordinary differential equations (6.1) for the time evolution of the system. We will assume an initial concentration of p_0 M (moles per litre) for P at time $t = 0$, and an initial concentration of 0 for P_2 . Although knowing the volume of the container (in this case the bacterium) is not strictly necessary for either a deterministic or stochastic analysis, it is required in order to compare the two. So here we assume a volume of V L (litres, or dm^3).

For the particular problem we are interested in, the initial concentration of P is $0.5 \mu\text{M}$, giving $p_0 = 5 \times 10^{-7}$, and the volume of the bacterium is $V = 10^{-15}$. It will be assumed that the values of k_1 and k_2 have been determined from a macroscopic experiment and found to be $k_1 = 5 \times 10^5$, $k_2 = 0.2$. The SBML-shorthand that en-

```

@model:2.1.1=DimerKineticsDet "Dimer Kinetics (deterministic)"
@compartments
  Cell=1e-15
@species
  Cell:[P]=5e-7
  Cell:[P2]=0
@reactions
@r=Dimerisation
  2P->P2
  Cell*k1*P*P : k1=5e5
@r=Dissociation
  P2->2P
  Cell*k2*P2 : k2=0.2

```

Figure 7.1 *SBML-shorthand for the dimerisation kinetics model (continuous deterministic version)*

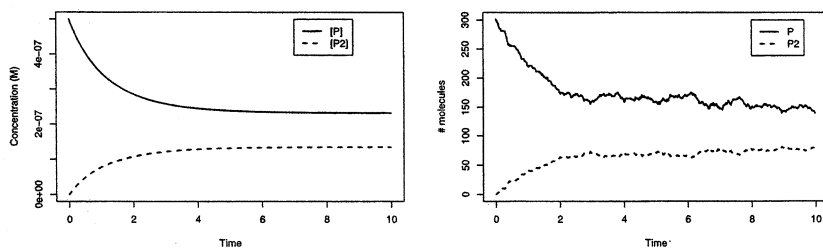


Figure 7.2 *Left: Simulated continuous deterministic dynamics of the dimerisation kinetics model. Right: A simulated realisation of the discrete stochastic dynamics of the dimerisation kinetics model.*

codes this model is given in Figure 7.1 (and the full SBML is listed in Appendix A.3). Note the additional factor of $Cell$ in the rate laws. This is interpreted as the volume of the container and is necessary because SBML rate laws are expected to be in units of substance (here moles) per unit time, and not concentration per unit time, which is how continuous deterministic rate laws are traditionally written.

The dynamics associated with this model can be simulated either by integrating the ODEs directly or by using an SBML-compliant simulator to give the dynamics shown in Figure 7.2 (left). However, we know from the stochastic kinetic theory developed in the previous chapter that for reactions involving species at low concentration in small volumes, the continuous deterministic formulation is a poor approximation to the true stochastic kinetic behaviour of the system. We therefore now turn our attention to recasting the above model in the stochastic kinetic framework in order to be able to study it from this perspective.

Algebraically, the initial amount of P is $n_{AP_0}V$ molecules, and the stochastic rate

```

@model:2.1.1=DimerKineticsStoch "Dimer Kinetics (stochastic)"
@units
  substance=item
@compartments
  Cell=1e-15
@species
  Cell:P=301 s
  Cell:P2=0 s
@reactions
@r=Dimerisation
  2P->P2
  c1*P*(P-1)/2 : c1=1.66e-3
@r=Dissociation
  P2->2P
  c2*P2 : c2=0.2
    
```

Figure 7.3 *SBML-shorthand for the dimerisation kinetics model (discrete stochastic version)*

constants are obtained as $c_1 = 2k_1/(n_A V)$, $c_2 = k_2$, using the results from Section 6.6. For our particular constants, this gives an initial value of 301 molecules of P (and no molecules of P_2) and stochastic rate constants $c_1 = 1.66 \times 10^{-3}$, $c_2 = 0.2$. The SBML-shorthand encoding of the discrete stochastic version of this model is given in Figure 7.3 (and the full SBML is listed in Appendix A.3.2). Note that in principle it should be possible for an SBML-aware software tool to automatically convert the SBML in Appendix A.3.2 into the SBML listed in Appendix A.3.1 (and possibly even the reverse, though this is harder). However, careful study of the two models reveals that this is not quite as trivial as it might first seem, as it will require the tool to have a fairly deep understanding of SBML units and semantics, as well as the relation between deterministic and stochastic rate laws, and the ability to recognise mass-action rate laws (possibly written in slightly different ways). Consequently, at the time of writing, it is typically easier to make the conversion by hand. It is also worth noting that there is nothing particularly discrete or stochastic about the discrete stochastic version of the model. A continuous deterministic simulator with a good understanding of SBML units should be able to correctly simulate the dynamics of the model described in Figure 7.3 to give output similar to that shown in Figure 7.2 (left). In practice, however, there are few (if any) simulator tools with this degree of SBML compliance at the time of writing.

We now have the model encoded in a suitable format for studying the stochastic dynamics. We will see shortly that this model is analytically tractable. However, we will ignore this fact for the present and instead use stochastic simulation as our primary investigative tool. The dynamics can be simulated using an SBML Level 2 stochastic simulator (such as `gillespie2`), or by encoding it as a SPN for `R` and using the `R` functions from Chapter 6. An appropriate SPN object for `R` can be built with the commands given in Figure 7.4.

Regardless of the tool used to conduct the simulation, the results should be the same. A single realisation of the process is given in Figure 7.2 (right). Note again

2015-11-11 11:11:11

```

N=list()
N$M=c(301,0)
N$Pre=matrix(c(2,0,0,1),ncol=2,byrow=TRUE)
N$Post=matrix(c(0,1,2,0),ncol=2,byrow=TRUE)
N$h=function(x,th=c(1.66e-3,0.2)) {
    return( c( th[1]*x[1]*(x[1]-1)/2 ,
              th[2]*x[2] ) )
}

```

Figure 7.4 R code to build an SPN object representing the dimerisation kinetics model

that a different realisation will be obtained each time the simulation is run (provided the random number seed is not fixed). Comparing this to Figure 7.2 (left) it is clear that the qualitative behaviour is somewhat similar, but that the stochastic fluctuations are very pronounced. Again it must be emphasised that these fluctuations are intrinsic to the system and have nothing to do with experimental measurement error (which we have not yet considered at all). Obviously the scales are different in the two cases, but if desired, the stochastic output can easily be mapped onto a concentration scale by dividing through by $n_A V$. A different realisation of the process is shown on this scale in Figure 7.5 (left). These two realisations are not really sufficient to give good insight into the range of behaviour that the model is likely to exhibit. This insight is typically obtained by running the simulation model many times and summarising the output in a sensible way. If we focus now on P (remember that P and P_2 are deterministically related, so it is not necessary to consider both), we begin to understand the range of dynamics it exhibits by running a relatively small number of simulations and overlaying the trajectories (Figure 7.5, right). A more sophisticated approach is to carry out many runs of the simulator and summarise the *distribution* of the level by computing appropriate statistics (such as the sample mean and variance). These can then be used to produce a plot such as Figure 7.6 (left), which is based on 1,000 runs of the simulator. Note that although it is the case here that the results of the deterministic analysis are consistent with the mean of the stochastic kinetic study, this is not true in general (and even here, the mean of the stochastic process is not exactly the deterministic solution, but it provides a reasonable approximation). In general the deterministic analysis provides no useful information about the stochastic kinetic analysis. If the marginal distribution at each time point was normal (which clearly is not the case, but is often a reasonable approximation), we would then expect over 99% of realisations to lie within 3 standard deviations of the mean. It therefore seems reasonable to use the sample mean plus/minus 3 sample standard deviations as a guide to the range of likely values at each time point. For a more detailed understanding of the distribution of values at particular time points, it is possible to study the full set of realisations at a given time. For example, Figure 7.6 (right) shows the empirical PMF for the distribution of P at time $t = 10$. Note the discrete nature of the distribution and the fact that only odd-numbered values are possible (the process was initialised at an odd number of molecules and can only change by a multiple of two). Also note how “normal” the distribution of realisations appears to be, justifying

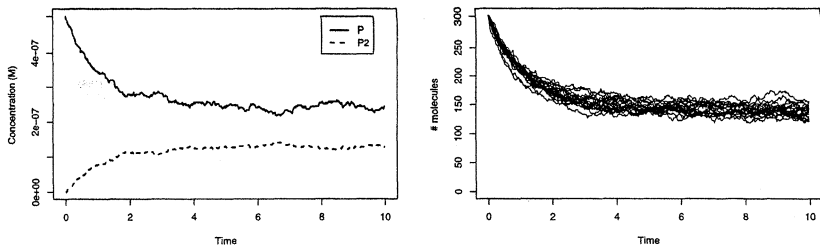


Figure 7.5 Left: A simulated realisation of the discrete stochastic dynamics of the dimerisation kinetics model plotted on a concentration scale. Right: The trajectories for levels of P from 20 runs overlaid.

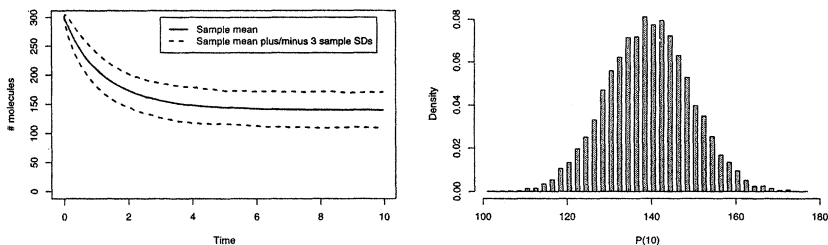


Figure 7.6 Left: The mean trajectory of P together with some approximate (point-wise) “confidence bounds” based on 1,000 runs of the simulator. Right: Density histogram of the simulated realisations of P at time $t = 10$ based on 10,000 runs, giving an estimate of the PMF for $P(10)$.

the use of a “mean plus/minus 3 SD” approach to summarising the output. Also note that Figure 7.6 (left) suggests that the system seems to have converged to a stationary distribution at time $t = 10$, and hence Figure 7.6 (right) is essentially the equilibrium distribution of the process.

For a problem as simple as this one, it is possible to make a direct analytic attack on the equilibrium probability distribution. In this case, it is made most straightforward by using the deterministic relationship between the number of molecules of P and P_2 to reduce the two-dimensional state space to a one-dimensional state space, which can then be analysed in a similar way to the immigration-death model from Chapter 5. So to start, we put

$$P + 2P_2 = n,$$

where n is the number of molecules of P that would be present if they were fully disassociated (so $n = 301$ in our example). Next we can rewrite the rate laws in

ACCEPTED MANUSCRIPT

terms of P_2 only, which we will denote by x ,

$$h_1(x, c_1) = c_1 \frac{(n-2x)(n-2x-1)}{2}$$

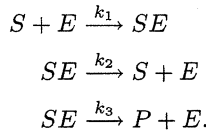
$$h_2(x, c_2) = c_2 x$$

$$h_0(x, c) = c_1 \frac{(n-2x)(n-2x-1)}{2} + c_2 x.$$

Since reaction 1 corresponds to increasing x by 1 and reaction 2 corresponds to decreasing x by 1, the transition rate matrix Q again has a tri-diagonal form. Note that here the state space (and corresponding Q matrix) is finite, as the size of x cannot exceed $n/2$. Numerical or analytical analysis of this matrix can then yield information regarding the dynamic behaviour of the system. This is essentially what physical scientists would refer to as a "Master equation approach" (although they probably would not approach the problem in exactly this way). However, the actual analysis is somewhat technical and not particularly relevant in the context of this book. It is important to bear in mind that although analytic analysis of simple processes is intellectually attractive and can sometimes give insight into more complex problems, the class of models where analytic approaches are possible is very restricted and does not cover any models of serious interest in the context of systems biology (where we are typically interested in the complex interactions between several intricate mechanisms). Therefore, computationally intensive study based on stochastic simulation and analysis is the only realistic way to gain insight into system dynamics in general.

7.3 Michaelis-Menten enzyme kinetics

Another reaction system worthy of special study is the Michaelis-Menten enzyme kinetic system. Here a substrate S is converted to a product P only in the presence of a catalyst E (for enzyme). A plausible model for this is



If we assume that k_1 , k_2 , and k_3 are deterministic mass-action kinetic rate constants, then the ODEs governing the deterministic dynamics are

$$\begin{aligned} \frac{d[S]}{dt} &= k_2[SE] - k_1[S][E] \\ \frac{d[E]}{dt} &= (k_2 + k_3)[SE] - k_1[S][E] \\ \frac{d[SE]}{dt} &= k_1[S][E] - (k_2 + k_3)[SE] \\ \frac{d[P]}{dt} &= k_3[SE]. \end{aligned}$$

```

@model:2.1.1=MMkineticsDet "M-M Kinetics (deterministic)"
@compartments
  Cell=1e-15
@species
  Cell:[S]=5e-7
  Cell:[E]=2e-7
  Cell:[SE]=0
  Cell:[P]=0
@reactions
@r=Binding
  S+E->SE
  Cell*k1*S*E : k1=1e6
@r=Dissociation
  SE->S+E
  Cell*k2*SE : k2=1e-4
@r=Conversion
  SE->P+E
  Cell*k3*SE : k3=0.1

```

Figure 7.7 *SBML-shorthand for the Michaelis-Menten kinetics model (continuous deterministic version)*

This ODE system is most easily constructed from its matrix representation,

$$\frac{d}{dt} \begin{pmatrix} [S] \\ [E] \\ [SE] \\ [P] \end{pmatrix} = \begin{pmatrix} -1 & 1 & 0 \\ -1 & 1 & 1 \\ 1 & -1 & -1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} k_1[S][E] \\ k_2[SE] \\ k_3[SE] \end{pmatrix}.$$

For a given set of rate constants and initial conditions, we can integrate this system numerically on a computer.

Again we will assume the setting of low concentrations in small volumes. The compartmental volume is $V = 10^{-15}$ L, the initial concentrations of S and E will be 5×10^{-7} M and 2×10^{-7} M respectively, and the initial concentrations of SE and P will be zero. The model specification is completed with the three rate constants $k_1 = 1 \times 10^{-6}$, $k_2 = 1 \times 10^{-4}$, $k_3 = 0.1$. The SBML-shorthand for this model is given in Figure 7.7. The simulated dynamics for this model are shown in Figure 7.8 (left).

It is clear from this plot that there are conservation laws in this system (it is particularly clear that the sum of $[E]$ and $[SE]$ is constant). Such laws can be used to reduce the dimensionality of the system under consideration. Recalling the Petri net theory from Chapter 2, the reaction matrix has the form

$$A = \begin{pmatrix} -1 & -1 & 1 & 0 \\ 1 & 1 & -1 & 0 \\ 0 & 1 & -1 & 1 \end{pmatrix}.$$

It is clear that this matrix is rank deficient (the first two rows sum to zero), so to find

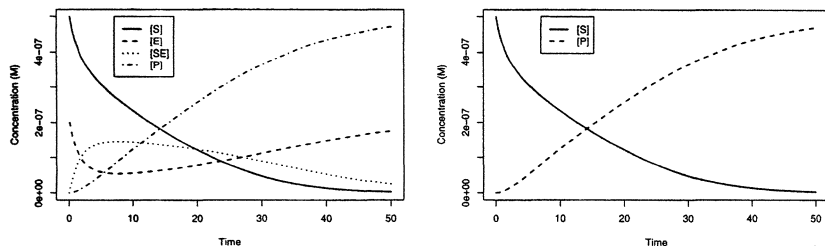


Figure 7.8 Left: Simulated continuous deterministic dynamics of the Michaelis-Menten kinetics model. Right: Simulated continuous deterministic dynamics of the Michaelis-Menten kinetics model based on the two-dimensional representation.

the P -invariants we just need to solve the linear system

$$\begin{pmatrix} 1 & 1 & -1 & 0 \\ 0 & 1 & -1 & 1 \end{pmatrix} y = 0.$$

Straightforward Gaussian elimination leads to the two invariants $y = (0, 1, 1, 0)'$ and $(1, 0, 1, 1)'$, corresponding to the conservation laws

$$\begin{aligned} [E] + [SE] &= e_0 \\ [S] + [SE] + [P] &= s_0, \end{aligned}$$

where the conservation constants are determined from the initial conditions of the system (here, e_0 is the initial value for $[E]$ and s_0 is the initial value for $[S]$). The first conservation law can be used to eliminate $[E]$ from the ODE system, then the second can be used to eliminate $[SE]$ giving

$$\begin{aligned} \frac{d[S]}{dt} &= k_2(s_0 - [S] - [P]) - k_1[S](e_0 - s_0 + [S] + [P]) \\ \frac{d[P]}{dt} &= k_3(s_0 - [S] - [P]). \end{aligned}$$

This two-dimensional system of ODEs is exactly equivalent to the original (apparently) four-dimensional system. To confirm this, the simulated dynamics for this system are given in Figure 7.8 (right) (note that the missing components can be easily reconstructed using the conservation laws if required). Such dimension-reduction techniques are particularly important in the continuous-deterministic context. First, mathematical analysis of the system in most cases requires an ODE system of full-rank. Second (of more direct practical relevance), reducing the dimension of the system will improve the speed, accuracy, and general numerical stability of the ODE-integration algorithm.

It is straightforward to convert the Michaelis-Menten system to a discrete stochastic model. The SBML-shorthand for the conversion is given in Figure 7.9, and a single realisation of the process is given in Figure 7.10 (left). Dimensionality reduc-

```

@model:2.1.1=MMKineticsStoch "M-M Kinetics (stochastic)"
@units
  substance=item
@compartments
  Cell=1e-15
@species
  Cell:S=301 s
  Cell:E=120 s
  Cell:SE=0 s
  Cell:P=0 s
@reactions
@r=Binding
  S+E->SE
  c1*S*E : c1=1.66e-3
@r=Dissociation
  SE->S+E
  c2*SE : c2=1e-4
@r=Conversion
  SE->P+E
  c3*SE : c3=0.1

```

Figure 7.9 *SBML-shorthand for the Michaelis-Menten kinetics model (discrete stochastic version)*

tion techniques can also be used in the context of discrete stochastic models. Again, the Petri net theory can be used to identify the conservation laws of the system, and these can be used to remove E and SE from the model. The SBML-shorthand for the reduced model obtained in this way is given in Figure 7.11, and a single realisation of the process is given in Figure 7.10 (right). It is clear that a software tool could be written to automatically reduce models in this way.* This model is somewhat noteworthy in that it is a valid discrete stochastic model with rate laws that are not immediately recognisable as mass-action.

Although dimensionality reduction is clearly applicable in the context of discrete stochastic modelling, it is somewhat less important than in the continuous deterministic case. This is for two main reasons. The first is that the speed improvement obtained by working with the reduced dimension system is not that significant. The second (arguably more fundamental) reason is that exact simulation algorithms such as the Gillespie algorithm are just that — exact. There is therefore no improvement in accuracy or numerical stability to be gained by working with the reduced dimension system. That said, for some of the fast approximate algorithms to be considered in Chapter 8 (notably those that exploit diffusion or ODE approximations), dimensionality reduction is just as important as in the continuous deterministic framework.

* Indeed, using the SBML construct of *assignment rules*, it is possible to do this simply by replacing redundant species with appropriate assignment rules based on the conservation laws.

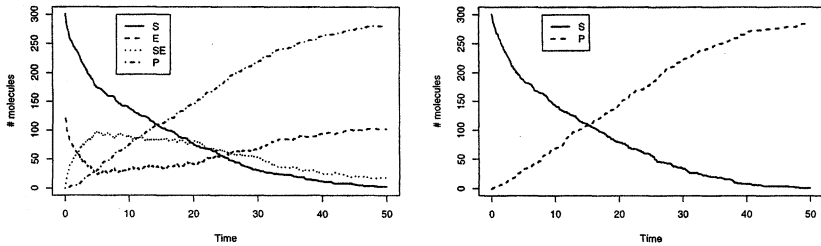


Figure 7.10 Left: A simulated realisation of the discrete stochastic dynamics of the Michaelis-Menten kinetics model. Right: A simulated realisation of the discrete stochastic dynamics of the reduced-dimension Michaelis-Menten kinetics model.

```
@model:2.1.1=RedMMKineticsStoch "Reduced M-M Kinetics (stoch)"
@units
  substance=item
@compartments
  Cell=1e-15
@species
  Cell:S=301 s
  Cell:P=0 s
@reactions
@r=Binding
  S->
  c1*S*(120-301+S+P) : c1=1.66e-3
@r=Dissociation
  ->S
  c2*(301-(S+P)) : c2=1e-4
@r=Conversion
  ->P
  c3*(301-(S+P)) : c3=0.1
```

Figure 7.11 SBML-shorthand for the reduced dimension Michaelis-Menten kinetics model (discrete stochastic version)

7.4 An auto-regulatory genetic network

The dimerisation kinetics and Michaelis-Menten kinetics models are interesting to study from a stochastic viewpoint, but both are somewhat unsatisfactory in the sense that unless the concentrations and volumes involved are really very small, the stochastic fluctuations are not particularly significant. In both cases the continuous deterministic treatment, while clearly an approximation to the truth, actually captures the most important aspects of the dynamics reasonably well. If all models of interest to systems biologists were of this nature, it would be quite proper to question whether the additional effort associated with discrete stochastic modelling is worthwhile. How-

ever, for any model where there can be only a handful (say, less than ten) of molecules of any of the key reacting species, then stochastic fluctuations can dominate, and the models can (and often do) exhibit behaviour that would be impossible to predict from the associated continuous deterministic analysis. A good example of this is the possibility of the Lotka-Volterra model to go extinct (or explode). Similar things can also happen in the context of molecular cell biology. As a trivial example, random events can trigger apoptosis or other forms of cell death. However, stochastic fluctuations are a normal part of life in the cell, which can have important consequences, and are not just associated with catastrophic events such as cell death.

A good example of a noisy process is gene expression (and its regulation). For this example we will return to the model of prokaryotic gene auto-regulation introduced in Section 1.5.7 and used as the main example throughout Chapter 2. The SBML shorthand for this model is given in Section 2.5.8, and the full SBML is listed in Appendix A.1. It should be noted that this is an artificial model, with rate constants in arbitrary units, chosen simply to make the model exhibit interesting behaviour. A simulated realisation of this process over a 5,000-second period is shown in Figure 7.12 (left). Only the three key “outputs” of the model are shown. The discrete bursty stochastic dynamics of the process are clear in this realisation. RNA transcript events are comparatively rare and random in their occurrence. The number of protein monomers oscillates wildly between 10 and 50 molecules, and the number of protein dimers jumps abruptly at random times and then gradually decays away. Looking more closely at the first 250 seconds of the same realisation (Figure 7.12, right), it is clear that the jumps in protein dimer levels coincide with the RNA transcript events. This illustrates an important point regarding stochastic variation in complex models: despite the fact that there are a relatively large number of protein dimer molecules, their behaviour is strongly stochastic due to the fact that they are affected by the number of RNA transcripts, and there are very few RNA transcript molecules in the model. Consequently, even if primary interest lies in a species with a relatively large number of molecules, a continuous deterministic model will not adequately capture its behaviour if it is affected by a species which can have a small number of molecules.

Figure 7.13 (left) shows (for the same realisation) the time-evolution of the number of molecules of protein monomers, P , over the first 10 seconds of the simulation. It is clear that even over this short time period the stochastic fluctuations are very significant. In order to understand this variation in more detail, let us now focus on the number of molecules of P at time $t = 10$. By running many simulations of the process it is possible to build up a picture of the probability distribution for the number of molecules, and this is shown in Figure 7.13 (right) (based on 10,000 runs). This clearly shows that there is an almost even chance that there will be no molecules of P at time 10 (and the most likely explanation for this is that there will not yet have been a transcription event). The distribution is clearly far from normal, so a mean plus/minus three SD summary of the distribution is unlikely to be adequate in this case.

Once a model becomes as complex as this one, there is likely to be some uncertainty regarding some quantitative aspects of the model specification. This could be,

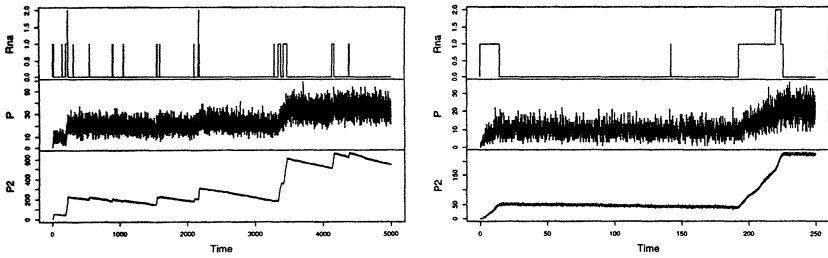


Figure 7.12 *Left: A simulated realisation of the discrete stochastic dynamics of the prokaryotic genetic auto-regulatory network model, for a period of 5,000 seconds. Right: A close-up on the first period of 250 seconds of the left plot.*

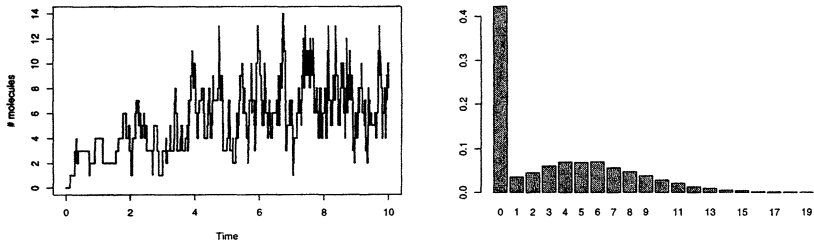


Figure 7.13 *Left: Close-up showing the time-evolution of the number of molecules of P over a 10-second period. Right: Empirical PMF for the number of molecules of P at time $t = 10$ seconds, based on 10,000 runs.*

for example, uncertainty about the initial conditions or the stochastic rate constants adopted. For concreteness, we will suppose here that there is a degree of uncertainty regarding the value of the gene transcription rate k_2 .[†] A value of $k_2 = 0.01$ was specified in the model, but let us suppose that any value between 0.005 and 0.03 is plausible. It is therefore natural to want to investigate the sensitivity of the model dynamics to this particular specification.

For continuous deterministic models, a sophisticated framework for sensitivity analysis is well established. However, the techniques from this domain do not transfer well to the stochastic modelling paradigm. One of the main motivations for thinking about model sensitivity is a desire to understand uncertainty in the true process dynamics. However, in the context of stochastic modelling, uncertainty about the process dynamics is integral to the whole approach. Consequently, even in the

[†] Note that in contradiction to the convention adopted this far in the book, k_2 is a stochastic rate constant and not a deterministic one. This is because in practice, stochastic modellers often use k rather than c for their rate constants, so it is best not to become over-reliant on this notational cue.

case of complete certainty regarding the model structure, rate laws, rate constants, and initial conditions, the time-evolution of the process is uncertain (or random, or stochastic, depending on choice of terminology). A natural way to incorporate uncertainty regarding model parameters is to adopt a subjective Bayesian interpretation of probability. In the Bayesian paradigm, uncertainty regarding model parameters is not fundamentally different to uncertainty regarding the time evolution of the process due to the stochastic kinetic dynamics. We have already constructed mechanisms for handling uncertainty in the process dynamics by trying to understand the probability distribution of the outcomes, rather than by simply looking at a particular realisation of the process. There is no reason why these probability distributions should not include uncertainty regarding the model parameters in addition to the uncertainty induced by the stochastic kinetics. This is best illustrated by example.

Figure 7.13 (right) shows our uncertainty about the level of P at time $t = 10$ based on a value of $k_2 = 0.01$. Similar plots can be produced based on other plausible values; Figure 7.14 (left) shows the plot corresponding to $k_2 = 0.02$, for example. In order to obtain a summary of our uncertainty regarding the level of P , we need to average over our uncertainty for k_2 in an appropriate way. In order to do this properly, we need to specify a probability distribution which reflects our uncertainty in k_2 . It was previously stated that all values between 0.005 and 0.03 are plausible. We will now make the much stronger assumption that all values in this range are equally plausible, which leads directly to the probability distribution $U(0.005, 0.3)$. Within the Bayesian framework, this is known as a *prior* probability distribution for a parameter.[†] Having specified the prior probability distribution (in practice, there is likely to be uncertainty regarding several parameters, but it is completely straightforward to assign independent prior probability distributions to each uncertain value), it is then straightforward to incorporate this into the subsequent analysis. Rather than running many simulations with the same parameters, each run begins by first picking uncertain parameters from their prior probability distributions. This has the effect of correctly embedding the parameter uncertainty into the model; all subsequent analysis then proceeds as normal. For example, if interest is in the level of P at time $t = 10$, the values can be recorded at the end of each run to build up a picture of the marginal uncertainty for the level of P . Such a distribution is depicted in Figure 7.14 (right). Unsurprisingly, it looks a bit like a compromise between Figure 7.13 (right) and Figure 7.14 (left).

In practice one is not simply interested in the uncertainty of the level of one particular biochemical species at one particular time, but it is clear that exactly the same approach can be applied to any numerical summary of the simulation output (for example, cell-cycle time, time to cell death, population doubling time, time for RNA expression levels to increase five-fold, etc.). When applied to derived simulation outputs of genuine biological interest (possibly directly experimentally measurable),

[†] It is known as a prior distribution because it is possible to use the model and experimental data to update this prior distribution into a *posterior* distribution, which describes the uncertainty regarding the parameter having utilised the information in the data. This is *Bayesian statistics*, and it will be discussed further in the final chapters of this book.

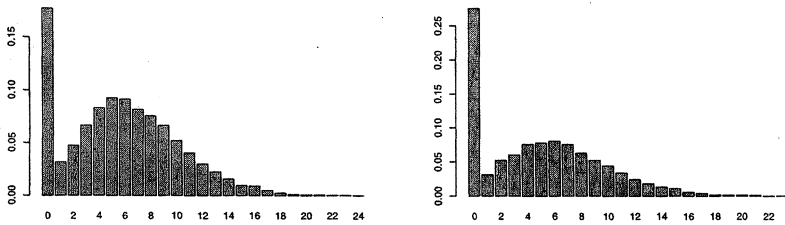


Figure 7.14 Left: Empirical PMF for the number of molecules of P at time $t = 10$ seconds when k_2 is changed from 0.01 to 0.02, based on 10,000 runs. Right: Empirical PMF for the prior predictive uncertainty regarding the observed value of P at time $t = 10$ based on the prior distribution $k_2 \sim U(0.005, 0.03)$.

Bayesian uncertainty analysis provides a powerful framework for model introspection.

7.5 The *lac* operon

We now return to the *lac* operon model introduced in Section 1.5.8. Given some appropriate rate constants, we are now in a position to completely specify this model and study its dynamics. Some SBML-shorthand that specifies (a very simplified version of) the model is given in Figure 7.15. The rate constants have been chosen to be biologically plausible (with a time unit of seconds), then fine-tuned to make the model behave sensibly. Again, the model is meant to be illustrative and does not represent a serious effort to accurately model the true dynamics of the *lac* regulation dynamics (or even the actual mechanism, as several simplifying assumptions have been made here as well). Like the auto-regulatory model, from a sensible set of initial conditions, this model will go through a transient phase then settle down to an equilibrium probability distribution which is not particularly interesting on its own. The most interesting aspect of the *lac* mechanism is the dynamic response to an influx of lactose. Whether or not such an external intervention should be regarded as being part of the model is actually somewhat controversial. However, SBML Level 2 provides a means for encoding interventions via the event element. SBML events were not discussed in Chapter 2, but the discussion in the SBML Level 2 specification document is quite readable. Events are also supported in SBML-shorthand from version 2.1.2 onward — again, see the specification document for further details. The event listed at the end of the shorthand model in Figure 7.15 has the effect of introducing 10,000 molecules of lactose into the cell at time $t = 20,000$.

Conceptually, simulating a model that includes a timed intervention of this kind is quite straightforward. In this particular case, it could be done by running the simulator until time 20,000, then recording the state at this time, adding 10,000 to the final level of Lactose, and then restarting the simulator from the new state. In practice,

```

@model:2.1.2=lacOperon "lac operon model (stochastic)"
@units
  substance=item
@compartments
  Cell=1e-15
@species
  Cell:Idna=1 s
  Cell:Irna=0 s
  Cell:I=50 s
  Cell:Op=1 s
  Cell:Rnap=100 s
  Cell:Rna=0 s
  Cell:Z=0 s
  Cell:Lactose=20 s
  Cell:ILactose=0 s
  Cell:IOp=0 s
  Cell:RnapOp=0 s
@reactions
@r=InhibitorTranscription
  Idna -> Idna + Irna
  c1*Idna : c1=0.02
@r=InhibitorTranslation
  Irna -> Irna + I
  c2*Irna : c2=0.1
@r=LactoseInhibitorBinding
  I + Lactose -> ILactose
  c3*I*Lactose : c3=0.005
@r=LactoseInhibitorDissociation
  ILactose -> I + Lactose
  c4*ILactose : c4=0.1
@r=InhibitorBinding
  I + Op -> IOp
  c5*I*Op : c5=1
@r=InhibitorDissociation
  IOp -> I + Op
  c6*IOp : c6=0.01
@r=RnapBinding
  Op + Rnap -> RnapOp
  c7*Op*Rnap : c7=0.1
@r=RnapDissociation
  RnapOp -> Op + Rnap
  c8*RnapOp : c8=0.01
@r=Transcription
  RnapOp -> Op + Rnap + Rna
  c9*RnapOp : c9=0.03
@r=Translation
  Rna -> Rna + Z
  c10*Rna : c10=0.1
@r=Conversion
  Lactose + Z -> Z
  c11*Lactose*Z : c11=1e-5
@r=InhibitorRnaDegradation
  Irna ->
  c12*Irna : c12=0.01
@r=InhibitorDegradation
  I ->
  c13*I : c13=0.002
@r=LactoseInhibitorDegradation
  ILactose -> Lactose
  c13*ILactose : c13=0.002
@r=RnaDegradation
  Rna ->
  c14*Rna : c14=0.01
@r=ZDegradation
  Z ->
  c15*Z : c15=0.001
@events
  Intervention = t>=20000 : Lactose=Lactose+10000

```

Figure 7.15 SBML-shorthand for the lac-operon model (discrete stochastic version)

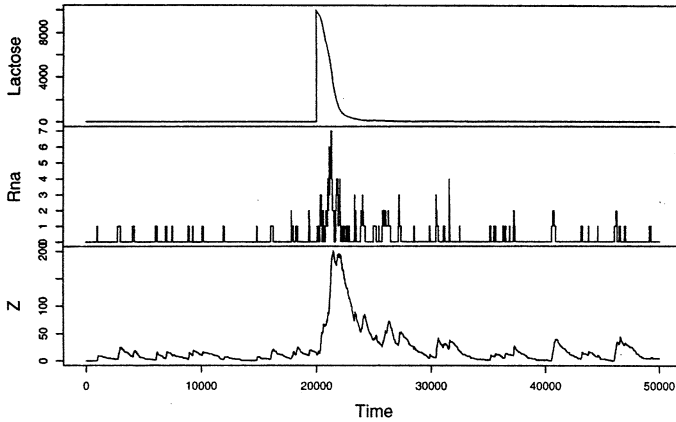


Figure 7.16 A simulated realisation of the discrete stochastic dynamics of the lac-operon model for a period of 50,000 seconds. An intervention is applied at time $t = 20,000$, when 10,000 molecules of lactose are added to the cell.

some simulators include built-in support for SBML events, which simplifies the process greatly. The realisation of the process plotted in Figure 7.16 was generated using the simulator `gillespie2`. The plot shows that under equilibrium conditions the expression level of the *lac-Z* protein is very low. However, in response to the introduction of lactose, the rate of transcription of the *lac* operon is increased, leading to a significant increase in the expression levels of the *lac-Z* protein. The result of this is that the lactose is quickly converted to something else, allowing the cell to gradually return to its equilibrium behaviour.

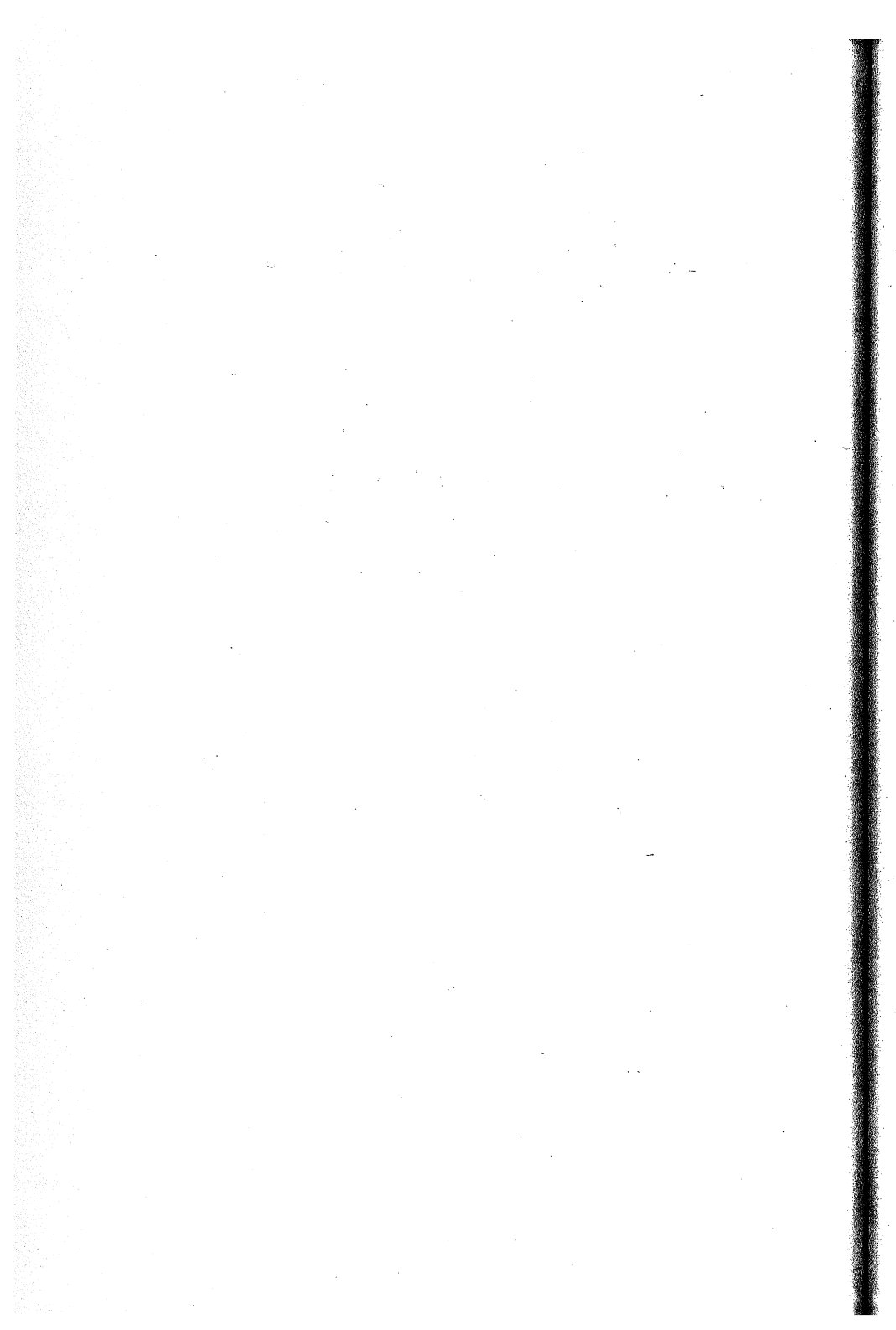
7.6 Exercises

1. Carry out the simulations described in each section of this chapter and reproduce all of the plots. This is probably the most important exercise in the entire book! Of course it will not be possible to *exactly* reproduce the plots showing individual stochastic realisations, but even here it should be possible to obtain plots showing qualitatively similar behaviour.
2. Use a Master equation approach to exactly compute and plot the exact version of Figure 7.6 (right). Also compute the exact equilibrium distribution and compare the two.
3. Identify the conservation laws in the auto-regulatory network and use them to reduce the dimension of the model. Simulate the dynamics of the reduced model to ensure they are consistent with the original.

4. Deduce the continuous deterministic version of the auto-regulatory model (assuming a container volume of 10^{-15}L) and simulate it. Compare it to the stochastic version. How well does it describe the mean behaviour of the stochastic version?
5. In the dimerisation kinetics model, the dissociation rate constant c_2 took the value 0.2. Suppose now that there is uncertainty regarding this value which can be well described by a $U(0.1, 0.4)$ probability distribution. Produce a new version of Figure 7.6 (right) which incorporates this uncertainty.

7.7 Further reading

This chapter provides only the briefest of introductions to stochastic modelling of genetic and biochemical networks, but it should provide sufficient background in order to render the literature in this area more accessible. To build large and complex stochastic models of interesting biological processes, it is necessary to read around the deterministic modelling literature in addition to the stochastic literature, as the deterministic literature is much more extensive, and can still provide useful information for building stochastic models. Good starting points include Bower & Bolouri (2000), Kitano (2001) and Klipp et al. (2005); also see the references therein. The existing literature on stochastic modelling is of course particularly valuable, and I have found the following articles especially interesting: McAdams & Arkin (1997), Arkin et al. (1998), McAdams & Arkin (1999), Goss & Peccoud (1998), Pinney et al. (2003), Hardy & Robillard (2004), Salis & Kaznessis (2005*b*). Readers with an interest in stochastic modelling of aging mechanisms might also like to read Proctor et al. (2005). In addition to the modelling literature, it is invariably necessary to trawl the wet-biology literature and online databases for information on mechanisms and kinetics — mastering this process is left as an additional exercise for the reader, though some guidance is provided in Klipp et al. (2005).



Beyond the Gillespie algorithm

8.1 Introduction

In Chapter 6 an algorithm for simulating the time-course behaviour of a stochastic kinetic model was introduced. This discrete-event simulation algorithm, usually referred to as the Gillespie algorithm, has the nice properties that it simulates every reaction event and is exact in the sense that it generates exact independent realisations of the underlying stochastic kinetic model. It is also reasonably efficient in terms of computation time among all such algorithms with those properties. However, it should be emphasised that the Gillespie algorithm is just one approach out of many that could be taken to simulating stochastic biochemical dynamics. In this chapter we will look at some other possible approaches, motivated by problems in applying the Gillespie algorithm to large models containing species with large numbers of molecules and fast reactions. In Section 8.2, refinements of the Gillespie algorithm will be considered that still generate exact realisations of the stochastic kinetic process. In Section 8.3, methods based on approximating the process by another that is faster to simulate will be examined. Then in Section 8.4, hybrid algorithms that combine both exact and approximate updating strategies will be considered.

8.2 Exact simulation methods

Let us begin by reconsidering discrete event simulation of a continuous time Markov process with a finite number of states, as described in Section 5.4 (p. 121). We saw how to do this using the *direct method* in Section 5.4.2. There is an alternative algorithm to this, known as the *first event method*, which also generates an exact realisation of the Markov process. It can be stated as follows:

1. Initialise the process at $t = 0$ with initial state i .
2. Call the current state i . For each potential other state k ($k \neq i$), compute a putative time to transition to that state $t_k \sim \text{Exp}(q_{ik})$.
3. Let j be the index of the smallest of the t_k .
4. Put $t := t + t_j$.
5. Let the new state be j .
6. Output the time t and state j .
7. If $t < T_{max}$, return to step 2.

It is perhaps not immediately obvious that this algorithm is exactly equivalent to the direct method. In order to see that it is, we need to recall some relevant properties

of the exponential distribution, presented in Section 3.9 (p. 77). By Proposition 3.18, the minimum of the putative times has the correct distribution for the time to the first event, and by Proposition 3.19, the index associated with that minimum time clearly has the correct probability mass function (PMF).

Although the first event method is just as correct as the direct method, the direct method is generally to be preferred, as it is more efficient. In particular, the direct method requires just two random numbers to be simulated per event, whereas the first event method requires r (where r is the number of states). The first event method is interesting, however, as it gives us another way of thinking about the simulation of the process, and forms the basis of a very efficient exact simulation algorithm for stochastic kinetic models.

8.2.1 First reaction method

Before looking at the Gibson-Bruck algorithm in detail, it is worth thinking briefly about how the first event method for finite state Markov processes translates to stochastic kinetic models with a potentially infinite state space. We have already seen how the direct method translates into the Gillespie algorithm (Section 6.4). The first event method translates into a variant of the Gillespie algorithm known as the *first reaction method* (Gillespie 1976), which can be stated as follows:

1. Initialise the starting point of the simulation with $t := 0$, rate constants $c = (c_1, \dots, c_v)$ and initial state $x = (x_1, \dots, x_u)$.
2. Calculate the reaction hazards $h_i(x, c_i)$, $i = 1, 2, \dots, v$.
3. Simulate a putative time to the next type i reaction, $t_i \sim \text{Exp}(h_i(x, c_i))$, $i = 1, 2, \dots, v$.
4. Let j be the index of the smallest t_i .
5. Put $t := t + t_j$.
6. Update the state x according to the reaction with index j . That is, set $x := x + \zeta^{(j)}$.
7. Output t and x .
8. If $t < T_{max}$ return to step 2.

As stated, this algorithm is clearly less efficient than Gillespie's direct method, but with a few clever tricks it can be turned into a very efficient algorithm. An R function to implement the first reaction method for a SPN is given in Figure 8.1.

8.2.2 The Gibson-Bruck algorithm

The so-called *next reaction method* (also known as the Gibson-Bruck algorithm) is a modification of the first reaction method which makes it much more efficient. A first attempt at presenting the basics of the Gibson-Bruck algorithm follows:

1. Initialise $t := 0$, c and x , and additionally calculate all of the initial reaction hazards $h_i(x, c_i)$, $i = 1, \dots, v$. Use these hazards to simulate putative first reaction times $t_i \sim \text{Exp}(h_i(x, c_i))$.

```

frm <- function(N, n, ...)
{
  tt=0
  x=N$M
  S=t(N$Post-N$Pre)
  u=nrow(S)
  v=ncol(S)
  tvec=vector("numeric",n)
  xmat=matrix(0,ncol=u,nrow=n+1)
  xmat[1,]=x
  for (i in 1:n) {
    h=N$h(x, ...)
    pu=rexp(v,h)
    j=which.min(pu)
    x=x+S[,j]
    tt=tt+pu[j]
    tvec[i]=tt
    xmat[i+1,]=x
  }
  return(list(t=tvec,x=xmat))
}

```

Figure 8.1 An R function to implement the first reaction method for a stochastic Petri net representation of a coupled chemical reaction system. It is to be used in the same way as the gillespie function from Figure 6.5.

2. Let j be the index of the smallest t_i .
3. Set $t := t_j$.
4. Update x according to reaction with index j .
5. Update $h_j(x, c_j)$ according to the new state x and simulate a new putative time $t_j := t + \text{Exp}(h_j(x, c_j))$.
6. For each reaction $i (\neq j)$ whose hazard is changed by reaction j :
 - (a) Update $h'_i = h_i(x, c_i)$ (but temporarily keep the old h_i).
 - (b) Set $t_i := t + (h_i/h'_i)(t_i - t)$.
 - (c) Forget the old h_i .
7. If $t < T_{max}$ return to step 2.

There are several things to note about this algorithm. The first is that it has moved from working with “relative” times (times from now until the next event) to “absolute” times (the time of the next event). The reason for doing this is that it saves generating new times for all of the reactions that are not affected by the reaction that has just taken place (thanks to the memoryless property, Proposition 3.16, this is okay). The second thing to note is that the times that are affected by the most recent reaction are “re-used” by appropriately rescaling the old variable (conditional

on it being greater than t). Again, a combination of the memoryless property and the rescaling property (Proposition 3.20) ensures that this is okay.

The next thing worth noting is that it is assumed that the algorithm “knows” which hazards are affected by each reaction. Gibson & Bruck (2000) suggest that this is done by creating a “dependency graph” for the system. The dependency graph has nodes corresponding to each reaction in the system. A directed arc joins node i to node j if a reaction event of type i induces a change of state that affects the hazard for the reaction of type j . These can be determined (automatically) from the forms of the associated reactions. Using this graph, if a reaction of type i occurs, the set of all children of node i in the graph gives the set of hazards that needs to be updated.

An interesting alternative to the dependency graph is to work directly on the Petri net representation of the system. Then, for a given reaction node, the set of all “neighbours” (species nodes connected to that reaction node), \mathcal{X} is the set of all species that can be altered. Then the set of all reaction nodes that are “children” of a node in \mathcal{X} is the set of all reaction nodes whose hazards may need updating. This approach is slightly conservative in that the resulting set of reaction nodes is a superset of the set which absolutely must be updated, but nevertheless represents a satisfactory alternative.

This algorithm is now “local” in the sense that all computations (bar one) involve only adjacent nodes on the associated Petri net representation of the problem. The only remaining “global” computation is the location of the index of the smallest reaction time. Gibson and Bruck’s clever solution to this problem is to keep all reaction times (and their associated indices) in an “indexed priority queue.” This is another graph, allowing searches and changes to be made using only fast and local operations; see the original paper for further details of exactly how this is done. The advantage of having local operations on the associated Petri net is that the algorithm becomes straightforward to implement in an event-driven object-oriented programming style, consistent with the ethos behind the Petri net approach. Further, such an implementation could be multi-threaded on an SMP or hyper-threading machine, and would also lend itself to a full message-passing implementation on a parallel computing cluster. For further information on parallel stochastic simulation, see Wilkinson (2005).

This algorithm is more efficient than Gillespie’s direct method in the sense that only one new random number needs to be simulated for each reaction event which takes place, as opposed to the two that are required for the Gillespie algorithm. Note however, that selective recalculation of the hazards, $h_i(x, c_i)$ (and the cumulative hazard $h_0(x, c)$), is also possible (and highly desirable) for the Gillespie algorithm, and could speed up that algorithm enormously for large systems. Given that the Gillespie algorithm otherwise requires fewer operations than the next reaction method, and does not rely on the ability to efficiently store putative reaction times in an indexed priority queue, the relative efficiency of a cleverly implemented direct method and the next reaction method is likely to depend on the precise structure of the model and the speed of the random number generator used.

8.2.3 Time-varying volume

A problem that has so far been overlooked is that of reaction hazards that vary continuously over time. The most common context for this to arise in a practical modelling situation is when a growing cell (or cellular compartment) has its volume modelled as a continuous deterministic function of time. For example, let us suppose that the container volume at time t , $V(t)$, is modelled as

$$V(t) = v_0 + \alpha t, \quad t \geq 0, \quad (8.1)$$

for some constant $\alpha > 0$. If the model contains any second-order reactions, the hazards of these should be inversely proportional to $V(t)$ (Section 6.2). The hazards of first-order reactions are independent of volume and hence unaffected. What to do about any zero-order reactions is somewhat unclear. As zero-order reactions are typically used to model “production” or “influx” in a simple-minded way, it is conceivable that at least some zero-order reactions should have volume dependence. Certain production rates might reasonably be considered to be directly proportional to $V(t)$, while influx equations might have hazards proportional to $V(t)^{2/3}$ (as surface area increases more slowly than volume). In general, zero-order reactions should be considered on a case-by-case basis.

In order to keep the presentation as straightforward as possible, we will just consider modifying the (inefficient) first reaction method to take account of time-varying reaction hazards.* It should be reasonably clear that since the only steps involving the hazards are steps 2 and 3 that only steps 2 and 3 need modification. Since the hazard is time varying, we should now write it $h_i(x, c_i, t)$, $i = 1, 2, \dots, u$. Now we could simply run the algorithm using these time-dependent hazards but otherwise unmodified. Unfortunately this will lead to an algorithm that is only approximately correct, as we would be essentially assuming that the hazards remain constant between each reaction event, which is not actually true. At this point it is helpful to recall the inhomogeneous Poisson process (Section 5.4.4), as this is exactly what is needed for dealing with non-constant hazards. Proposition 5.4 and the subsequent discussion tell us exactly how to simulate the time of the next event of an inhomogeneous Poisson process, but note that the lower limit of the integral defining the cumulative hazard must be the current simulation time, and not zero.

For concreteness consider a reaction with hazard $h(t) = a/V(t)$, where a will be a function of x and c_i , but constant with respect to t . Suppose further that $V(t)$ is given by (8.1), the current simulation time is t_0 , and we wish to simulate the time t' of the next reaction event. We begin by computing the cumulative hazard

$$H(t) = \int_{t_0}^t h(t) dt = \int_{t_0}^t \frac{a dt}{v_0 + \alpha t} = \frac{a}{\alpha} \log \left(\frac{v_0 + \alpha t}{v_0 + \alpha t_0} \right), \quad t \geq t_0,$$

and then compute the cumulative distribution function (CDF) of the time of the next

* The next reaction method (Gibson-Bruck algorithm) is also quite straightforward to modify. The direct method (Gillespie algorithm) is actually a bit awkward to modify, and so the desire to work with time-varying hazards is one reason why some people prefer to use an algorithm in the Gibson and Bruck style. These issues are discussed in some detail in Gibson & Bruck (2000).

event as

$$F(t) = 1 - \exp\{-H(t)\} = 1 - \left(\frac{v_0 + \alpha t}{v_0 + \alpha t_0}\right)^{-\alpha/\alpha}, \quad t \geq t_0.$$

Once we have the CDF we can simulate $u \sim U(0, 1)$ and solve $u = F(t')$ for t' , the time of the next event, to obtain

$$t' = \frac{1}{\alpha} \left[(v_0 + \alpha t_0) u^{-\alpha/\alpha} - v_0 \right].$$

Returning to the problem of modifying the first reaction method, one simply simulates putative times to the next event using the above strategy for any reactions with time varying rates, and the rest of the algorithm remains untouched.

This provides an example of coupling a discrete stochastic process with a variable that changes continuously in time. Essentially the same strategy is used in several of the hybrid algorithms to be considered in Section 8.4, where some variables are treated as discrete and others as varying continuously in time. A good understanding of the above technique is a necessary pre-requisite for understanding hybrid simulation algorithms.

8.3 Approximate simulation strategies

8.3.1 Time discretisation

Gibson and Bruck's next reaction method is regarded by many to be the best available method for *exact* simulation of a stochastic kinetic model. However, if one is prepared to sacrifice the exactness of the simulation procedure, there is a potential for huge speed-up at the expense of a little accuracy. These fast approximate methods are all based on a time discretisation of the Markov process.

The essential idea is that the time axis is divided into (small) discrete chunks, and the underlying kinetics are approximated so that advancement of the state from the start of one chunk to another can be made in one go. Most of the methods work on the assumption that the time intervals have been chosen to be sufficiently small that the reaction hazards can be assumed constant over the interval. We know that a point process with constant hazard is a (homogeneous) Poisson process (Section 3.6.5). Based on the definition of the Poisson process, we assume that the number of reactions (of a given type) occurring in a short time interval has a Poisson distribution (independently of other reaction types). We can then simulate Poisson numbers of reaction events and update the system accordingly.

For a fixed (small) time step Δt , we can present an approximate simulation algorithm as follows (we use the matrix notation from Section 2.3.2):

1. Initialise the problem with time $t := 0$, rate constants c , state x , and stoichiometry matrix S .
2. Calculate $h_i(x, c_i)$, for $i = 1, \dots, v$, and simulate the u -dimensional reaction vector r , with i th entry a $Po(h_i(x, c_i)\Delta t)$ random quantity.
3. Update the state according to $x := x + Sr$.

```
pts <- function(N, T=100, dt=1, ...)
{
  tt=0
  n=T/%dt
  x=N$M
  S=t(N$Post - N$Pre)
  u=nrow(S)
  v=ncol(S)
  xmat=matrix(0, ncol=u, nrow=n)
  for (i in 1:n) {
    h=N$h(x, ...)
    r=rpois(v, h*dt)
    x = x + (S %*% r)
    xmat[i,]=x
  }
  ts(xmat, start=0, deltat=dt)
}
```

Figure 8.2 An R function to implement the Poisson timestep method for a stochastic Petri net representation of a coupled chemical reaction system. It is to be used in the same way as the gillespie function from Figure 6.11.

4. Update $t := t + \Delta t$.
5. Output t and x .
6. If $t < T_{max}$ return to step 2.

We could call this the *Poisson timestep method*. Note that step 3 should (ideally) be accomplished with a sparse matrix update. An R function to implement the Poisson timestep method is given in Figure 8.2. The problem with the above method is that of choosing an appropriate timestep Δt so that the method is fast but reasonably accurate. Clearly the smaller Δt , the more accurate, and the larger Δt , the faster. Another problem is that although one particular Δt may be good enough for one part of a simulation, it may not be appropriate for another. This motivates the idea of stepping ahead a variable amount of time τ , based on the rate constants, c and the current state of the system, x . This is the idea behind Gillespie's τ -leap algorithm.

8.3.2 Gillespie's τ -leap method

The τ -leap method (Gillespie 2001) is an adaptation of the Poisson timestep method to allow stepping ahead in time by a variable amount τ , where at each timestep τ is chosen in an appropriate way in order to ensure a sensible trade-off between accuracy and speed. This is achieved by making τ as large (and hence fast) as possible while still satisfying some constraint designed to ensure accuracy. In this context, the accuracy is determined by the extent to which the assumption of constant hazard over the interval is appropriate. Clearly whenever any reaction occurs some of the reaction hazards change, and so an assessment needs to be made of the magnitude of change

of the hazards $h_i(x, c_i)$. Essentially, the idea is to choose τ so that the (proportional) change in all of the $h_i(x, c_i)$ is small.

The simplest way to check that a chosen τ is satisfactory is to apply a *post-leap* check. That is, after a leap of τ , check that $|h_i(x', c_i) - h_i(x, c_i)|$ is sufficiently small for each i (where x and x' represent the state of the system before and after the leap). If any of the differences are too large, try again with a smaller value of τ . One of the problems with this method is that it biases the system away from large yet legitimate state changes.

A *pre-leap* check seems more promising. Here we can calculate the expected new state as $E(x') = x + E(r)A$, where the i th element of $E(r)$ is just $h_i(x, c_i)\tau$. We can then calculate the change in hazard at this “expected” new state and see if this is acceptably small (it should be noted that this is *not* necessarily the expected change in hazard, due to the potential non-linearity of $h_i(x, c_i)$). It is suggested that the magnitude of acceptable change should be a fraction of the cumulative hazard $h_0(x, c)$, i.e.,

$$|h_i(x', c_i) - h_i(x, c_i)| \leq \epsilon h_0(x, c), \quad \forall i.$$

Gillespie provides an approximate method for calculating the largest τ satisfying this property (Gillespie 2001). Note that if the resulting τ is as small (or almost as small) as the expected time leap associated with an exact single reaction update, then it is preferable to do just that. Since the time to the first event is $Exp(h_0(x, c))$, which has expectation $1/h_0(x, c)$, one should always prefer an exact update if the suggested τ is less than (say) $2/h_0(x, c)$.

Gillespie & Petzold (2003) consider refinements of this basic τ selection algorithm that improve its behaviour somewhat. However, in my opinion, the “pure” τ -leap method is always likely to be somewhat unsatisfactory in the context of biochemical networks with very small numbers of molecules of some species (say, zero or one copy of an activated gene). On the other hand, a hybrid algorithm known as the *maximal timestep method* uses the τ -leap method for some variables and not others. This algorithm, which seems quite promising, will be briefly described in Section 8.4.2.

8.3.3 Diffusion approximation (chemical Langevin equation)

Another way of speeding up simulation is to simulate from the diffusion approximation to the true process (Section 5.5). A formal discussion of this procedure is beyond the scope of this book. However, we shall here content ourselves with deriving the form of the diffusion approximation using an intuitive procedure (which can be formalised with a little effort; see, for example Gillespie (2000)).

It is clear from the discussion of the Poisson timestep method that the change in state of the process, dX_t in an infinitesimally small time interval dt is $S dR_t$, where dR_t is a u -vector whose i th element is a $Po(h_i(X_t, c_i)dt)$ random quantity. Matching the mean and variance, we put

$$dR_t \simeq h(x, c)dt + \text{diag} \left\{ \sqrt{h(x, c)} \right\} dW_t,$$

where $h(x, c) = (h_1(x, c_1), \dots, h_v(x, c_v))'$, dW_t is the increment of a v -d Wiener process, and for a p -d vector v , $\text{diag} \{v\}$ denotes the $p \times p$ matrix with the elements

of v along the leading diagonal and zeros elsewhere. We now have the diffusion approximation

$$\begin{aligned} dX_t &= S dR_t \\ &= S \left(h(X_t, c)dt + \text{diag} \left\{ \sqrt{h(X_t, c)} \right\} dW_t \right) \\ \Rightarrow dX_t &= Sh(X_t, c)dt + S \text{diag} \left\{ \sqrt{h(X_t, c)} \right\} dW_t. \end{aligned} \quad (8.2)$$

Equation (8.2) is one way of writing the *chemical Langevin equation* (CLE) for a stochastic kinetic model.[†] One slightly inconvenient feature of this form of the equation is that the dimension of X_t (u) is different from the dimension of the driving process W_t (v). Since we will typically have $v > u$, there will be unnecessary redundancy in the formulation associated with this representation. However, using some straightforward multi-variate statistics (not covered in Chapter 3), it is easy to see that the variance-covariance matrix for dX_t is $S \text{diag} \{h(X_t, c)\} S' dt$ (or $A' \text{diag} \{h(X_t, c)\} A dt$). So, (8.2) can be rewritten

$$dX_t = Sh(X_t, c)dt + \sqrt{S \text{diag} \{h(X_t, c)\} S'} dW_t, \quad (8.3)$$

where dW_t now denotes the increment of a u -d Wiener process, and we use a common convention in statistics for the square root of a matrix.[‡] In some ways (8.3) represents a more efficient description of the CLE. However, there can be computational issues associated with calculating the square root of the diffusion matrix, particularly when S is rank degenerate (as is typically the case, due to conservation laws in the system), so both (8.2) and (8.3) turn out to be useful representations of the CLE, depending on the precise context of the problem, and both provide the basis for approximate simulation algorithms. An R function to integrate the CLE using the Euler method is given in Figure 8.3. Again, this method can work extremely well if there are more than (say) ten molecules of each reacting species throughout the course of the simulation. However, if the model contains species with a very small number of molecules, simulation based on a pure Langevin approximation is likely to be unsatisfactory. Again, however, a hybrid algorithm can be constructed that uses discrete updating for the low copy-number species and a Langevin approximation for the rest; such an algorithm will be discussed in Section 8.4.3.

Having used a relatively informal method to derive the CLE, it is worth taking time to understand more precisely its relationship with the true discrete stochastic kinetic model (represented by a Markov jump process). In Section 6.7 we derived the chemical Master equation as Kolmogorov's forward equation for the Markov jump process. It is possible to develop a corresponding forward equation for the chemical Langevin equation (8.3). For a general k -dimensional Langevin equation

$$dX_t = \mu(X_t)dt + \beta^{1/2}(X_t)dW_t,$$

[†] Note that setting the term in dW_t to zero gives the ODE system corresponding to the continuous deterministic approximation.

[‡] Here, for a $p \times p$ matrix M (typically symmetric), \sqrt{M} (or $M^{1/2}$) denotes any $p \times p$ matrix N satisfying $NN' = M$. Common choices for N include the symmetric square root matrix and the Cholesky factor; see Golub & Van Loan (1996) for computational details.

```

cle <- function(N, T=100, dt=0.1, ...)
{
  sdt=sqrt(dt)
  tt=0
  n=T%/dt
  x=N$M
  S=t(N$Post-N$Pre)
  u=nrow(S)
  v=ncol(S)
  xmat=matrix(0,ncol=u,nrow=n)
  for (i in 1:n) {
    h=N$h(x, ...)
    dw=rnorm(v,0,sdt)
    dx=(S%*%h)*dt + S %**% (sqrt(h)*dw)
    x = x + dx
    xmat[i,]=x
  }
  ts(xmat,start=0,deltat=dt)
}

```

Figure 8.3 An R function to integrate the CLE using an Euler method for a stochastic Petri net representation of a coupled chemical reaction system. It is to be used in the same way as the `pts` function from Figure 8.2.

this takes the form

$$\frac{\partial}{\partial t} p(x, t) = - \sum_{i=1}^k \frac{\partial}{\partial x_i} \{ \mu_i(x) p(x, t) \} + \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k \frac{\partial^2}{\partial x_i \partial x_j} \{ \beta_{ij}(x) p(x, t) \},$$

and is known as the *Fokker-Planck equation* corresponding to the Langevin. Here $p(x, t)$ represents (as a function of x) the (marginal) probability density function for X_t . Now substituting in $\mu(x) = Sh(x, c)$ and $\beta(x) = S \text{diag} \{ h(x, c) \} S'$ gives the Fokker-Planck equation for the CLE, and this can be shown to correspond to a second-order approximation to the chemical master equation. Thus the CLE in some sense represents the “best” SDE approximation to the true stochastic kinetic model. See Gillespie (2000) for further details.

8.4 Hybrid simulation strategies

Hybrid algorithms aim to bridge the gap between the exact algorithms considered in Section 8.2 and the approximate algorithms considered in Section 8.3. The essential idea is to partition model species into two (or possibly more) groups that can be treated in a similar way. The algorithms considered recently in the literature all partition the species into two groups; a group of low-copy number species that need to be treated as discrete and updated with an “exact” method, and a group of species that can be treated with an approximate algorithm of some description. We will denote these two groups X_D and X_A . It then turns out that it is necessary to partition

the set of reactions similarly into two groups, R_D and R_A (sometimes referred to as “stochastic partitioning”). This is typically done by assigning to R_D any reaction that can change the level of any species assigned to X_D (though there are several variations of this idea). Also note that the reactions in R_D are often referred to as “slow,” while those in R_A are “fast.”

8.4.1 Discrete/ODE models

We begin by looking at a method of combining discrete event simulation with conventional ODE models, based loosely on the algorithm of Kiehl, Mattheyses & Simmons (2004), to which the reader is referred for further details. Also see Alfonsi et al. (2005) for a more formal discussion of this approach to hybrid simulation. These methods are appealing to conventional mathematical modellers as they allow the approximate part of the system to be modelled with the usual continuous deterministic techniques that they are most familiar with.

The algorithm begins with a partitioning of species and reactions, and the fixing of a step size for the numerical integration algorithm for the ODE system. Note that it is possible to use any ODE solver in conjunction with this algorithm (Euler, Runge-Kutta, etc.). The basic idea is that an ODE integration step will take place on the assumption that no discrete reaction takes place. Once this has been done, time-varying hazards for the discrete reaction hazards are determined (if an Euler method is being used, the hazards will be linear in time, and if a Runge-Kutta method is used, they will be a higher-order polynomial). It is then possible to use the methods discussed in Section 5.4.4 and Section 8.2.3 to check if a discrete event occurred. If not, the timestep is done. If there is, the time of the event should be identified, the discrete update should take place, and the continuous variables should be updated over this new shorter timestep. There are various possible modifications of this procedure. In particular, the algorithm is greatly simplified if it is decided reasonable to assume that the hazards for the discrete regimes are approximately constant over the time interval of interest (though this obviously introduces additional error). Another possibility is to reduce the size of the numerical integration step dynamically in order to reduce the number of integration steps that contain a discrete reaction event. The procedure is summarised in the algorithm below:

1. Initialise the system and set $t := 0$.
2. Calculate the discrete reaction hazards at time t .
3. On the basis of these, decide on an appropriate timestep for the ODE integration, Δt .
4. Use an ODE integration step to compute the trajectory of the continuous variables over the interval $[t, t + \Delta t]$.
5. On the basis of these, compute the corresponding time-varying discrete reaction hazards, and decide if a discrete reaction event has taken place.
6. If no discrete reaction has taken place, put $t := t + \Delta t$, and update the continuous variables to the values appropriate for this new time.

7. If a discrete reaction has occurred, find the time, t_1 and type of the (first) reaction. Then put $t := t_1$, update the continuous values to those appropriate for time t_1 , and update the discrete variables according the reaction type that has occurred.
8. If $t < T_{max}$, return to step 2.

The algorithm has been presented assuming that Gillespie's direct method will be used for the discrete updating, but is easily adapted to other discrete updating schemes, such as the next reaction method. There are a number of techniques that could be used for determining an appropriate value of Δt . In Kiehl et al. (2004) it is assumed that an appropriate timestep for the ODE solver can be determined a priori as δt . Then an appropriate "discrete step size" $\delta \tau$ can be chosen, possibly as the expected time until the first discrete reaction, giving $\delta t = 1/h_0(x, c)$. Then $\Delta t = \min\{\delta t, \delta \tau\}$. In my opinion, it is likely to be beneficial to make $\delta \tau$ a bit smaller than the expected time until the first discrete reaction, in order to reduce the number of time intervals containing a discrete event.

This algorithm is fairly straightforward to understand and implement and has a similar structure to the other hybrid algorithms we will consider. It therefore provides a useful starting point for thinking about hybrid simulation strategies. The main problem with the algorithm is that the ODE approximation is too crude in the context of stochastic modelling of genetic and biochemical networks. As Kiehl et al. (2004) fully acknowledge, the algorithm has the effect of suppressing the intrinsic variation of variables assigned to the continuous regime, and this can have important consequences when the study of stochastic dynamics is of direct interest. This therefore motivates the study of hybrid algorithms that better preserve the stochastic variation of the system dynamics.

8.4.2 Maximal timestep algorithm

The *maximal timestep method* is a hybrid algorithm proposed by Puchalka & Kierzek (2004) that combines an exact updating procedure for the low concentration species with a τ -leaping approximate updating algorithm for the other species. Puchalka & Kierzek (2004) chose to use the next reaction method for the exact updating, but it would be straightforward to reformulate the algorithm using Gillespie's direct method. A greatly simplified version of the full algorithm could be described as follows:

1. Initialise the system and set $t := 0$.
2. Calculate the fast reaction hazards and use these to select an appropriate τ (for the τ -leap updating scheme).
3. Assuming a constant hazard for the slow reactions, decide if a slow reaction has taken place in $[t, t + \tau]$.
4. If no slow reaction has taken place, perform a τ -leap update on the fast reactions to $t := t + \tau$.
5. If a slow reaction has taken place, identify the time, t_1 , and type of the (first) reaction. Then put $t := t_1$, perform a τ -leap update of the fast reactions to t_1 , and update the slow variables according to the slow reaction that has occurred.

6. If $t < T_{max}$, return to step 2.

This algorithm has a similar structure to the one considered earlier, but in the context of τ -leaping, it is more difficult to take account of the time-varying nature of the reaction hazards, so it is just hoped that the timestep is small enough that a constant hazard assumption is valid. The full algorithm presented in Puchalka & Kierzek (2004) is considerably more complex than the simple overview given above, but most of the additional details are associated with dynamic repartitioning of reactions into fast and slow categories. The reader is referred to the original paper for further information. Note that these authors also make use of the quasi-steady-state assumption of Rao & Arkin (2003) in order to utilise reactions with rate laws that do not fit into the mass-action class.

The maximal timestep algorithm is quite appealing in that it is reasonably clear that provided a sufficiently small τ is used and sufficiently strict criteria are imposed for admission to the class of fast reactions, then this algorithm is capable of capturing the true stochastic kinetic dynamics arbitrarily well. However, it is not clear how well such an algorithm will scale-up to the very large complex systems that it is intended to tackle.

8.4.3 Discrete/Langevin methods

The final hybrid algorithm we will consider here is based on combining a discrete updating scheme for slow reactions with a diffusion (or Langevin) approach to updating the fast reactions. The presentation will be loosely based on that of Salis & Kaznessis (2005a), which is itself a refinement of the basic strategy described by Haseltine & Rawlings (2002). Salis & Kaznessis (2005a) also recommend an automatic partitioning and dynamic repartitioning of fast and slow reactions. Based on a timestep of Δt for the updating of the fast reactions, they suggest that reaction R_j should only be classified as fast if both

$$h_j(x, c_j)\Delta t \geq \lambda$$

and

$$x \geq \varepsilon |S^{(j)}|$$

are true for the current system state x , and suitably chosen λ and ε . The authors suggest using $\lambda = 10$ and $\varepsilon = 100$. Obviously then, as the system state x evolves over time, the set of reactions classified as fast will also change.

While Haseltine & Rawlings (2002) use a direct method for updating the slow reactions, Salis & Kaznessis (2005a) use a next reaction method, and hence refer to their combined scheme as the *next reaction hybrid* (NRH) algorithm. Again, there is nothing particularly fundamental as to the choice of slow updating scheme.

A greatly simplified version of the basic algorithm structure (which is sufficiently general to apply to both the Haseltine & Rawlings (2002) and Salis & Kaznessis (2005a) algorithms) could be stated as follows; see the original papers for further details.

1. Initialise the system and set $t := 0$.

2. Calculate the hazards for the fast reactions.
3. Numerically integrate the CLE for the fast reactions from t to $t + \Delta t$ to obtain a "sample path" for the continuous variables over the interval $[t, t + \Delta t]$.[§]
4. Using the time-dependent hazards for the slow reactions, decide whether or not a slow reaction has happened in $[t, t + \Delta t]$.
5. If no slow reaction has occurred, set $t := t + \Delta t$ and update the continuous variables to their proposed values at this time.
6. If one slow reaction has occurred, identify the time, t_1 and type, set $t := t_1$, and update the system to t_1 (using the type of the discrete reaction and the appropriate continuous state).
7. If more than one slow reaction has occurred, reduce Δt and return to step 3.
8. If $t < T_{max}$, return to step 2.

Here, steps 6 and 7 could (in principle) be replaced with 6 and 7. Identify the time and type of the first reaction and update accordingly.

The reason the algorithm is presented this way is that it is slow and difficult to identify the times of slow reactions using numerical integration and solution methods, because the reaction hazards correspond to sample paths of unknown stochastic processes. For this reason, Salis & Kaznessis (2005a) use approximate numerical techniques to estimate the reaction times that are more accurate if the time is first narrowed down to a small interval where it is known that only this one reaction occurs. Salis & Kaznessis (2005a) also present another algorithm (the ANRH algorithm), which is faster but less accurate than the NRH, as it allows multiple slow reactions to fire without affecting the fast reactions; again, see the paper for further details.

The NRH and related algorithms that couple a Markov jump process with a Langevin approximation are (in principle) very attractive as they are more accurate than the ODE hybrid algorithms and are likely to scale-up to very large models much better than the maximal timestep algorithm.

8.4.4 Discussion

There is a clear and pressing need for accurate stochastic simulation algorithms that are much faster than the exact algorithms. Although the "pure" approximate algorithms are elegant and appealing, they fail to adequately cope with the common problem of multiple reactions happening on differing time scales. This then motivates the development of hybrid algorithms that can address exactly this issue using a combination of exact and approximate techniques. Unfortunately, all three of the hybrid algorithms presented (and there are other, related algorithms that have not been explicitly covered here) are somewhat "crude" in terms of their statistical sophistication. The ODE model is unsatisfactory as the jump in scale from discrete

[§] Of course the true sample path is a diffusion process, which cannot be obtained in its entirety, and this is the main complication for algorithms of this type; see the text immediately following the algorithm for further details.

stochastic to continuous deterministic is too great. The maximal timestep method is hard to couple without assuming constant hazard over the integration timestep, which makes it heavily reliant on a small integration step size for adequate accuracy, and in any case it is not clear how well it will cope with vastly differing time scales.

Of the approaches presented, the discrete/CLE partitioning adopted by Haseltine & Rawlings (2002) seems to have the greatest potential, and the refinements proposed by Salis & Kaznessis (2005a) lead to an algorithm which should perform reasonably well in practice. However, there are many improvements that could be made to their algorithm to improve both its accuracy and its computational efficiency. Obviously, as noted by Salis & Kaznessis (2005a), better algorithms than the Euler-Maruyama method could be used for numerically integrating the fast reactions. At a more fundamental level, however, improvements could be made to the way that the discrete and continuous regimes are coupled. It turns out that exact sampling strategies can be used to decide whether or not a slow reaction has occurred in the interval, and if so, at exactly what time the first reaction takes place (this can be done efficiently without any approximate numerical methods); see Wilkinson (2006) for further details. Once this has been done, the only sources of error will be the error in adopting the CLE approximation of the fast reactions and the error in the numerical integration step. There is not much that can be done about the error in the CLE approximation within this framework (other than ensuring that only appropriate reactions are designated "fast"), but it *might* be possible to eliminate the discretisation error completely. Using techniques for exact sampling of non-linear diffusions (Beskos & Roberts 2005), it could turn out to be possible to directly and efficiently simulate the time to the first slow reaction and the state of the fast system at that time (and any other finite number of intermediate times) exactly, without any time discretisation error at all. We would then be in a situation more comparable with that of the exact discrete algorithms, where the simulated realisations are exact, conditional on the model. However, it is far from clear that the techniques currently being explored for retrospective exact sampling of diffusions will be applicable to non-linear multi-variate diffusion processes as general as the CLE.

It is worth emphasising that the literature on stochastic simulation is vast, and that there has not been any attempt here to give a comprehensive survey. An interesting area that has not been covered is that of spatial modelling of stochastic kinetic systems. Here, one can formulate the model as a stochastic reaction-diffusion system and simulate the time course behaviour using an algorithm such as the *next sub-volume method* (Elf & Ehrenberg 2004), as implemented in the software MesoRD (Hattne, Fange & Elf 2005); see also Ander et al. (2004) and Lemerle, Di Ventura & Serrano (2005) for a related approach. Such algorithms will not be particularly efficient if the number of molecules being considered is significantly smaller than the number of subvolumes. In this case, a single-particle-based method is likely to be more appropriate, such as (approximately) implemented by StochSim (Le Novère & Shimizu 2001). Note that the algorithm used by StochSim is not particularly attractive in the context of regular systems biology models such as those considered in Chapter 7, as it will typically be more computationally intensive than the Gillespie algorithm as well as inexact. However, algorithms that represent each molecule as a

single software object can be attractive in the context of modelling complex chemical species where the potential number of distinct molecular species is vast, yet only a small number are expected to be present during the course of a given simulation run.

8.5 Exercises

1. Read the original source papers for the various different algorithms discussed in this chapter, and work through some simple examples by hand on paper to see how they really work.
2. Code up some of the algorithms and compare them for accuracy on some of the example models from Chapter 7. This can be done by comparing sample means, standard deviations, and full empirical probability distributions at a collection of time points, using data from many independent simulation runs; see Salis & Kaznessis (2005a) for further details.
3. After tuning the different algorithms so that they have reasonable accuracy, compare them for speed. Note that the algorithms discussed in this chapter are designed to be efficient for large models with large numbers of species, reactions, and molecules, so it should not be too surprising if the algorithms turn out to be slower than Gillespie's direct method on small models.
4. Download some software packages that implement stochastic simulation using algorithms other than the direct method, and compare them in terms of flexibility, accuracy, and performance. Links to some appropriate software are given on this book's website.

8.6 Further reading

The main source papers for this chapter are Gibson & Bruck (2000), Gillespie (2001), Gillespie & Petzold (2003), Gillespie (2000), Kiehl et al. (2004), Puchalka & Kierzek (2004), Rao & Arkin (2003), Haseltine & Rawlings (2002), and Salis & Kaznessis (2005a). All are required reading for a complete understanding of how the various algorithms work. Also see Wilkinson (2006) for a discussion of the issues involved in coupling Markov jump and Langevin processes.

Bayesian inference and MCMC

9.1 Likelihood and Bayesian inference

The final part of this book turns attention away from probabilistic modelling and stochastic simulation tools appropriate for systems biology models and toward statistical inference for the parameters of such models on the basis of experimental data. Although most of the earlier parts of the book are intended to be largely self-contained, these last chapters are likely to be more accessible to readers with a basic familiarity with ideas of estimation and statistical inference, including the concepts of likelihood, Bayesian analysis, and multi-variate statistics. A standard text on mathematical statistics such as Rice (1993) should provide most of the necessary background. However, any reader who has found the rest of the book straightforward should not have too much difficulty with the final chapters.

Chapter 10 examines the particular problem of inference for stochastic kinetic models using time-course experimental data. In this chapter we develop the methods of Bayesian inference and Markov chain Monte Carlo (MCMC) that will be required in Chapter 10. We will begin by examining the essential concepts of Bayesian inference and the reasons why analytic approaches to Bayesian inference in complex models are generally intractable, before moving on to MCMC and its application to Bayesian inference.

9.1.1 Bayesian inference

Often we are able to understand the probability of some *outcome*, X conditional on various possible *hypotheses*, H_i , $i = 1, 2, \dots, n$, where the H_i form a partition (Definition 3.6). We can then compute probabilities of the form $P(X = x' | H_i)$, $i = 1, \dots, n$, $x' \in S_X$. However, when we actually *observe* some outcome $X = x$, we are interested in the probabilities of the hypotheses *conditional* on the outcome, $P(H_i | X = x)$. Bayes Theorem (Corollary 3.1) tells us how to compute these, but the answer also depends on the prior probabilities for the hypotheses, $P(H_i)$, and hence to use Bayes Theorem, these too must be specified. Thus Bayes Theorem provides us with a coherent way of updating our prior beliefs about the hypotheses $P(H_i)$ to $P(H_i | X = x)$, our posterior beliefs based on the occurrence of $X = x$, as

$$P(H_i | X = x) = \frac{P(X = x | H_i) P(H_i)}{\sum_{j=1}^n P(X = x | H_j) P(H_j)}, \quad i = 1, \dots, n.$$

Note that the probabilities $P(X = x|H_i)$ are known as *likelihoods*, and are often written $L(H_i; x)$, as they tend to be regarded as a function of the H_i for given fixed outcome x . Note, however, that the *likelihood function* does not represent a PMF for the H_i ; in particular, there is no reason to suppose that it will sum to 1.

This is how it all works for purely discrete problems, but some adaptation is required before it can be used with continuous or mixed problems. Let us first stay with discrete outcome X and consider a continuum of hypotheses represented by a continuous parameter Θ . Our prior beliefs must now be represented by a density function, traditionally written, $\pi(\theta)$. Taking the continuous limit in the usual way, Bayes Theorem becomes

$$\pi(\theta|X = x) = \frac{\pi(\theta) P(X = x|\theta)}{\int_{\Theta} P(X = x|\theta') \pi(\theta') d\theta'}$$

In this case the likelihood function is $L(\theta; x) = P(X = x|\theta)$, regarded as a function of θ for given fixed x . Again, note that the likelihood function is not a density for θ , as it does not integrate to 1. Using this notation we can rewrite Bayes Theorem as

$$\pi(\theta|X = x) = \frac{\pi(\theta)L(\theta; x)}{\int_{\Theta} \pi(\theta')L(\theta'; x)d\theta'}$$

and this is the way it is usually written in the context of Bayesian statistics, though the likelihood function $L(\theta; x)$ means slightly different things depending on the context. Note that the integral on the bottom line of Bayes Theorem is not a function of θ , and so simply represents a constant of proportionality. Thus, we can rewrite Bayes Theorem in the simpler form

$$\pi(\theta|X = x) \propto \pi(\theta)L(\theta; x), \quad (9.1)$$

giving rise to the Bayesian mantra “*the posterior is proportional to the prior times the likelihood.*”

Example

Suppose that for a particular gene in a particular cell, transcription events occur according to a Poisson process with rate θ per minute. Prior to carrying out an experiment, a biological expert specifies his opinion regarding θ in the form of a $\Gamma(a, b)$ distribution (Section 3.11). Suppose that for our expert, $a = 2$, $b = 1$. Counts of the number of transcript events are gathered from n separate one-minute intervals to get

data $x = (x_1, x_2, \dots, x_n)'$. In this case the likelihood for θ is

$$\begin{aligned} L(\theta; x) &= P(x|\theta) \\ &= \prod_{i=1}^n P(x_i|\theta) \\ &= \prod_{i=1}^n \frac{\theta^{x_i} e^{-\theta}}{x_i!} \\ &\propto \prod_{i=1}^n \theta^{x_i} e^{-\theta} \\ &= \theta^{\sum_{i=1}^n x_i} e^{-n\theta}. \end{aligned}$$

The second line follows from the first because the data are independent (given θ). The likelihood depends on the data only through n and \bar{x} , so n and \bar{x} are said to be *sufficient statistics* for the likelihood function. Then since θ is gamma, we have

$$\pi(\theta) \propto \theta^{a-1} e^{-b\theta}$$

giving

$$\begin{aligned} \pi(\theta|x) &\propto \pi(\theta)L(\theta; x) \\ &\propto \theta^{a+\sum_{i=1}^n x_i-1} e^{-(b+n)\theta}. \end{aligned}$$

In other words,

$$\theta|x \sim \Gamma\left(a + \sum_{i=1}^n x_i, b + n\right).$$

So in this case, starting with a gamma prior results in a gamma posterior. Problems of this nature are said to be *conjugate*, and so in this case the gamma prior is said to be conjugate for the Poisson likelihood. In the context of our example, observing the data $x = (4, 2, 3)$ leads to a $\Gamma(11, 4)$ posterior distribution. This distribution (which has an expectation of $11/4$ and a variance of $11/16$) represents our belief about the value of θ having observed the data, and is shown in Figure 9.1.

The analysis is essentially the same when X is continuous rather than discrete, except that here the likelihood is evaluated using the PDF rather than the PMF.

9.1.2 Bayesian computation

In principle, the previous section covers everything we need to know about Bayesian inference — the posterior is nothing more (or less) than a conditional distribution for the parameters given the data. In practice, however, this may not be entirely trivial to work with.

The first problem one encounters is choosing the constant of proportionality so that the density integrates to 1. If the density is non-standard (as is usually the case for non-trivial problems), then the problem reduces to integrating the product of the likelihood and the prior (known as the *kernel* of the posterior) over the support of Θ .

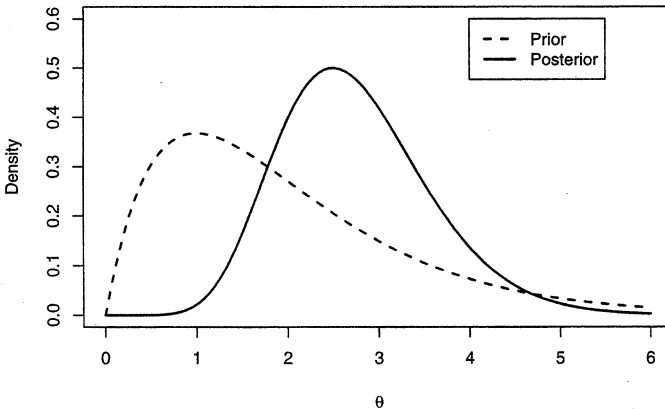


Figure 9.1 Plot showing the prior and posterior for the Poisson rate example. Note how the prior is modified to give a posterior more consistent with the data (which has a sample mean of 3).

If the support is infinite in extent, and/or multi-dimensional, then this is a highly non-trivial numerical problem.

Even if we have the constant of integration, if the parameter space is multi-dimensional, we will want to know what the marginal distribution of each component looks like. For each component, we have a very difficult numerical integration problem.

Example

Consider the case where we have a collection of observations, X_i , which we believe to be independent identically distributed (iid) normal with unknown mean and precision (the reciprocal of variance). We write

$$X_i | \mu, \tau \sim N(\mu, 1/\tau).$$

The likelihood for a single observation is

$$L(\mu, \tau; x_i) = f(x_i | \mu, \tau) = \sqrt{\frac{\tau}{2\pi}} \exp \left\{ -\frac{\tau}{2} (x_i - \mu)^2 \right\}$$

and so for n independent observations, $x = (x_1, \dots, x_n)'$ is

$$\begin{aligned} L(\mu, \tau; x) &= f(x|\mu, \tau) = \prod_{i=1}^n \sqrt{\frac{\tau}{2\pi}} \exp\left\{-\frac{\tau}{2}(x_i - \mu)^2\right\} \\ &= \left(\frac{\tau}{2\pi}\right)^{\frac{n}{2}} \exp\left\{-\frac{\tau}{2}[(n-1)s^2 + n(\bar{x} - \mu)^2]\right\} \\ &\propto \tau^{\frac{n}{2}} \exp\left\{-\frac{\tau}{2}[(n-1)s^2 + n(\bar{x} - \mu)^2]\right\} \end{aligned}$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

For a Bayesian analysis, we need also to specify prior distributions for the parameters, μ and τ . There is a conjugate analysis for this problem based on the specifications

$$\tau \sim \Gamma(a, b), \quad \mu|\tau \sim N\left(c, \frac{1}{d\tau}\right).$$

However, this specification is rather unsatisfactory — μ and τ are not independent, and in many cases our prior beliefs for μ and τ will separate into independent specifications. For example, we may prefer to specify independent priors for the parameters:

$$\tau \sim \Gamma(a, b), \quad \mu \sim N\left(c, \frac{1}{d}\right).$$

However, this specification is no longer conjugate, making analytic analysis intractable. Let us see why. We have

$$\pi(\mu) = \sqrt{\frac{d}{2\pi}} \exp\left\{-\frac{d}{2}(\mu - c)^2\right\} \propto \exp\left\{-\frac{d}{2}(\mu - c)^2\right\}$$

and

$$\pi(\tau) = \frac{b^a}{\Gamma(a)} \tau^{a-1} \exp\{-b\tau\} \propto \tau^{a-1} \exp\{-b\tau\},$$

so

$$\pi(\mu, \tau) \propto \tau^{a-1} \exp\left\{-\frac{d}{2}(\mu - c)^2 - b\tau\right\},$$

giving

$$\begin{aligned} \pi(\mu, \tau|x) &\propto \tau^{a-1} \exp\left\{-\frac{d}{2}(\mu - c)^2 - b\tau\right\} \times \tau^{\frac{n}{2}} \exp\left\{-\frac{\tau}{2}[(n-1)s^2 \right. \\ &\quad \left. + n(\bar{x} - \mu)^2]\right\} \\ &= \tau^{a+\frac{n}{2}-1} \exp\left\{-\frac{\tau}{2}[(n-1)s^2 + n(\bar{x} - \mu)^2] - \frac{d}{2}(\mu - c)^2 - b\tau\right\}. \end{aligned}$$

The posterior density for μ and τ certainly will not factorise (μ and τ are not independent *a posteriori*), and will not even separate into the form of the conditional normal-gamma conjugate form mentioned earlier.

So, we have the kernel of the posterior for μ and τ , but it is not in a standard form. We can gain some idea of the likely values of (μ, τ) by plotting the bivariate surface (the integration constant is not necessary for that), but we cannot work out the posterior mean or variance, or the forms of the marginal posterior distributions for μ or τ , since we cannot integrate out the other variable. We need a way of understanding posterior densities which does not rely on being able to analytically integrate the posterior density.

In fact, there is nothing particularly special about the fact that the density represents a Bayesian posterior. Given any complex non-standard multi-variate probability distribution, we need ways to understand it, to calculate its moments, and to compute its conditional and marginal distributions and their moments. Markov chain Monte Carlo (MCMC) algorithms such as the Gibbs sampler and the Metropolis-Hastings method provide a possible solution.

9.2 The Gibbs sampler

9.2.1 Introduction

The Gibbs sampler is a way of simulating from multi-variate distributions based only on the ability to simulate from conditional distributions. In particular, it is appropriate when sampling from marginal distributions is not convenient or possible.

Example

Reconsider the problem of Bayesian inference for the mean and variance of a normally distributed random sample. In particular, consider the non-conjugate approach based on independent prior distributions for the mean and variance. The posterior took the form

$$\pi(\mu, \tau | x) \propto \tau^{a + \frac{n}{2} - 1} \exp \left\{ -\frac{\tau}{2} [(n-1)s^2 + n(\bar{x} - \mu)^2] - \frac{d}{2} (\mu - c)^2 - b\tau \right\}.$$

As explained previously, this distribution is not in a standard form. However, while clearly not conjugate, this problem is often referred to as *semi-conjugate*, because the two *full conditional* distributions $\pi(\mu | \tau, x)$ and $\pi(\tau | \mu, x)$ are of standard form, and further, are of the same form as the independent prior specifications. That is, $\tau | \mu, x$ is gamma distributed and $\mu | \tau, x$ is normally distributed. In fact, by picking out terms in the variable of interest and regarding everything else as a constant of proportionality, we get

$$\begin{aligned} \tau | \mu, x &\sim \Gamma \left(a + \frac{n}{2}, b + \frac{1}{2} [(n-1)s^2 + n(\bar{x} - \mu)^2] \right), \\ \mu | \tau, x &\sim N \left(\frac{cd + n\tau\bar{x}}{n\tau + d}, \frac{1}{n\tau + d} \right). \end{aligned}$$

3. Change counter j to $j + 1$, and return to step 2.

This clearly defines a homogeneous Markov chain, as each simulated value depends only on the previous simulated value, and not on any other previous values or the iteration counter j . However, we need to show that $\pi(\theta)$ is a stationary distribution of this chain. The transition kernel of the chain is

$$p(\theta, \phi) = \prod_{i=1}^d \pi(\phi_i | \phi_1, \dots, \phi_{i-1}, \theta_{i+1}, \dots, \theta_d).$$

Therefore, we just need to check that $\pi(\theta)$ is the stationary distribution of the chain with this transition kernel. Unfortunately, the traditional *fixed-sweep* Gibbs sampler just described is *not* reversible, and so we cannot check stationarity by checking for detailed balance (as detailed balance fails). We need to do a direct check of the stationarity of $\pi(\theta)$, that is, we need to check that

$$\pi(\phi) = \int_S p(\theta, \phi) \pi(\theta) d\theta.$$

See Section 5.3 for a recap of the relevant concepts. For the bivariate case, we have

$$\begin{aligned} \int_S p(\theta, \phi) \pi(\theta) d\theta &= \int_S \pi(\phi_1 | \theta_2) \pi(\phi_2 | \phi_1) \pi(\theta_1, \theta_2) d\theta_1 d\theta_2 \\ &= \pi(\phi_2 | \phi_1) \int_{S_1} \int_{S_2} \pi(\phi_1 | \theta_2) \pi(\theta_1, \theta_2) d\theta_1 d\theta_2 \\ &= \pi(\phi_2 | \phi_1) \int_{S_2} \pi(\phi_1 | \theta_2) d\theta_2 \int_{S_1} \pi(\theta_1, \theta_2) d\theta_1 \\ &= \pi(\phi_2 | \phi_1) \int_{S_2} \pi(\phi_1 | \theta_2) \pi(\theta_2) d\theta_2 \\ &= \pi(\phi_2 | \phi_1) \pi(\phi_1) \\ &= \pi(\phi_1, \phi_2) \\ &= \pi(\phi). \end{aligned}$$

The general case is similar. So, $\pi(\theta)$ is a stationary distribution of this chain. Discussions of uniqueness and convergence are beyond the scope of this book. In particular, these issues are complicated somewhat by the fact that the sampler described is not reversible.

9.2.4 Reversible Gibbs samplers

While the fixed-sweep Gibbs sampler itself is not reversible, each *component update* is, and hence there are many variations on the fixed-sweep Gibbs sampler which are reversible and do satisfy detailed balance. Let us start by looking at why each component update is reversible.

Suppose we wish to update component i , that is, we update θ by replacing θ_i with ϕ_i drawn from $\pi(\phi_i | \theta_{-i})$. All other components will remain unchanged. The

transition kernel for this update is

$$p(\theta, \phi) = \pi(\phi_i | \theta_{-i}) I(\theta_{-i} = \phi_{-i})$$

where

$$I(E) = \begin{cases} 1 & \text{if } E \text{ is true,} \\ 0 & \text{if } E \text{ is false.} \end{cases}$$

Note that the density is zero for any transition changing the other components. Now we may check for detailed balance:

$$\begin{aligned} \pi(\theta) p(\theta, \phi) &= \pi(\theta) \pi(\phi_i | \theta_{-i}) I(\theta_{-i} = \phi_{-i}) \\ &= \pi(\theta_{-i}) \pi(\theta_i | \theta_{-i}) \pi(\phi_i | \theta_{-i}) I(\theta_{-i} = \phi_{-i}) \\ &= \pi(\phi_{-i}) \pi(\theta_i | \phi_{-i}) \pi(\phi_i | \phi_{-i}) I(\theta_{-i} = \phi_{-i}) \quad (\text{as } \theta_{-i} = \phi_{-i}) \\ &= \pi(\phi) \pi(\theta_i | \phi_{-i}) I(\theta_{-i} = \phi_{-i}) \\ &= \pi(\phi) p(\phi, \theta). \end{aligned}$$

Therefore, detailed balance is satisfied, and hence the update is reversible with stationary distribution $\pi(\cdot)$.

If this particular update is reversible and it preserves the equilibrium distribution of the chain, why bother updating any other component? The reason is that the chain defined by a single update is *reducible*, and hence will not converge to the stationary distribution from an arbitrary starting point. In order to ensure *irreducibility* of the chain, we need to make sure that we update each component sufficiently often. As we have seen, one way to do this is to update each component in a fixed order. The drawback of this method is that *reversibility* is lost.

An alternative to the *fixed-sweep* strategy is to pick a component at random at each stage and update that. This gives a reversible chain with the required stationary distribution, and is known as the *random scan* Gibbs sampler.

An even simpler way to restore the reversibility of the chain is to first scan through the components in fixed order, and then scan backward through the components. This does define a reversible Gibbs sampler. We can check that it works in the bivariate case as follows. The algorithm starts with (θ_1, θ_2) and then generates (ϕ_1, ϕ_2) as follows:

$$\begin{aligned} \theta'_1 &\sim \pi(\theta'_1 | \theta_2) \\ \phi_2 &\sim \pi(\phi_2 | \theta'_1) \\ \phi_1 &\sim \pi(\phi_1 | \phi_2). \end{aligned}$$

Here θ'_1 is an auxiliary variable that we are not interested in *per se* and which needs to be integrated out of the problem. The full transition kernel for a move from (θ_1, θ_2) to $(\theta'_1, \phi_1, \phi_2)$ is

$$p(\theta, (\theta'_1, \phi)) = \pi(\theta'_1 | \theta_2) \pi(\phi_2 | \theta'_1) \pi(\phi_1 | \phi_2),$$

and integrating out the auxiliary variable gives

$$\begin{aligned} p(\theta, \phi) &= \int \pi(\theta'_1 | \theta_2) \pi(\phi_2 | \theta'_1) \pi(\phi_1 | \phi_2) d\theta'_1 \\ &= \pi(\phi_1 | \phi_2) \int \pi(\theta'_1 | \theta_2) \pi(\phi_2 | \theta'_1) d\theta'_1. \end{aligned}$$

We can now check for detailed balance:

$$\begin{aligned} \pi(\theta) p(\theta, \phi) &= \pi(\theta) \pi(\phi_1 | \phi_2) \int \pi(\theta'_1 | \theta_2) \pi(\phi_2 | \theta'_1) d\theta'_1 \\ &= \pi(\theta_2) \pi(\theta_1 | \theta_2) \pi(\phi_1 | \phi_2) \int \pi(\theta'_1 | \theta_2) \pi(\phi_2 | \theta'_1) d\theta'_1 \\ &= \pi(\theta_1 | \theta_2) \pi(\phi_1 | \phi_2) \int \pi(\theta_2) \pi(\theta'_1 | \theta_2) \pi(\phi_2 | \theta'_1) d\theta'_1 \\ &= \pi(\theta_1 | \theta_2) \pi(\phi_1 | \phi_2) \int \pi(\theta'_1, \theta_2) \pi(\phi_2 | \theta'_1) d\theta'_1 \\ &= \pi(\theta_1 | \theta_2) \pi(\phi_1 | \phi_2) \int \pi(\theta'_1) \pi(\theta_2 | \theta'_1) \pi(\phi_2 | \theta'_1) d\theta'_1, \end{aligned}$$

and, as this is symmetric in θ and ϕ , we must have

$$\pi(\theta) p(\theta, \phi) = \pi(\phi) p(\phi, \theta).$$

This chain is therefore reversible with stationary distribution $\pi(\cdot)$.

We have seen that there are ways of adapting the standard fixed-sweep Gibbs sampler in ways which ensure reversibility. However, reversibility is not a requirement of a useful algorithm — it simply makes it easier to determine the properties of the chain. In practice, the fixed-sweep Gibbs sampler often has as good as or better convergence properties than its reversible cousins. Given that it is slightly easier to implement and debug, it is often simpler to stick with the fixed-sweep scheme than to implement a more exotic version of the sampler.

9.2.5 Simulation and analysis

Suppose that we are interested in a multivariate distribution $\pi(\theta)$ (which may be a Bayesian posterior distribution), and that we are able to simulate from the full conditional distributions of $\pi(\theta)$. Simulation from $\pi(\theta)$ is possible by first initialising the sampler somewhere in the support of θ , and then running the Gibbs sampler. The resulting chain should be monitored for convergence, and the “burn-in” period should be discarded for analysis. After convergence, the simulated values are all from $\pi(\theta)$. In particular, the values for a particular component will be simulated values from the marginal distribution of that component. A histogram of these values will give an idea of the “shape” of the marginal distribution, and summary statistics such as the mean and variance will be approximations to the mean and variance of the marginal distribution. The accuracy of the estimates can be gauged using the techniques from the end of Chapter 3.

```

normgibbs <- function(N, n, a, b, cc, d, xbar, ssquared)
{
  mat = matrix(ncol = 2, nrow = N)
  mu = cc
  tau = a/b
  mat[1, ] = c(mu, tau)
  for (i in 2:N) {
    muprec = n*tau + d
    mumean = (d*cc + n*tau*xbar)/muprec
    mu = rnorm(1, mumean, sqrt(1/muprec))
    taub = b + 0.5*((n - 1)*ssquared + n*(xbar -
      mu)^2)
    tau = rgamma(1, a + n/2, taub)
    mat[i, ] = c(mu, tau)
  }
  mat
}

```

Figure 9.2. An R function to implement a Gibbs sampler for the simple normal random sample model. Example code for using this function is given in Figure 9.3.

Example

Returning to the case of the posterior distribution for the normal model with unknown mean and precision, we have already established that the full conditional distributions are

$$\tau | \mu, x \sim \Gamma \left(a + \frac{n}{2}, b + \frac{1}{2} [(n-1)s^2 + n(\bar{x} - \mu)^2] \right),$$

$$\mu | \tau, x \sim N \left(\frac{cd + n\tau\bar{x}}{n\tau + d}, \frac{1}{n\tau + d} \right).$$

We can initialise the sampler anywhere in the half-plane where the posterior (and prior) has support, but convergence will be quicker if the chain is not started in the tails of the distribution. One possibility is to start the sampler near the posterior mode, though this can make convergence more difficult to diagnose. A simple strategy which is often easy to implement for problems in Bayesian inference is to start off the sampler at a point simulated from the prior distribution, or even at the mean of the prior distribution. Here, the prior mean for (τ, μ) is $(a/b, c)$. Once initialised, the sampler proceeds with alternate simulations from the full conditional distributions. The first few (hundred?) values should be discarded, and the rest can give information about the joint posterior distribution and marginals.

An R function to implement a simple Gibbs sampler is given in Figure 9.2, some example code that uses it is shown in Figure 9.3, and the results of running the example code are shown in Figure 9.4; see the figure legends for further details.

This example, though useful for illustrative purposes, is not a particularly compelling motivation for the deployment of MCMC methodology, due to the small


```

postmat=normgibbs (N=11000, n=15, a=3, b=11, cc=10, d=1/100, xbar=25,
                  ssquared=20)

postmat=postmat [1001:11000,]
op=par (mfrow=c (3, 3))
plot (postmat, col=1:10000)
plot (postmat, type="l")
plot.new ()
plot (ts (postmat [, 1]))
plot (ts (postmat [, 2]))
plot (ts (sqrt (1/postmat [, 2])))
hist (postmat [, 1], 40)
hist (postmat [, 2], 40)
hist (sqrt (1/postmat [, 2]), 40)
par (op)

```

Figure 9.3 Example R code illustrating the use of the function `normgibbs` from Figure 9.2. The plots generated by running this code are shown in Figure 9.4. In this example the prior took the form $\mu \sim N(10, 100)$, $\tau \sim \Gamma(3, 11)$, and the sufficient statistics for the data were $n = 15$, $\bar{x} = 25$, $s^2 = 20$. The sampler was run for 11,000 iterations with the first 1,000 discarded as burn-in, and the remaining 10,000 iterations used for the main monitoring run.

number of dimensions and relatively simple structure. Before leaving the topic of Gibbs sampling, a slightly more substantial example will be examined.

Example

This example will be motivated by considering a biological experiment to estimate (the logarithm of) a biochemical reaction rate constant. The details of the experiment will not concern us here; we will assume simply that the experiment results in the generation of a single number, representing a (sensible) estimate of the rate constant. Several labs will conduct this experiment, and each lab will replicate the experiment several times. Thus, in totality, we will have a collection of estimates from a collection of labs. We will consider the problem of inference for the true rate constant on the basis of prior knowledge and all available data.

Consider the following simple *hierarchical* (or *one-way random effects*) model,

$$\begin{aligned}
 Y_{ij} | \theta_i, \tau &\sim N(\theta_i, 1/\tau), & \text{independently, } i = 1, \dots, m, j = 1, \dots, n_i \\
 \theta_i | \mu, \nu &\sim N(\mu, 1/\nu), & i = 1, \dots, m.
 \end{aligned}$$

Here there are m labs, and the i th lab replicates the experiment n_i times. Y_{ij} is the measurement made on the j th experiment by lab i . We assume that the measurements from lab i have mean θ_i and that the measurements are normally distributed. We also assume that the θ_i are themselves normally distributed around the true rate μ . Essentially, the model has the effect of inducing a correlation between replicates from a particular lab, due to the fact that we expect replicates from one lab to be more similar than replicates from different labs (due to hidden factors that are not being properly controlled and accounted for). Note that this generic scenario can be

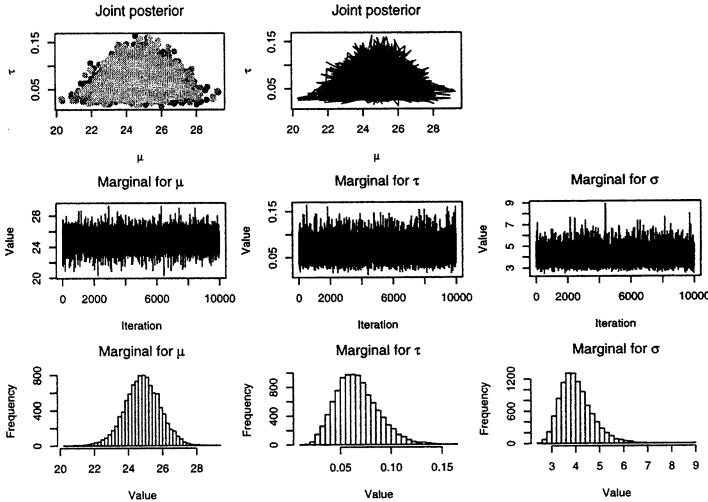


Figure 9.4 Figure showing the Gibbs sampler output resulting from running the example code in Figure 9.3. The top two plots give an indication of the bivariate posterior distribution. The second row shows trace plots of the marginal distributions of interest, indicating a rapidly mixing MCMC algorithm. The final row shows empirical marginal posterior distributions for the parameters of interest.

applied to a range of different scenarios and is particularly popular in the context of meta-analysis.

We will consider the most general (and quite typical) case where μ , τ and ν are all unknown. We wish to make inferences about these parameters in addition to the unknown θ_i . Thus there are $m + 3$ parameters of interest in this model.

The specification of the model is completed with independent priors for μ , τ , and ν ,

$$\begin{aligned} \mu &\sim N(a, 1/b) \\ \tau &\sim \Gamma(c, d) \\ \nu &\sim \Gamma(e, f). \end{aligned}$$

In principle we have now completely specified the model and can compute the posterior distribution. Of course, the posterior distribution is very high dimensional and, more importantly, not of a standard form. Here MCMC techniques can be used to study the posterior distribution.* The likelihood contribution for each observation y_{ij}

* It turns out that it is possible to use the linear Gaussian structure of the model to integrate out all variables except the two variance components, reducing the problem to a two-dimensional one; see Wilkinson & Yeung (2002) and Wilkinson & Yeung (2004) for details. However, for models without this linear Gaussian structure, such techniques are not always convenient or possible, and simple MCMC algorithms like the one about to be described are often the only straightforward way to work with the posterior distribution.

is

$$L(\theta_i, \tau; y_{ij}) = \sqrt{\frac{\tau}{2\pi}} \exp \left\{ -\frac{\tau}{2} (y_{ij} - \theta_i)^2 \right\}$$

and so the full likelihood is

$$\begin{aligned} L(\theta, \tau; y) &= \prod_{i=1}^m \prod_{j=1}^{n_i} L(\theta_i, \tau; y_{ij}) \\ &= \left(\frac{\tau}{2\pi} \right)^{N/2} \exp \left\{ -\frac{\tau}{2} \sum_{i=1}^m [(n_i - 1)s_i^2 + n_i(y_{i\cdot} - \theta_i)^2] \right\} \end{aligned}$$

where

$$N = \sum_{i=1}^m n_i, \quad y_{i\cdot} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}, \quad s_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (y_{ij} - y_{i\cdot})^2.$$

The prior takes the form

$$\pi(\mu, \tau, \nu, \theta) = \pi(\mu)\pi(\tau)\pi(\nu)\pi(\theta|\mu, \nu)$$

where

$$\begin{aligned} \pi(\mu) &\propto \exp \left\{ -\frac{b}{2} (\mu - a)^2 \right\} \\ \pi(\tau) &\propto \tau^{c-1} \exp\{-d\tau\} \\ \pi(\nu) &\propto \nu^{e-1} \exp\{-f\nu\} \\ \pi(\theta_i|\mu, \nu) &= \sqrt{\frac{\nu}{2\pi}} \exp \left\{ -\frac{\nu}{2} (\theta_i - \mu)^2 \right\} \\ \Rightarrow \pi(\theta|\mu, \nu) &\propto \nu^{m/2} \exp \left\{ -\frac{\nu}{2} \sum_{i=1}^m (\theta_i - \mu)^2 \right\} \end{aligned}$$

and therefore,

$$\pi(\mu, \tau, \nu, \theta) \propto \nu^{e+m/2-1} \tau^{c-1} \exp \left\{ -\frac{1}{2} \left[2d\tau + 2f\nu + b(\mu - a)^2 + \nu \sum_{i=1}^m (\theta_i - \mu)^2 \right] \right\}$$

Now we have the likelihood and the prior, and we can write down the posterior distribution,

$$\begin{aligned} \pi(\mu, \tau, \nu, \theta|y) &\propto \\ &\tau^{c+N/2-1} \nu^{e+m/2-1} \exp \left\{ -\frac{1}{2} \left[2d\tau + 2f\nu + b(\mu - a)^2 \right. \right. \\ &\quad \left. \left. + \sum_{i=1}^m (\nu(\theta_i - \mu)^2 + \tau(n_i - 1)s_i^2 + \tau n_i (y_{i\cdot} - \theta_i)^2) \right] \right\}. \end{aligned}$$

It is difficult to do anything analytic with this, so we will try to construct a Gibbs sampler in order to investigate it. This is fairly straightforward, and the full conditionals are as follows. Picking out terms in μ , the full conditional for μ is

$$\begin{aligned}\pi(\mu|\cdot) &\propto \exp \left\{ -\frac{1}{2} \left[b(\mu - a)^2 + \sum_{i=1}^m \nu(\theta_i - \mu)^2 \right] \right\} \\ &\propto \exp \left\{ -\frac{1}{2} (b + m\nu) \left(\mu - \frac{ba + m\nu\bar{\theta}}{b + m\nu} \right)^2 \right\}, \quad \text{where } \bar{\theta} = \frac{1}{m} \sum_{i=1}^m \theta_i,\end{aligned}$$

that is

$$\mu|\cdot \sim N \left(\frac{ba + m\nu\bar{\theta}}{b + m\nu}, \frac{1}{b + m\nu} \right).$$

The conditional for τ is

$$\pi(\tau|\cdot) \propto \tau^{c+N/2-1} \exp \left\{ -\tau \left[d + \frac{1}{2} \sum_{i=1}^m ((n_i - 1)s_i^2 + n_i(y_i - \theta_i)^2) \right] \right\},$$

that is

$$\tau|\cdot \sim \Gamma \left(c + N/2, d + \frac{1}{2} \sum_{i=1}^m [(n_i - 1)s_i^2 + n_i(y_i - \theta_i)^2] \right).$$

The conditional for ν is

$$\pi(\nu|\cdot) \propto \nu^{e+m/2-1} \exp \left\{ -\nu \left[f + \frac{1}{2} \sum_{i=1}^m (\theta_i - \mu)^2 \right] \right\},$$

that is

$$\nu|\cdot \sim \Gamma \left(e + m/2, f + \frac{1}{2} \sum_{i=1}^m (\theta_i - \mu)^2 \right).$$

The conditional for θ_i is

$$\begin{aligned}\pi(\theta_i|\cdot) &\propto \exp \left\{ -\frac{1}{2} [\nu(\theta_i - \mu)^2 + \tau n_i(y_i - \theta_i)^2] \right\} \\ &\propto \exp \left\{ -\frac{1}{2} (\nu + n_i\tau) \left(\theta_i - \frac{\nu\mu + n_i y_i \tau}{\nu + n_i\tau} \right)^2 \right\},\end{aligned}$$

that is

$$\theta_i|\cdot \sim N \left(\frac{\nu\mu + n_i y_i \tau}{\nu + n_i\tau}, \frac{1}{\nu + n_i\tau} \right), \quad i = 1, \dots, m.$$

Therefore, we have a Gibbs sampler with $m + 3$ components. We just have to specify

the prior parameters a, b, c, d, e, f and compute the data summaries m, n_i, N, y_i, s_i^2 , $i = 1, \dots, m$. Then we initialise the sampler by simulating from the prior, or by starting off each component at its prior mean. The sampler is then run to convergence, and samples from the stationary distribution are used to understand the marginals of the posterior distribution. This model is of sufficient complexity that assessing convergence of the sampler to its stationary distribution is a non-trivial task. At the very least, multiple large simulation runs are required, with different starting points, and the first portion (say, a third) of any run should be discarded as “burn-in.” As the complexity of the model increases, problems with assessment of the convergence of the sampler also increase. There are many software tools available for MCMC convergence diagnostics. R-CODA is an excellent package for R which carries out a range of output analysis and diagnostic tasks, but its use is beyond the scope of this book.

It is clear that in principle at least, it ought to be possible to automate the construction of a Gibbs sampler from a specification containing the model, the prior, and the data. There are several freely available software packages that are able to do this for relatively simple models. Examples include, WinBUGS, OpenBugs, and JAGS; see this book’s website for links. Unfortunately, it turns out to be difficult to use these software packages for the stochastic kinetic models that will be considered in Chapter 10.

Of course, the Gibbs sampler tacitly assumes that we have some reasonably efficient mechanism for simulating from the full conditional distributions, and yet this is not always the case. Fortunately, the Gibbs sampler can be combined with Metropolis-Hastings algorithms when the full conditionals are difficult to simulate from.

9.3 The Metropolis-Hastings algorithm

Suppose that $\pi(\theta)$ is the density of interest. Suppose further that we have some (arbitrary) transition kernel $q(\theta, \phi)$ (known as the *proposal distribution*) which is easy to simulate from, but does not (necessarily) have $\pi(\theta)$ as its stationary density. Consider the following algorithm:

1. Initialise the iteration counter to $j = 1$, and initialise the chain to $\theta^{(0)}$.
2. Generate a *proposed* value ϕ using the kernel $q(\theta^{(j-1)}, \phi)$.
3. Evaluate the *acceptance probability* $\alpha(\theta^{(j-1)}, \phi)$ of the proposed move, where

$$\alpha(\theta, \phi) = \min \left\{ 1, \frac{\pi(\phi)q(\phi, \theta)}{\pi(\theta)q(\theta, \phi)} \right\}.$$

4. Put $\theta^{(j)} = \phi$ with probability $\alpha(\theta^{(j-1)}, \phi)$, and put $\theta^{(j)} = \theta^{(j-1)}$ otherwise.
5. Change the counter from j to $j + 1$ and return to step 2.

In other words, at each stage, a new value is generated from the proposal distribution. This is either accepted, in which case the chain moves, or rejected, in which case the chain stays where it is. Whether or not the move is accepted or rejected depends on an acceptance probability which itself depends on the relationship between the density

of interest and the proposal distribution. Note that the density of interest, $\pi(\cdot)$, only enters into the acceptance probability as a ratio, and so the method can be used when the density of interest is only known up to a scaling constant. This algorithm is essentially that of Hastings (1970), which is a generalisation of the algorithm introduced by Metropolis et al. (1953).

The Markov chain defined in this way is reversible and has stationary distribution $\pi(\cdot)$ irrespective of the choice of proposal distribution, $q(\cdot, \cdot)$. Let us see why. The transition kernel is clearly given by

$$p(\theta, \phi) = q(\theta, \phi)\alpha(\theta, \phi), \quad \text{if } \theta \neq \phi.$$

But there is also a finite probability that the chain will remain at θ . This is 1 minus the probability that the chain moves, and thus is given by

$$1 - \int q(\theta, \phi)\alpha(\theta, \phi) d\phi.$$

So, the transition kernel is part continuous and part discrete. We can easily write down the cumulative distribution form of the transition kernel as

$$P(\theta, \phi) = \int_{-\infty}^{\phi} q(\theta, \phi)\alpha(\theta, \phi) d\phi + I(\phi \geq \theta) \left[1 - \int q(\theta, \phi)\alpha(\theta, \phi) d\phi \right].$$

We then get the full density form of the kernel by differentiating with respect to ϕ as

$$p(\theta, \phi) = q(\theta, \phi)\alpha(\theta, \phi) + \delta(\theta - \phi) \left[1 - \int q(\theta, \phi)\alpha(\theta, \phi) d\phi \right],$$

where $\delta(\cdot)$ is the Dirac δ -function. Now we have the transition kernel we can check whether detailed balance is satisfied:

$$\begin{aligned} \pi(\theta)p(\theta, \phi) &= \pi(\theta)q(\theta, \phi) \min \left\{ 1, \frac{\pi(\phi)q(\phi, \theta)}{\pi(\theta)q(\theta, \phi)} \right\} \\ &\quad + \delta(\theta - \phi) \left[\pi(\theta) - \int \pi(\theta)q(\theta, \phi) \min \left\{ 1, \frac{\pi(\phi)q(\phi, \theta)}{\pi(\theta)q(\theta, \phi)} \right\} d\phi \right] \\ &= \min \{ \pi(\theta)q(\theta, \phi), \pi(\phi)q(\phi, \theta) \} \\ &\quad + \delta(\theta - \phi) \left[\pi(\theta) - \int \min \{ \pi(\theta)q(\theta, \phi), \pi(\phi)q(\phi, \theta) \} d\phi \right]. \end{aligned}$$

The first term is clearly symmetric in θ and ϕ . Also, the second term must be symmetric in θ and ϕ , because it is only non-zero precisely when $\theta = \phi$. Consequently, detailed balance is satisfied, and the Metropolis-Hastings algorithm defines a reversible Markov chain with stationary distribution $\pi(\cdot)$, irrespective of the form of $q(\cdot, \cdot)$.

Complete freedom in the choice of the proposal distribution $q(\cdot, \cdot)$ leaves us wondering what kinds of choices might be good, or generally quite useful. Some commonly used special cases are discussed below.

9.3.1 Symmetric chains (Metropolis method)

The simplest case is the Metropolis sampler, which is based on the use of a symmetric proposal with $q(\theta, \phi) = q(\phi, \theta)$, $\forall \theta, \phi$. We see then that the acceptance probability simplifies to

$$\alpha(\theta, \phi) = \min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\},$$

and hence does not involve the proposal density at all. Consequently proposed moves which will take the chain to a region of higher density are always accepted, while moves which take the chain to a region of lower density are accepted with probability proportional to the ratio of the two densities — moves which will take the chain to a region of very low density will be accepted with very low probability. Note that any proposal of the form $q(\theta, \phi) = f(|\theta - \phi|)$ is symmetric, where $f(\cdot)$ is an arbitrary density. In this case, the proposal will represent a symmetric displacement from the current value. This also motivates random walk chains.

9.3.2 Random walk chains

In this case, the proposed value ϕ at stage j is $\phi = \theta^{(j-1)} + w_j$ where the w_j are iid random variables (completely independent of the state of the chain). Suppose that the w_j have density $f(\cdot)$, which is easy to simulate from. We can then simulate an *innovation*, w_j , and set the *candidate* point to $\phi = \theta^{(j-1)} + w_j$. The transition kernel is then $q(\theta, \phi) = f(\phi - \theta)$, and this can be used to compute the acceptance probability. Of course, if $f(\cdot)$ is symmetric about zero, then we have a symmetric chain, and the acceptance probability does not depend on $f(\cdot)$ at all.

So, suppose that it is decided to use a symmetric random walk chain with proposed mean zero innovations. There is still the question of how they should be distributed, and what variance they should have. A simple, easy to simulate from distribution is always a good idea, such as uniform or normal (normal is generally better, but is a bit more expensive to simulate). The choice of variance will affect the acceptance probability, and hence the overall proportion of accepted moves. If the variance of the innovation is too low, then most proposed values will be accepted, but the chain will move very slowly around the space — the chain is said to be too “cold.” On the other hand, if the variance of the innovation is too large, very few proposed values will be accepted, but when they are, they will often correspond to quite large moves — the chain is said to be too “hot.” Experience suggests that an overall acceptance rate of around 30% is desirable, and so it is possible to “tune” the variance of the innovation distribution to get an acceptance rate of around this level. This should be done using a few trial short runs, and then a single fixed value should be adopted for the main monitoring run.[†]

An R function implementing a simple Metropolis random walk sampler is given

[†] Although it sounds appealing to adaptively change the tuning parameter during the main monitoring run, this usually affects the stationary distribution of the chain, and hence should be avoided (unless you *really* know what you are doing).

```

metrop <- function(n, alpha)
{
  vec=vector("numeric", n)
  x=0
  vec[1]=x
  for (i in 2:n) {
    can=x+runif(1, -alpha, alpha)
    apro=dnorm(can) / dnorm(x)
    u=runif(1)
    if (u < apro)
      x=can
    vec[i]=x
  }
  vec
}

```

Figure 9.5 An R function to implement a Metropolis sampler for a standard normal random quantity based on $U(-\alpha, \alpha)$ innovations. So, `metrop(10000, 1)` will execute a run of length 10,000 with an α of 1. This α is close to optimal. Running with $\alpha = 0.1$ gives a chain that is too cold, and $\alpha = 100$ gives a chain that is too hot.

in Figure 9.5. In this example the target distribution is a $N(0, 1)$ random quantity, and the innovations are $U(-\alpha, \alpha)$.

9.3.3 Independence chains

In this case, reminiscent of the envelope rejection method (Section 4.6.1) and importance sampling (Section 4.2), the proposed transition is formed independent of the previous position of the chain, and so $q(\theta, \phi) = f(\phi)$ for some density $f(\cdot)$. Here the acceptance probability becomes

$$\alpha(\theta, \phi) = \min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \frac{f(\phi)}{f(\theta)} \right\},$$

and we see that the acceptance probability can be increased by making $f(\cdot)$ as similar to $\pi(\cdot)$ as possible (in this case, the higher the acceptance probability, the better).

Bayes Theorem via independence chains

In the context of Bayesian inference, one possible choice for the proposal density is the prior density. The acceptance probability then becomes

$$\alpha(\theta, \phi) = \min \left\{ 1, \frac{L(\phi; x)}{L(\theta; x)} \right\},$$

and hence depends only on the likelihood ratio of the candidate point and the current value. This is attractive because it is usually very straightforward to simulate from the prior distribution and very difficult to simulate from the posterior. Also, the acceptance probability depends only on the likelihood ratio, which is often quite

straightforward to calculate. Unfortunately this method tends to result in a badly mixing chain if the problem is high dimensional and the data are not in strong accordance with the prior.

9.4 Hybrid MCMC schemes

We have seen how we can use the Gibbs sampler to sample from multi-variate distributions provided that we can simulate from the full conditionals. We have also seen how we can use Metropolis-Hastings methods to sample from awkward distributions (perhaps full conditionals). If we wish, we can combine these in order to form hybrid Markov chains whose stationary distribution is a distribution of interest.

Componentwise transition: Given a multivariate distribution with full conditionals that are awkward to sample from directly, we can define a Metropolis-Hastings scheme for each full conditional and apply them to each component in turn for each iteration. This is like the Gibbs sampler, but each component update is a Metropolis-Hastings update, rather than a direct simulation from the full conditional. This is in fact the original form of the Metropolis algorithm.

Metropolis within Gibbs: Given a multivariate distribution with full conditionals, some of which may be simulated from directly, and others which have Metropolis-Hastings updating schemes, the Metropolis within Gibbs algorithm goes through each in turn, and simulates directly from the full conditional, or carries out a Metropolis-Hastings update as necessary.

Blocking: The components of a Gibbs sampler, and those of Metropolis-Hastings chains, can be vectors (or matrices) as well as scalars. For many high-dimensional problems, it can be helpful to group related parameters into blocks and use multivariate simulation techniques to update those together if possible. This can greatly improve the mixing of the chain, at the expense of increasing the computational cost of each iteration.

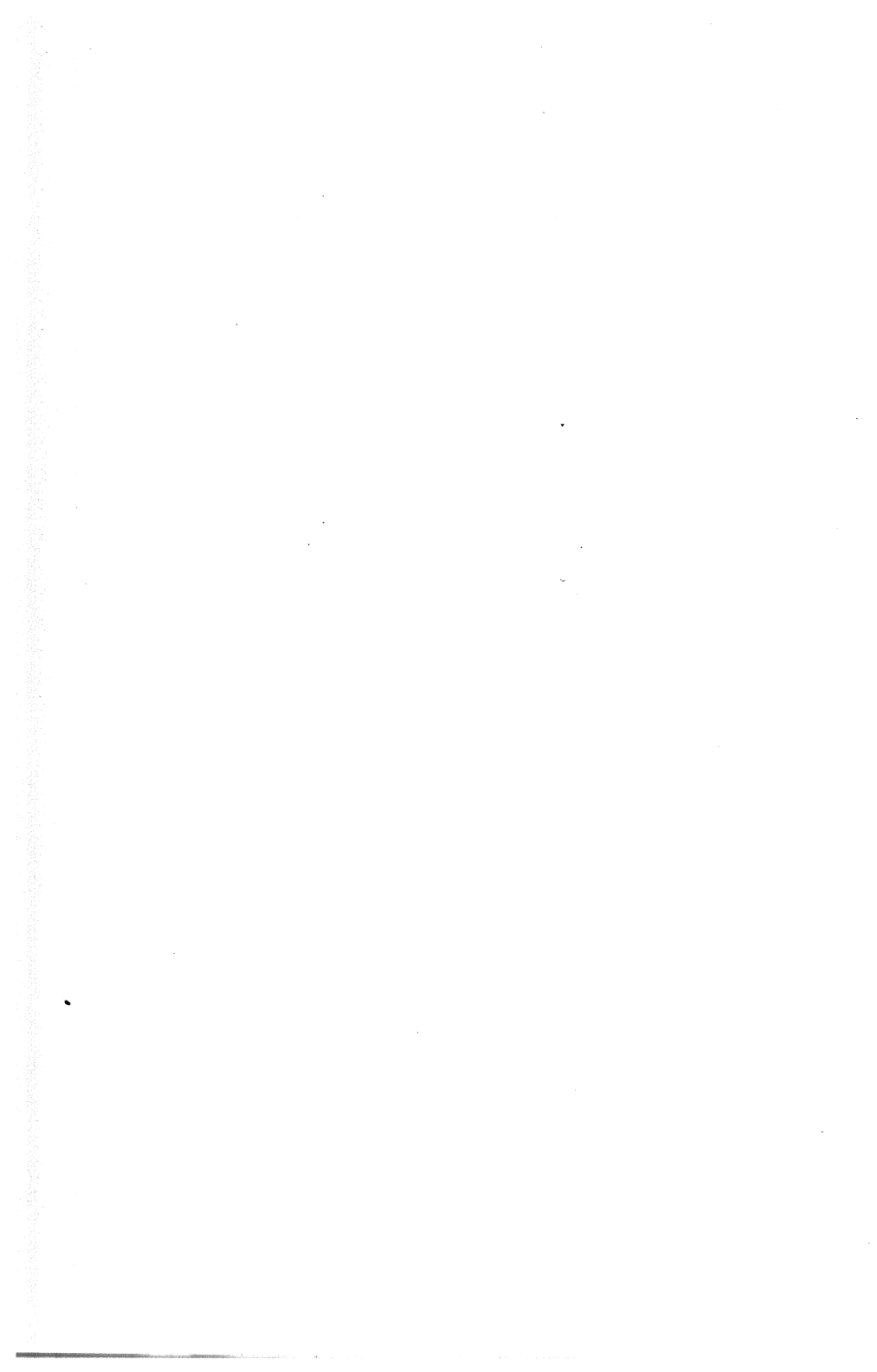
Some of the methods discussed in this section will be illustrated in practice in Chapter 10.

9.5 Exercises

1. Modify the simple Metropolis code given in Figure 9.5 in order to compute the overall acceptance rate of the chain. Write another R function which uses this modified function in order to automatically find a tuning parameter giving an overall acceptance rate of around 30%.
2. Rewrite the Gibbs sampling code from Figure 9.2 in a faster language such as C, Fortran, or Java. Real MCMC algorithms run too slowly in R, so it is necessary to build up an MCMC code base in a more efficient language.
3. Install the R-CODA package for MCMC output analysis and diagnostics and learn how it works. Try it out on the examples you have been studying.
4. Download some automatic MCMC software (linked from this book's website). Learn how these packages work and try them out on some simple models.

9.6 Further reading

For a more leisurely and more comprehensive introduction to Bayesian inference and MCMC, it is worth starting with O'Hagan & Forster (2004) for the background on Bayesian inference, and then moving on to Gamerman (1997) for further information regarding MCMC.



Inference for stochastic kinetic models

10.1 Introduction

This chapter provides an introduction to the computational statistical techniques that can be used in order to identify stochastic kinetic models from experimental data. It is worth mentioning at this point the experimental framework envisioned. Essentially, the techniques to be developed require high-quality, calibrated, high-resolution time-course measurements of levels of a reasonably large subset of model species at the single-cell level. At the time of writing such data are relatively difficult to obtain, so it is reasonable to question the practical utility of such methods to the average systems biologist. The reality is that such data are going to be *required* to identify stochastic kinetic models, and so experimental biology needs to change to make such data more readily available. In the meantime, however, the development of these inference techniques sheds light on the issues involved in identifying stochastic kinetic models, and can in any case be adapted to other (less perfect) data sources, including population-averaged data.* The methods are particularly useful for highlighting model identification issues. For example, it might be intended to identify a small network containing 15 unknown rate constants using data on three key protein species. The inference techniques here described might tell you that it is impossible to identify some of the rate constants even using *perfect* information on those three proteins. The methodology should also be able to generate a minimal set of species that need to be measured in order to identify all 15 rate constants satisfactorily and give some indication of the amount of data on these species that will need to be collected in order to have reliable estimates. Thus, even in the absence of high-quality single-cell data, the methods allow investigation of the vital but difficult experimental design issues that need to be addressed in order to get maximum benefit from costly wet-biology programmes. In any case, work on data acquisition at the single-cell level is progressing rapidly (Bahcall 2005), and so it seems likely that by the time this book is in print, good single-cell data may well be a realistically attainable target for well-funded systems biology centres.

This chapter can provide only a brief introduction to the ideas and techniques required. The aim of the chapter is not to make the reader an expert in Bayesian inference for stochastic kinetic models, but simply to make the (fairly technical)

* The problem with population-averaged data is that they essentially provide information about the mean of the stochastic process, and effectively mask all other information about the process behaviour. We saw in Chapter 1 an example of a stochastic process (the linear birth-death process) that cannot be identified by its mean. Because it involves only first-order reactions, the deterministic solution is the mean of the stochastic process (Section 6.7), and it was explained in Chapter 1 why it is not possible to use the deterministic model to identify both rate constants.

literature in this area more accessible. Appropriate further reading will be highlighted where relevant.

10.2 Inference given complete data

It turns out to be helpful to consider first the problem of inference given perfect data on the state of the model over a finite time interval $[0, T]$. That is, we will assume that the entire sample path of each species in the model is known over the time period $[0, T]$. This is equivalent to assuming that we have been given discrete-event output from a Gillespie simulator, and are then required to figure out the rate constants that were used on the basis of the output. Although it is completely unrealistic to assume that experimental data of this quality will be available in practice, understanding this problem is central to understanding the more general inference problem. In any case, it is clear that if we cannot solve even this problem, then inference from data sources of lower quality will be beyond our reach.

It will be helpful to assume the model notation from Chapter 6, with species x_1, \dots, x_u , reactions R_1, \dots, R_v , rate constants $c = (c_1, \dots, c_v)'$, reaction hazards $h_1(x, c_1), \dots, h_v(x, c_v)$, and combined hazard

$$h_0(x, c) = \sum_{i=1}^v h_i(x, c_i).$$

It is now necessary to explicitly consider the state of the system at a given time, and this will be denoted $x(t) = (x_1(t), \dots, x_u(t))'$. Our observed sample path will be written

$$x = \{x(t) : t \in [0, T]\}.$$

As we have complete information on the sample path, we also know the time and type of each reaction event (in fact, this is what we really mean by complete data). It is helpful to use the notation r_j for the number of reaction events of type R_j that occurred in the sample path x , $j = 1, \dots, v$, and to define $n = \sum_{j=1}^v r_j$ to be the total number of reaction events occurring in the interval $[0, T]$. We will now consider the time and type of each reaction event, (t_i, ν_i) , $i = 1, \dots, n$, where the t_i are assumed to be in increasing order and $\nu_i \in \{1, \dots, v\}$. It is notationally convenient to make the additional definitions $t_0 = 0$ and $t_{n+1} = T$.

In order to carry out model-based inference for the process, we need the likelihood function. A formal approach to the development of a rigorous theory of likelihood for continuous sample paths is beyond the scope of a text such as this, but it is straightforward to compute the likelihood in an informal way by considering the terms in the likelihood that arise from constructing the sample path according to Gillespie's direct method. Here, the term in the likelihood corresponding to the i th event is just the joint density of the time and type of that event. That is,

$$\begin{aligned} h_0(x(t_{i-1}), c) \exp\{-h_0(x(t_{i-1}), c)[t_i - t_{i-1}]\} &\times \frac{h_{\nu_i}(x(t_{i-1}), c_i)}{h_0(x(t_{i-1}), c)} \\ &= \exp\{-h_0(x(t_{i-1}), c)[t_i - t_{i-1}]\} h_{\nu_i}(x(t_{i-1}), c_i). \end{aligned}$$

The full likelihood is the product of these terms, together with a final term representing the information in the fact that there is no event in the final interval $(t_n, T]$. This is just given by the probability of that event, which is

$$\exp\{-h_0(x(t_n), c)[T - t_n]\},$$

giving the combined likelihood

$$\begin{aligned} L(c; \mathbf{x}) &= \pi(\mathbf{x}|c) \\ &= \left\{ \prod_{i=1}^n h_{\nu_i}(x(t_{i-1}), c_i) \exp\{-h_0(x(t_{i-1}), c)[t_i - t_{i-1}]\} \right\} \\ &\quad \times \exp\{-h_0(x(t_n), c)[T - t_n]\} \\ &= \left\{ \prod_{i=1}^n h_{\nu_i}(x(t_{i-1}), c_i) \right\} \left\{ \prod_{i=1}^{n+1} \exp\{-h_0(x(t_{i-1}), c)[t_i - t_{i-1}]\} \right\} \\ &= \left\{ \prod_{i=1}^n h_{\nu_i}(x(t_{i-1}), c_i) \right\} \exp \left\{ \sum_{i=1}^{n+1} -h_0(x(t_{i-1}), c)[t_i - t_{i-1}] \right\} \\ &= \left\{ \prod_{i=1}^n h_{\nu_i}(x(t_{i-1}), c_i) \right\} \exp \left\{ - \sum_{i=0}^n h_0(x(t_i), c)[t_{i+1} - t_i] \right\}. \quad (10.1) \end{aligned}$$

This expression for the likelihood (10.1) can be used directly for computing the likelihood from complete data and is therefore known as the *complete-data likelihood*. It is interesting to note that due to the piecewise constant nature of the combined hazard function, the sum in the exponential term is actually an integral, and so the complete-data likelihood can be written more neatly in the following way.

Theorem 10.1 *The complete-data likelihood for a stochastic kinetic model on the time interval $[0, T]$ takes the form*

$$L(c; \mathbf{x}) = \pi(\mathbf{x}|c) = \left\{ \prod_{i=1}^n h_{\nu_i}(x(t_{i-1}), c_i) \right\} \exp \left\{ - \int_0^T h_0(x(t), c) dt \right\}. \quad (10.2)$$

From this informal perspective, the occurrence of the integral of the combined hazard function in (10.2) seems slightly mysterious. In fact, when viewed from a more advanced perspective, developing the notion of likelihood through stochastic calculus and the Radon-Nikodym derivative, the form of (10.2) is seen to be completely intuitive and typical of all Markov jump processes; the first (product) term represents all of the information in the jumps, and the second (integral) term represents all of the information in the period of full measure where no jumps occur.

For completely general hazard functions, equations (10.1) and (10.2) represent the complete-data likelihood in its simplest form. However, in the case of the simple mass-action kinetic rate laws typically used in this context, it turns out that the likelihood factorises in a particularly convenient way. The key requirement is that the hazard functions can be written in the form $h_j(x, c_j) = c_j g_j(x)$, $j = 1, \dots, v$.

Substituting into (10.2) and simplifying then gives

$$\begin{aligned}
 L(\mathbf{c}; \mathbf{x}) &= \left\{ \prod_{i=1}^n c_{\nu_i} g_{\nu_i}(x(t_{i-1})) \right\} \exp \left\{ - \int_0^T \sum_{j=1}^v c_j g_j(x(t)) dt \right\} \\
 &\propto \left\{ \prod_{j=1}^v c_j^{r_j} \right\} \exp \left\{ - \sum_{j=1}^v \int_0^T c_j g_j(x(t)) dt \right\} \\
 &= \prod_{j=1}^v c_j^{r_j} \exp \left\{ -c_j \int_0^T g_j(x(t)) dt \right\} \\
 &= \prod_{j=1}^v L_j(c_j; \mathbf{x}),
 \end{aligned}$$

where the component likelihoods are defined by

$$L_j(c_j; \mathbf{x}) = c_j^{r_j} \exp \left\{ -c_j \int_0^T g_j(x(t)) dt \right\}, \quad j = 1, \dots, v. \quad (10.3)$$

This factorisation of the complete-data likelihood has numerous important consequences for inference. It means that in the complete data scenario, information regarding each rate constant is independent of the information regarding the other rate constants. That is, inference may be carried out for each rate constant separately. For example, in a maximum likelihood framework (where parameters are chosen to make the likelihood as large as possible), the likelihood can be maximised for each parameter separately. So, by partially differentiating (10.3) with respect to c_j and equating to zero, we obtain the maximum likelihood estimate of c_j as

$$\hat{c}_j = \frac{r_j}{\int_0^T g_j(x(t)) dt}, \quad j = 1, \dots, v. \quad (10.4)$$

In the context of Bayesian inference, the factorisation means that if independent prior distributions are adopted for the rate constants, then this independence will be retained *a posteriori*. It is also clear from the form of (10.3) that the complete-data likelihood is conjugate to an independent gamma prior for the rate constants. Thus, adopting priors for the rate constants of the form

$$\pi(\mathbf{c}) = \prod_{j=1}^v \pi(c_j), \quad c_j \sim \Gamma(a_j, b_j), \quad j = 1, \dots, v,$$

we can use Bayes Theorem (9.1) to obtain

$$\begin{aligned} \pi_j(c_j|\mathbf{x}) &\propto \pi(c_j)L(c_j; \mathbf{x}) \\ &\propto c_j^{a_j-1} \exp\{-b_j c_j\} c_j^{r_j} \exp\left\{-c_j \int_0^T g_j(x(t))dt\right\} \\ &= c_j^{a_j+r_j-1} \exp\left\{-c_j \left[b_j + \int_0^T g_j(x(t))dt\right]\right\}, \end{aligned}$$

that is

$$c_j|\mathbf{x} \sim \Gamma\left(a_j + r_j, b_j + \int_0^T g_j(x(t))dt\right), \quad j = 1, \dots, v. \quad (10.5)$$

So, in the context of complete data, the inference problem is straightforward. However, even in the context of less than complete data, the distributions in (10.5) represent full-conditionals for the rate constants, and thus form an important part of an MCMC algorithm for inferring the rate constants given discrete-time observations.

10.3 Discrete-time observations of the system state

It is, of course, unrealistic to suppose that it is possible to perfectly observe the entire sample path of the process, and so this assumption can be gradually weakened, and its effects on the inference process considered. It seems sensible to begin by considering the case of perfect observation of the system state at a finite collection of equally spaced times. In order to simplify notation, we will assume without loss of generality that we observe the process at integer times $t = 0, 1, \dots, T$. That is, we have a total of $T + 1$ observations on $x(t)$. We have seen in the previous section how to sample from the full-conditionals for the rate constants given a complete sample path on $[0, T]$, so an obvious strategy is to develop an MCMC algorithm which includes all of the missing parts of the sample path, and then simulate and average over the set of all plausible sample paths in the correct way.

Conditional on the $T + 1$ observations (and the rate constants), the stochastic process breaks up into T conditionally independent processes on unit intervals with given fixed end points. The idea is to construct an MCMC algorithm which cycles through each interval in turn, sampling an appropriate path from its relevant full-conditional distribution, and then completing the iteration by sampling the rates conditional on the full sample path. We therefore need to be able to sample paths on an interval of unit length given its end points. For notational convenience we will consider without loss of generality the interval $[0, 1]$. Of course we know several exact methods for sampling the path conditional only on the left-hand end-point, $x(0)$ (including Gillespie's direct method). However, here we need to sample from the interval given both $x(0)$ and $x(1)$. This turns out to be *much* more difficult to do directly, and so a simpler strategy is adopted. The idea is that because we are working in the context of an MCMC algorithm, we can use a Metropolis-Hastings move to update the sample path, resulting in a Metropolis-within-Gibbs style overall algorithm. We therefore need only a method that will sample from a distribution over the

space of bridging sample paths that has the same support as the true bridging process. Then we can use an appropriate Metropolis-Hastings acceptance probability in order to correct for the approximate step. An outline of the proposed MCMC algorithm can be stated as follows.

1. Initialise the algorithm with a valid sample path consistent with the observed data.
2. Sample rate constants from their full conditionals given the current sample path.
3. For each of the T intervals, propose a new sample path consistent with its end-points and accept/reject it with a Metropolis-Hastings step.
4. Output the current rate constants.
5. Return to step 2.

In order to make progress with this problem, some notation is required. To keep the notation as simple as possible, we will now redefine some notation for the unit interval $[0, 1]$ which previously referred to the entire interval $[0, T]$. So now

$$\mathbf{x} = \{x(t) : t \in [0, 1]\}$$

denotes the “true” sample path that is only observed at times $t = 0$ and $t = 1$, and

$$\mathbf{X} = \{X(t) : t \in [0, 1]\}$$

represents the stochastic process that gives rise to \mathbf{x} as a single observation. Our problem is that we would like to sample directly from the distribution $(\mathbf{X} | x(0), x(1), c)$, but this is difficult, so instead we will content ourselves with constructing a Metropolis-Hastings update that has $\pi(\mathbf{x} | x(0), x(1), c)$ as its target distribution. Let us also re-define $\mathbf{r} = (r_1, \dots, r_v)'$ to be the numbers of reaction events in the interval $[0, 1]$, and $n = \sum_{j=1}^v r_j$. It is clear that knowing both $x(0)$ and $x(1)$ places some constraints on \mathbf{r} , but it will not typically determine it completely. It turns out to be easiest to sample a new interval in two stages: first pick an \mathbf{r} consistent with the end constraints and then sample a new interval conditional on $x(0)$ and \mathbf{r} . So, ignoring the problem of sampling \mathbf{r} for the time being, we would ideally like to be able to sample from $\pi(\mathbf{x} | x(0), \mathbf{r}, c)$, but this is still quite difficult to do directly. At this point it is helpful to think of the u -component sample path \mathbf{X} as being a function of the v -component point process of reaction events. This point process is hard to simulate directly as its hazard function is random, but the hazards are known at the end-points $x(0)$ and $x(1)$, and so they can probably be reasonably well approximated by v -independent inhomogeneous Poisson processes whose rates vary linearly between the rates at the end points. In order to make this work, we need to be able to sample from an inhomogeneous Poisson process conditional on the number of events. This requires some Poisson process theory not covered in Chapter 5.

Lemma 10.1 *For given fixed $\lambda, p > 0$, consider $N \sim Po(\lambda)$ and $X | N \sim B(N, p)$. Then marginally we have*

$$X \sim Po(\lambda p).$$

Proof.

$$\begin{aligned}
 P(X = k) &= \sum_{i=0}^{\infty} P(X = k|N = i) P(N = i) \\
 &= \sum_{i=k}^{\infty} P(X = k|N = i) P(N = i) \\
 &= \sum_{i=k}^{\infty} \frac{i!}{k!(i-k)!} p^k (1-p)^{i-k} \times \frac{\lambda^i e^{-\lambda}}{i!} \\
 &= \sum_{i=k}^{\infty} \frac{\lambda^i e^{-\lambda} p^k (1-p)^{i-k}}{k!(i-k)!} \\
 &= \sum_{i=0}^{\infty} \frac{\lambda^{i+k} e^{-\lambda} p^k (1-p)^i}{k!i!} \\
 &= \frac{\lambda^k e^{-\lambda} p^k}{k!} \sum_{i=0}^{\infty} \frac{\lambda^i (1-p)^i}{i!} \\
 &= \frac{\lambda^k e^{-\lambda} p^k}{k!} e^{\lambda(1-p)} \\
 &= \frac{(\lambda p)^k e^{-\lambda p}}{k!}, \quad k = 0, 1, \dots
 \end{aligned}$$

□

Proposition 10.1 Consider a homogeneous Poisson process with rate λ on the finite interval $[0, T]$. Let $R \sim Po(\lambda T)$ be the number of events. Then conditional on R , the (unsorted) event times are $U(0, T)$ random variables. In other words, the ordered event times correspond to R uniform order statistics.

Proof. We will just give a sketch proof of this fact that is anyway intuitively clear. The key property of the Poisson process is that in a finite interval $[a, b] \subseteq [0, T]$, the number of events is $Po(\lambda[b-a])$. So we will show that if a point process is constructed by first sampling R and then scattering R points uniformly, then the correct number of events will be assigned to the interval $[a, b]$. It is clear that for any particular event, the probability that it falls in the interval $[a, b]$ is $(b-a)/T$. Now, letting X denote the number of events in the interval $[a, b]$, it is clear that

$$(X|R = r) \sim B(r, (b-a)/T)$$

and so by the previous Lemma,

$$\begin{aligned}
 X &\sim Po(\lambda T \times (b-a)/T) \\
 &= Po(\lambda[b-a]).
 \end{aligned}$$

□

So we now have a way of simulating a homogeneous Poisson process conditional on the number of events. However, one way of thinking about an inhomogeneous

Poisson process is as a homogeneous Poisson process with time rescaled in a non-linear way.

Proposition 10.2 *Let X be a homogeneous Poisson process on the interval $[0, 1]$ with rate $\mu = (h_0 + h_1)/2$, and let Y be an inhomogeneous Poisson process on the same interval with rate $\lambda(t) = (1-t)h_0 + th_1$, for given fixed $h_0 \neq h_1$, $h_0, h_1 > 0$. A realisation of the process Y can be obtained from a realisation of the process X by applying the time transformation*

$$t := \frac{\sqrt{h_0^2 + \{h_1^2 - h_0^2\}t} - h_0}{h_1 - h_0}$$

to the event times of the X process.

Proof. Process X has cumulative hazard $M(t) = t(h_0 + h_1)/2$, while process Y has cumulative hazard

$$\Lambda(t) = \int_0^t [(1-t)h_0 + th_1]dt = h_0t + \frac{t^2}{2}(h_1 - h_0).$$

Note that the cumulative hazards for the two processes match at both $t = 0$ and $t = 1$, and so one process can be mapped to the other by distorting time to make the cumulative hazards match also at intermediate times. Let the local time for the X process be s and the local time for the Y process be t . Then setting $M(s) = \Lambda(t)$ gives

$$\begin{aligned} \frac{s}{2}(h_0 + h_1) &= h_0t + \frac{t^2}{2}(h_1 - h_0) \\ \Rightarrow 0 &= \frac{t^2}{2}(h_1 - h_0) + h_0t - \frac{s}{2}(h_0 + h_1) \\ \Rightarrow t &= \frac{-h_0 + \sqrt{h_0^2 + (h_1 - h_0)(h_0 + h_1)s}}{h_1 - h_0}. \end{aligned}$$

□

So, we can sample an inhomogeneous Poisson process conditional on the number of events by first sampling a homogeneous Poisson process with the average rate conditional on the number of events and then transforming time to get the correct inhomogeneity.

In order to correct for the fact that we are not sampling from the correct bridging process, we will need a Metropolis-Hastings acceptance probability that will depend both on the likelihood of the sample path under the true model and the likelihood of the sample path under the approximate model. We have already calculated the likelihood under the true model (the complete-data likelihood). We now need the likelihood under the inhomogeneous Poisson process model.

Proposition 10.3 *The complete data likelihood for a sample path x on the interval $[0, 1]$ under the approximate inhomogeneous Poisson process model is given by*

$$L_A(c; \mathbf{x}) = \left\{ \prod_{i=1}^n \lambda_{\nu_i}(t_i) \right\} \exp \left\{ -\frac{1}{2} [h_0(x(0), c) + h_0(x(1), c)] \right\},$$

where $\lambda_j(t) = (1-t)h_j(x(0), c) + th_j(x(1), c)$, $j = 1, \dots, v$.

Proof. Again, we will compute this in an informal way. Define the cumulative rates

$$\Lambda_j(t) = \int_0^t \lambda_j(t) dt,$$

the combined rate $\lambda_0(t) = \sum_{j=1}^v \lambda_j(t)$, and the cumulative combined rate $\Lambda_0(t) = \sum_{j=1}^v \Lambda_j(t)$. Considering the i th event, the density of the time and type is given by

$$\begin{aligned} \lambda_0(t_i) \exp\{-[\Lambda_0(t_i) - \Lambda_0(t_{i-1})]\} &\times \frac{\lambda_{\nu_i}(t_i)}{\lambda_0(t_i)} \\ &= \lambda_{\nu_i}(t_i) \exp\{-[\Lambda_0(t_i) - \Lambda_0(t_{i-1})]\}, \end{aligned}$$

and the probability of no event after time t_n is given by

$$\exp\{-[\Lambda_0(T) - \Lambda_0(t_n)]\}.$$

This leads to a combined likelihood of the form

$$\begin{aligned} L_A(c; \mathbf{x}) &= \left\{ \prod_{i=1}^n \lambda_{\nu_i}(t_i) \exp\{-[\Lambda_0(t_i) - \Lambda_0(t_{i-1})]\} \right\} \times \exp\{-[\Lambda_0(T) - \Lambda_0(t_n)]\} \\ &= \left\{ \prod_{i=1}^n \lambda_{\nu_i}(t_i) \right\} \left\{ \prod_{i=1}^{n+1} \exp\{-[\Lambda_0(t_i) - \Lambda_0(t_{i-1})]\} \right\} \\ &= \left\{ \prod_{i=1}^n \lambda_{\nu_i}(t_i) \right\} \exp\left\{-\sum_{i=1}^{n+1} [\Lambda_0(t_i) - \Lambda_0(t_{i-1})]\right\} \\ &= \left\{ \prod_{i=1}^n \lambda_{\nu_i}(t_i) \right\} \exp\{-\Lambda_0(T)\} \\ &= \left\{ \prod_{i=1}^n \lambda_{\nu_i}(t_i) \right\} \exp\left\{-\frac{1}{2}[h_0(x(0), c) + h_0(x(1), c)]\right\}. \end{aligned}$$

□

As we will see, the complete-data likelihoods occur in the Metropolis-Hastings acceptance probability in the form of the ratio $L(c; \mathbf{x})/L_A(c; \mathbf{x})$. This ratio is clearly just

$$\begin{aligned} \frac{L(c; \mathbf{x})}{L_A(c; \mathbf{x})} &= \left\{ \prod_{i=1}^n \frac{h_{\nu_i}(x(t_{i-1}), c_i)}{\lambda_{\nu_i}(t_i)} \right\} \\ &\times \exp\left\{\frac{1}{2}[h_0(x(0), c) + h_0(x(1), c)] - \int_0^1 h_0(x(t), c) dt\right\}. \end{aligned}$$

From a more advanced standpoint, this likelihood ratio is seen to be the Radon-Nikodym derivative $\frac{d\mathbb{P}}{d\mathbb{Q}}(\mathbf{x})$ of the true Markov jump process (\mathbb{P}) with respect to the linear inhomogeneous Poisson process approximation (\mathbb{Q}), and may be derived primitively and directly from the properties of the jump processes in an entirely rigorous

way. The Radon-Nikodym derivative measures the “closeness” of the approximating process to the true process, in the sense that the more closely the processes match, the closer the derivative will be to 1.

We are now in a position to state the basic form of the Metropolis-Hastings update of the interval $[0, 1]$. First a proposed new r vector will be sampled from an appropriate proposal distribution with PMF $f(r^*|r)$ (we will discuss appropriate ways of constructing this later). Then conditional on r^* , sample a proposed sample path \mathbf{x}^* from the approximate process and accept the pair (r^*, \mathbf{x}^*) with probability $\min\{1, A\}$ where

$$\begin{aligned} A &= \frac{\pi(\mathbf{x}^*|x(0), x(1), c)}{\pi(\mathbf{x}|x(0), x(1), c)} \bigg/ \frac{f(r^*|r)\pi_A(\mathbf{x}^*|x(0), r^*, c)}{f(r|r^*)\pi_A(\mathbf{x}|x(0), r, c)} \\ &= \frac{\pi(\mathbf{x}^*|x(0), c)}{\pi_A(\mathbf{x}^*|x(0), c)} \times \frac{q(r^*)}{f(r^*|r)} \\ &= \frac{\pi(\mathbf{x}|x(0), c)}{\pi_A(\mathbf{x}|x(0), c)} \times \frac{q(r)}{f(r|r^*)} \\ &= \frac{L(c; \mathbf{x}^*)}{L_A(c; \mathbf{x}^*)} \times \frac{q(r^*)}{f(r^*|r)}, \\ &= \frac{L(c; \mathbf{x})}{L_A(c; \mathbf{x})} \times \frac{q(r)}{f(r|r^*)}, \end{aligned}$$

where $q(r)$ is the PMF of r under the approximate model. That is,

$$q(r) = \prod_{j=1}^v q_j(r_j),$$

where $q_j(r_j)$ is the PMF of a Poisson with mean $[h_j(x(0), c) + h_j(x(1), c)]/2$. Again, we could write this more formally as

$$A = \frac{\frac{d\mathbb{P}}{d\mathbb{Q}}(\mathbf{x}^*)}{\frac{d\mathbb{P}}{d\mathbb{Q}}(\mathbf{x})} \times \frac{\frac{q(r^*)}{f(r^*|r)}}{\frac{q(r)}{f(r|r^*)}}.$$

So now the only key aspect of the MCMC algorithm that has not yet been discussed is the choice of the proposal distribution $f(r^*|r)$. Again, ideally we would like to sample directly from the true distribution of r given $x(0)$ and $x(1)$, but this is not straightforward. Instead we simply want to pick a proposal that effectively explores the space of r s consistent with the end points. Recalling the discussion of Petri nets from Section 2.3, to a first approximation the set of r that we are interested in is the set of all non-negative integer solutions in r to

$$\begin{aligned} x(1) &= x(0) + Sr \\ \Rightarrow Sr &= x(1) - x(0). \end{aligned} \tag{10.6}$$

There will be some solutions to this equation that do not correspond to possible sample paths, but there will not be many of these. Note that given a valid solution r , then $r + x$ is another valid solution, where x is any T -invariant of the Petri net. Thus,

the set of all solutions is closely related to the set of all T -invariants of the associated Petri net. Assuming that S is of full rank (if S is not of full rank, the dimension reducing techniques from Section 7.3 can be used to reduce the model until it is), then the space of solutions will have dimension $v - u$. One way to explore this space is to permute the columns of S so that the first u columns represent an invertible $u \times u$ matrix, \tilde{S} . Denoting the remaining columns of S by \vec{S} , we partition S as $S = (\tilde{S}, \vec{S})$. We similarly partition $r = \begin{pmatrix} \tilde{r} \\ \vec{r} \end{pmatrix}$, with \tilde{r} representing the first u elements of r . We then have

$$\begin{aligned} Sr &= (\tilde{S}, \vec{S}) \begin{pmatrix} \tilde{r} \\ \vec{r} \end{pmatrix} = \tilde{S}\tilde{r} + \vec{S}\vec{r} = x(1) - x(0) \\ \Rightarrow \tilde{r} &= \tilde{S}^{-1} [x(1) - x(0) - \vec{S}\vec{r}]. \end{aligned} \quad (10.7)$$

Equation (10.7) suggests a possible strategy for exploring the solution space. Starting from a valid solution r , one can perturb the elements corresponding to \vec{r} in an essentially arbitrary way, and then set the elements of \tilde{r} accordingly. Of course, there is a chance that some element will go negative, but such moves will be immediately rejected. One possible choice of symmetric proposal for updating the elements of \vec{r} is to add to each element the difference between two independent Poisson random quantities with the same mean, ω . If the tuning parameter ω is chosen to be independent of the current state, then the proposal distribution is truly symmetric and the relevant PMF terms cancel out of the Metropolis-Hastings acceptance probability. However, it is sometimes useful to allow ω to be a function of the current state of the chain (for example, allowing ω to be larger when the state is larger), and in this case the PMF for this proposal is required. The PMF is well known and is given by $p(y) = e^{-2\omega} I_y(2\omega)$, where $I_y(\cdot)$ is a regular modified Bessel function of order y ; see Johnson & Kotz (1969) and Abramowitz & Stegun (1984) for further details.

Now, given a correctly initialised Markov chain, we have everything we need to implement a MCMC algorithm that will have as its equilibrium distribution the exact posterior distribution of the rate constants given the (perfect) discrete time observations. However, it turns out that even the problem of initialising the chain is not entirely straightforward. Here we need a valid solution to (10.6) for each of the T unit intervals. The constraints encoded by (10.6) correspond to the constraints of an integer linear programming problem. By adopting an (essentially arbitrary) objective function $\min\{\sum_j r_j\}$, we can use standard algorithms for the solution of integer linear programming problems in order to initialise the chain.

Some “proof-of-concept” software that implements an approximate version of this algorithm is available, called `stochInf`. Note, however, that this software is not intended for practical application, as it assumes perfect observations on all species in the model. For details of the approximation used by the software, see the discussion of the approximate algorithm in Boys, Wilkinson & Kirkwood (2004). Also see Boys et al. (2004) for a detailed discussion of the application of these techniques to the stochastic Lotka-Volterra model. The paper also discusses an alternative MCMC algorithm and the extension of these ideas to deal with partial observation of the system state. Similar extensions can also be used to incorporate measurement error into the models.

10.4 Diffusion approximations for inference

The discussion in the previous section demonstrates that it is possible to construct exact MCMC algorithms for inference in discrete stochastic kinetic models based on discrete time observations (and it is possible to extend the techniques to more realistic data scenarios than those directly considered). The discussion gives great insight into the nature of the inferential problem and its conceptual solution. However, there is a slight problem with the techniques discussed there in the context of the relatively large and complex models of genuine interest to systems biologists. It should be clear that each iteration of the MCMC algorithm described in the previous section is more computationally demanding than simulating the process exactly using Gillespie's direct method (for the sake of argument, let us say that it is one order of magnitude more demanding). For satisfactory inference, a large number of MCMC iterations will be required. For models of the complexity discussed in the previous section, it is not uncommon for 10^7 – 10^8 iterations to be required for satisfactory convergence to the true posterior distribution. Using such methods for inference therefore has a computational complexity of 10^8 – 10^9 times that required to simulate the process. As if this were not bad enough, it turns out that MCMC algorithms are particularly difficult to parallelise effectively (Wilkinson 2005). One possible approach to improving the situation is to approximate the algorithm with a much faster one that is less accurate, as discussed in Boys et al. (2004). Unfortunately even that approach does not scale up well to genuinely interesting problems, so a different approach is required.

A similar problem was considered in Chapter 8, from the viewpoint of simulation rather than inference. We saw there how it was possible to approximate the true Markov jump process by the chemical Langevin equation (CLE), which is the diffusion process that behaves most like the true jump process. It was seen there how simulation of the CLE can be many orders of magnitude faster than an exact algorithm. This suggests the possibility of using the CLE as an approximate model for inferential purposes. It turns out that the CLE provides an excellent model for inference, even in situations where it does not perform particularly well as a simulation model. This observation at first seems a little counter-intuitive, but the reason is that in the context of inference, one is conditioning on data from the true model, and this helps to calibrate the approximate model and stop MCMC algorithms from wandering off into parts of the space that are plausible in the context of the approximate model, but not in the context of the true model.

What is required is a method for inference for general non-linear multivariate diffusion processes observed partially, discretely and with error. Unfortunately this too turns out to be a highly non-trivial problem, and is still the subject of a great deal of ongoing research. Such inference problems arise often in financial mathematics and econometrics, and so much of the literature relating to this problem can be found in that area; see Durham & Gallant (2002) for an overview.

The problem with diffusion processes is that any finite sample path contains an infinite amount of information, and so the concept of a complete-data likelihood does not exist in general. We will illustrate the problem in the context of high-resolution time-course data on the CLE. Starting with the CLE in the form of (8.3), define

$\mu(x, c) = Sh(x, c)$ and $\beta(x, c) = S \text{diag} \{h(x, c)\} S'$ to get the u -dimensional diffusion process

$$dX_t = \mu(X_t, c)dt + \sqrt{\beta(X_t, c)}dW_t.$$

We will assume that for some small timestep Δt we have data $x = (x_0, x_{\Delta t}, x_{2\Delta t}, \dots, x_{n\Delta t})$ (that is, $n + 1$ observations), and that the timestep is sufficiently small for the Euler-Maruyama approximation to be valid, leading to the difference equation

$$\Delta X_t \equiv X_{t+\Delta t} - X_t \simeq \mu(X_t, c)\Delta t + \sqrt{\beta(X_t, c)}\Delta W_t. \tag{10.8}$$

At this point we require some multivariate normal theory (not explicitly covered in the text). Equation (10.8) corresponds to the distributional statement

$$X_{t+\Delta t}|X_t, c \sim N(X_t + \mu(X_t, c)\Delta t, \beta(X_t, c)\Delta t),$$

where $N(\cdot, \cdot)$ here refers to the multivariate normal distribution, parameterised by its mean vector and covariance matrix. The probability density associated with this increment is given by

$$\begin{aligned} \pi(x_{t+\Delta t}|x_t, c) &= N(x_{t+\Delta t}; x_t + \mu(x_t, c)\Delta t, \beta(x_t, c)\Delta t) \\ &= (2\pi)^{-u/2} |\beta(x_t, c)\Delta t|^{-1/2} \\ &\quad \times \exp \left\{ -\frac{1}{2} (\Delta x_t - \mu(x_t, c)\Delta t)' [\beta(x_t, c)\Delta t]^{-1} (\Delta x_t - \mu(x_t, c)\Delta t) \right\} \\ &= (2\pi\Delta t)^{-u/2} |\beta(x_t, c)|^{-1/2} \\ &\quad \times \exp \left\{ -\frac{\Delta t}{2} \left(\frac{\Delta x_t}{\Delta t} - \mu(x_t, c) \right)' \beta(x_t, c)^{-1} \left(\frac{\Delta x_t}{\Delta t} - \mu(x_t, c) \right) \right\}. \end{aligned}$$

If S is not of full rank, then $\beta(x_t, c)$ will not be invertible. This can be tackled either by reducing the dimension of the model so that S is of full rank, or by using the Moore-Penrose generalised inverse of $\beta(x_t, c)$ wherever the inverse occurs in a multivariate normal density. It is important to understand that diffusion sample paths are not differentiable, so the quantity $\Delta x_t/\Delta t$ does not have a limit as Δt tends to zero. It is now possible to derive the likelihood associated with this set of

observations as

$$\begin{aligned}
 L(c; x) &= \pi(x|c) \\
 &= \pi(x_0|c) \prod_{i=0}^{n-1} \pi(x_{(i+1)\Delta t} | x_{i\Delta t}, c) \\
 &= \pi(x_0|c) \prod_{i=0}^{n-1} (2\pi\Delta t)^{-u/2} |\beta(x_{i\Delta t}, c)|^{-1/2} \\
 &\quad \times \exp \left\{ -\frac{\Delta t}{2} \left(\frac{\Delta x_{i\Delta t}}{\Delta t} - \mu(x_{i\Delta t}, c) \right)' \beta(x_{i\Delta t}, c)^{-1} \left(\frac{\Delta x_{i\Delta t}}{\Delta t} \right. \right. \\
 &\quad \left. \left. - \mu(x_{i\Delta t}, c) \right) \right\}.
 \end{aligned}$$

Now assuming that $\pi(x_0|c)$ is in fact independent of c , we can simplify the likelihood as

$$\begin{aligned}
 L(c; x) &\propto \left\{ \prod_{i=0}^{n-1} |\beta(x_{i\Delta t}, c)|^{-1/2} \right\} \times \\
 &\exp \left\{ -\frac{1}{2} \sum_{i=0}^{n-1} \left(\frac{\Delta x_{i\Delta t}}{\Delta t} - \mu(x_{i\Delta t}, c) \right)' \beta(x_{i\Delta t}, c)^{-1} \left(\frac{\Delta x_{i\Delta t}}{\Delta t} - \mu(x_{i\Delta t}, c) \right) \right. \\
 &\quad \left. \Delta t \right\}. \quad (10.9)
 \end{aligned}$$

Equation (10.9) is the closest we can get to a complete-data likelihood for the CLE, as it does not have a limit as Δt tends to zero.[†] In the case of perfect high-resolution observations on the system state, (10.9) represents the likelihood for the problem, which could be maximised (numerically) in the context of maximum likelihood estimation or combined with a prior to form the kernel of a posterior distribution for c . In this case there is no convenient conjugate analysis, but it is entirely straightforward to implement a Metropolis random walk MCMC sampler to explore the posterior distribution of c .

Of course it is unrealistic to assume perfect observation of all states of a model, and in the biological context, it is usually unrealistic to assume that the sampling frequency will be sufficiently high to make the Euler approximation sufficiently accurate. So, just as for the discrete case, MCMC algorithms can be used in order to “fill-in” all of the missing information in the model.

MCMC algorithms for multivariate diffusions can be considered conceptually in a similar way to those discussed in the previous section. Given (perfect) discrete-time

[†] If it were the case that $\beta(x, c)$ was independent of c , then we could drop terms no longer involving c to get an expression that is well behaved as Δt tends to zero. In this case, the complete data likelihood is the exponential of the sum of two integrals, one of which is a regular Riemann integral, and the other is an Itô stochastic integral. Unfortunately this rather elegant result is of no use to us here, as the diffusion matrix of the CLE depends on c in a fundamental way.

observations, the diffusion process breaks into a collection of multi-variate diffusion bridges that need to be sampled and averaged over in the correct way. However, there is a complication in the context of diffusion processes that does not occur in the context of Markov jump processes due to the infinite amount of information in a finite continuous diffusion sample path. This means both that it is impossible to fully impute a sample path, and also that even if it were possible, the full conditionals for the diffusion parameters would be degenerate, leading to a reducible MCMC algorithm. An informative discussion of this problem, together with an elegant solution in the context of univariate diffusions, can be found in Roberts & Stramer (2001). Fortunately, it turns out that by working with a time discretisation of the CLE (such as the Euler discretisation), both of these problems can be side-stepped (but not actually solved). Here the idea is to introduce a finite number of time points between each pair of observations, and implement an MCMC algorithm to explore the space of “skeleton” bridges between the observations. This is relatively straightforward, and there are a variety of different approaches that can be used for implementation purposes. Then because the sample path is represented by just a finite number of points, the full conditionals for the diffusion parameters are not degenerate, and the MCMC algorithm will have as its equilibrium distribution the exact posterior distribution of the rate parameters given the data, conditional on the Euler approximation to the CLE being the true model. Of course neither the CLE nor the Euler approximation to it are actually the true model, but this will hopefully have little effect on inferences. Such an approach to inferring the rate constants of stochastic kinetic models is discussed in Golightly & Wilkinson (2005a), to which the reader is referred for further details. This paper also discusses the case of partial observation of the system state and includes an application to a genetic regulatory network.

As mentioned in the previous paragraph, the technique of discretising time has the effect of side-stepping the technical problems of inference, but does not actually solve them. In particular, it is desirable to impute many latent points between each pair of observations, so that the discretisation bias is minimised and the skeleton bridges represent a good approximation to the true process on the same skeleton of points. Unfortunately, using a fine discretisation has the effect of making the full-conditionals for the diffusion parameters close to singular, which in turn makes the MCMC algorithm mix extremely poorly. This puts the modeller in the unhappy position of having to choose between a poor approximation to the CLE model or a poorly mixing MCMC algorithm. Clearly a more satisfactory solution is required.

It turns out that by adopting a sequential MCMC algorithm (where the posterior distribution is updated one time point at a time, rather than with a global algorithm that incorporates the information from all time points simultaneously), it is possible to solve the mixing problems by jointly updating the diffusion parameters and diffusion bridges, thereby breaking the dependence between them. The sampling mechanism relies heavily on the *modified diffusion bridge* construct of Durham & Gallant (2002). A detailed discussion of this algorithm is given in Golightly & Wilkinson (2005b), and its application to inference for stochastic kinetic models is explored in Golightly & Wilkinson (2006). This last paper discusses the most realistic setting of multiple data sets that are partial, discrete, and include measurement error. One of the

advantages of using a sequential algorithm is that it is very convenient to use in the context of multiple data sets from different cells or experiments, possibly measuring different species in each experiment. A discussion of this issue in the context of an illustrative example is given in Golightly & Wilkinson (2006).

Given this discussion, it is clearly possible to develop a generic piece of software for inferring the rate parameters of the CLE for realistic experimental data (at the single-cell level) using MCMC techniques. At the time of writing, I am currently developing such an application (`stochInf2`), and a link to it will be posted on the website for this book when it becomes available.

10.5 Network inference

At this point it is worth saying a few words regarding network inference. It is currently fashionable in some parts of the literature to attempt to deduce the structure of biochemical networks *ab initio* from routinely available experimental data. While it is clearly possible to develop computational algorithms to do this, the utility of doing so is not at all clear due to the fact that there will typically be a very large number of distinct network structures all of which are consistent with the available experimental data (large numbers of these will be biologically implausible, but a large number will also be quite plausible). In this case the “best fitting” network is almost certainly incorrect. For the time being it seems more prudent to restrict attention to the less ambitious goal of comparing the experimental support for a small number of competing network structures. Typically this will concern a relatively well-characterised biochemical network where there is some uncertainty as to whether one or two of the potential reaction steps actually take place. In this case, deciding whether or not a given reaction is present is a problem of discrimination between two competing network structures. There are a number of ways that this problem could be tackled. The simplest approach would be to fit all competing models using the MCMC techniques outlined in the previous sections and compute the marginal likelihoods associated with each (Chib 1995) in order to compute Bayes factors. More sophisticated strategies might adopt reversible jump MCMC techniques (Green 1995) in order to simultaneously infer parameters and structure. In principle, the reversible jump techniques could also be used for *ab initio* network inference, but note well the previous caveat. In particular, note that inferences are sensitive not only to the prior adopted over the set of models to be considered, but also to the prior structure adopted for the rate constants conditional on the model. Simultaneous inference for parameters and network structure is currently an active research area.

10.6 Exercises

1. Pick a simple model (say, from Chapter 7) and simulate some complete data from it using the Gillespie algorithm.
 - (a) Use these data to form the complete-data likelihood (10.2).
 - (b) Maximise the complete data likelihood using (10.4).

- (c) How well are you able to recover the true rate constants? How much data do you need? Does this vary according to the true rate constants originally chosen?
 - (d) Compute the Bayesian posterior distributions for the rate constants starting with a fairly diffuse prior (say, a $\Gamma(0.1, 0.1)$ distribution). Compare the mean of the posterior to the maximum likelihood estimates (MLEs).
 - (e) Consider the variance of the posterior. How does this change as the amount of available data changes? How does this shed light on this issue of how many data are required for reliable parameter inference?
2. Simulate some data from the CLE approximation to a simple model on a fairly fine time grid.
- (a) Calculate the likelihood of the simulated data and maximise it numerically to find the MLEs. How well can you recover the true rate constants?
 - (b) Implement a simple Metropolis-Hastings algorithm to compute the posterior distribution of the rate constants. How many data are needed to reduce the posterior uncertainty to an acceptable level?
 - (c) Investigate the relationship between time and sampling frequency. Is it better to have 1,000 observations over a 10-second period or a 10-hour period? How does this vary over models and rate constants?
 - (d) Redo this exercise using data simulated exactly using the Gillespie algorithm and then discretised onto a regular grid. How much bias does the CLE approximation introduce?
3. Read the referenced papers on MCMC algorithms for rate parameter inference given time-course data, and try coding up one of the algorithms for the special case of the linear birth-death process. Does the algorithm correctly recover both parameters of the process? How many data are required?

10.7 Further reading

Application of Bayesian inference to estimation of discrete stochastic kinetic models is examined in Boys et al. (2004). Note that this builds on previous work for the estimation of stochastic compartmental models; see Gibson & Renshaw (1998), for example. An overview of the problem of inference for multivariate diffusion processes is given in Durham & Gallant (2002), and an application to stochastic kinetic models is discussed in Golightly & Wilkinson (2005a). An effective MCMC algorithm for inference is presented in Golightly & Wilkinson (2005b), and is applied to stochastic kinetic models in Golightly & Wilkinson (2006). This is an active research area that I am particularly interested in. I keep an up-to-date publication list on the web (linked from this book's website), which might be worth consulting for recent developments.

Conclusions

This book has provided an introduction to the concepts of stochastic modelling relevant for systems biology applications and illustrated those concepts in a systems biology context. The main area of application has been the modelling and simulation of genetic and biochemical networks, but it should be reasonably clear that the techniques can also be applied to a range of other modelling scenarios. Most interesting biological processes exhibit stochastic variation, and understanding the nature and effect of the randomness can often be fundamental to understanding the process.

Although Chapters 9 and 10 have touched on issues relating to inference, the primary focus of the text has been on the development of models that adequately capture system dynamics and algorithms for simulation, rather than on statistical issues associated with wet-lab systems biology work. Indeed, there is clearly scope for an entire book on the design and analysis of statistical experiments for systems biology. Such a book would consider general problems of statistical inference for networks and rate constants from routinely available experimental data and would also tackle the difficult experimental design issues that arise in this context. Most existing work in that area builds on very simple descriptive models of the underlying processes, rather than the detailed “generative” models that have been considered in this text. This is positive in that it renders the associated design and analysis problems more tractable, but also places limitations on what can be achieved. For example, a simple dynamic Bayesian network (DBN) fitted to some time-course micro-array data is not attempting to capture the true system dynamics, and so cannot be translated back to a detailed generative model of the kind we would like to be able to simulate. Essentially, there is currently a “gap” between the detailed modelling work that has been considered here and the applied statistical work currently being carried out by bioinformaticians.

Although Chapter 10 has hinted at some of the issues relating to inference for stochastic kinetic models, it is not completely straightforward to use these techniques directly in conjunction with routinely available high-throughput bioinformatics data. There is therefore a pressing need to develop techniques that allow routine experimental data to be used for “calibrating” detailed simulation models. While this is reasonably straightforward for simple models that are deterministic and fast to simulate, calibration of stochastic models (or deterministic models that are slow to simulate) is not easy. For deterministic models that are slow to simulate, there is a theory of Bayesian model calibration (Kennedy & O’Hagan 2001, Santner, Williams & Notz 2003) that can be applied (though it is worth noting that most existing literature in this area is concerned with models having a low-dimensional output, which is typically not the case in the systems biology context). Although it is in principle

possible to extend model calibration technology for the fitting of stochastic models, it is not straightforward to do so in practice. The development of fast and approximate techniques for calibrating stochastic simulation models is currently an active research area.

In conclusion, this text has covered only the essential concepts required for beginning to think seriously about systems biology from a stochastic viewpoint. The really interesting statistical problems concerned with marrying stochastic systems biology models to routinely available high-throughput experimental data remain largely unsolved and the subject of a great deal of ongoing research.

SBML Models

This appendix contains a few examples of complete SBML documents describing models discussed in the text. All of these models, and many others, can be downloaded from this book's website.

A.1 Auto-regulatory network

The model below is the SBML for the auto-regulatory network discussed in Chapter 2.

```
<?xml version="1.0" encoding="UTF-8"?>
<sbml xmlns="http://www.sbml.org/sbml/level2" level="2" version="1">
  <model id="AutoRegulatoryNetwork" name="Auto-regulatory network">
    <listOfUnitDefinitions>
      <unitDefinition id="substance">
        <listOfUnits>
          <unit kind="item" multiplier="1" offset="0"/>
        </listOfUnits>
      </unitDefinition>
    </listOfUnitDefinitions>
    <listOfCompartments>
      <compartment id="Cell"/>
    </listOfCompartments>
    <listOfSpecies>
      <species id="Gene" compartment="Cell" initialAmount="10"
        hasOnlySubstanceUnits="true"/>
      <species id="P2Gene" name="P2.Gene" compartment="Cell"
        initialAmount="0" hasOnlySubstanceUnits="true"/>
      <species id="Rna" compartment="Cell" initialAmount="0"
        hasOnlySubstanceUnits="true"/>
      <species id="P" compartment="Cell" initialAmount="0"
        hasOnlySubstanceUnits="true"/>
      <species id="P2" compartment="Cell" initialAmount="0"
        hasOnlySubstanceUnits="true"/>
    </listOfSpecies>
    <listOfReactions>
      <reaction id="RepressionBinding" name="Repression binding"
        reversible="false">
        <listOfReactants>
          <speciesReference species="Gene"/>
          <speciesReference species="P2"/>
        </listOfReactants>
        <listOfProducts>
          <speciesReference species="P2Gene"/>
        </listOfProducts>
        <kineticLaw>
          <math xmlns="http://www.w3.org/1998/Math/MathML">
            <apply>
              <times/>
              <ci> k1 </ci>
              <ci> Gene </ci>
              <ci> P2 </ci>
            </apply>
          </math>
        </kineticLaw>
      </reaction>
    </listOfReactions>
  </model>
</sbml>
```



```

</math>
<listOfParameters>
  <parameter id="k1" value="1"/>
</listOfParameters>
</kineticLaw>
</reaction>
<reaction id="ReverseRepressionBinding" name="Reverse repression binding"
  reversible="false">
  <listOfReactants>
    <speciesReference species="P2Gene"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="Gene"/>
    <speciesReference species="P2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k1r </ci>
        <ci> P2Gene </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k1r" value="10"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="Transcription" reversible="false">
  <listOfReactants>
    <speciesReference species="Gene"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="Gene"/>
    <speciesReference species="Rna"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k2 </ci>
        <ci> Gene </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k2" value="0.01"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="Translation" reversible="false">
  <listOfReactants>
    <speciesReference species="Rna"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="Rna"/>
    <speciesReference species="P"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k3 </ci>
        <ci> Rna </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k3" value="10"/>

```

```

    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="Dimerisation" reversible="false">
  <listOfReactants>
    <speciesReference species="P" stoichiometry="2"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="P2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k4 </ci>
        <cn> 0.5 </cn>
        <ci> P </ci>
        <apply>
          <minus/>
          <ci> P </ci>
          <cn type="integer"> 1 </cn>
        </apply>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k4" value="1"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="Dissociation" reversible="false">
  <listOfReactants>
    <speciesReference species="P2"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="P" stoichiometry="2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k4r </ci>
        <ci> P2 </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k4r" value="1"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="RnaDegradation" name="RNA Degradation" reversible="false">
  <listOfReactants>
    <speciesReference species="Rna"/>
  </listOfReactants>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k5 </ci>
        <ci> Rna </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k5" value="0.1"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="ProteinDegradation" name="Protein degradation">

```

```

reversible="false">
  <listOfReactants>
    <speciesReference species="P"/>
  </listOfReactants>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> k6 </ci>
        <ci> P </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k6" value="0.01"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
</listOfReactions>
</model>
</sbml>

```

A.2 Lotka-Volterra reaction system

The model below is the SBML for the stochastic version of the Lotka-Volterra system, discussed in Chapter 6.

```

<?xml version="1.0" encoding="UTF-8"?>
<sbml xmlns="http://www.sbml.org/sbml/level2" level="2" version="1">
  <model id="LotkaVolterra">
    <listOfUnitDefinitions>
      <unitDefinition id="substance">
        <listOfUnits>
          <unit kind="item" multiplier="1" offset="0"/>
        </listOfUnits>
      </unitDefinition>
    </listOfUnitDefinitions>
    <listOfCompartments>
      <compartment id="Cell"/>
    </listOfCompartments>
    <listOfSpecies>
      <species id="Prey" compartment="Cell" initialAmount="50"
        hasOnlySubstanceUnits="true"/>
      <species id="Predator" compartment="Cell" initialAmount="100"
        hasOnlySubstanceUnits="true"/>
    </listOfSpecies>
    <listOfReactions>
      <reaction id="PreyReproduction" reversible="false">
        <listOfReactants>
          <speciesReference species="Prey"/>
        </listOfReactants>
        <listOfProducts>
          <speciesReference species="Prey" stoichiometry="2"/>
        </listOfProducts>
        <kineticLaw>
          <math xmlns="http://www.w3.org/1998/Math/MathML">
            <apply>
              <times/>
              <ci> c1 </ci>
              <ci> Prey </ci>
            </apply>
          </math>
          <listOfParameters>
            <parameter id="c1" value="1"/>
          </listOfParameters>
        </kineticLaw>
      </reaction>

```

```

<reaction id="PredatorPreyInteraction" reversible="false">
  <listOfReactants>
    <speciesReference species="Prey"/>
    <speciesReference species="Predator"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="Predator" stoichiometry="2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> c2 </ci>
        <ci> Prey </ci>
        <ci> Predator </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="c2" value="0.005"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="PredatorDeath" reversible="false">
  <listOfReactants>
    <speciesReference species="Predator"/>
  </listOfReactants>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> c3 </ci>
        <ci> Predator </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="c3" value="0.6"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
</listOfReactions>
</model>
</sbml>

```

A.3 Dimerisation-kinetics model

A.3.1 Continuous deterministic version

The model below is the SBML for the continuous deterministic version of the dimerisation kinetics model, discussed in Section 7.2.

```

<?xml version="1.0" encoding="UTF-8"?>
<sbml xmlns="http://www.sbml.org/sbml/level2" level="2" version="1">
  <model id="DimerKineticsDet" name="Dimerisation Kinetics (deterministic)">
    <listOfCompartments>
      <compartment id="Cell" size="1e-15"/>
    </listOfCompartments>
    <listOfSpecies>
      <species id="P" compartment="Cell" initialConcentration="5e-07"/>
      <species id="P2" compartment="Cell" initialConcentration="0"/>
    </listOfSpecies>
    <listOfReactions>
      <reaction id="Dimerisation" reversible="false">
        <listOfReactants>
          <speciesReference species="P" stoichiometry="2"/>
        </listOfReactants>

```

```

<listOfProducts>
  <speciesReference species="P2"/>
</listOfProducts>
<kineticLaw>
  <math xmlns="http://www.w3.org/1998/Math/MathML">
    <apply>
      <times/>
      <ci> Cell </ci>
      <ci> k1 </ci>
      <ci> P </ci>
      <ci> P </ci>
    </apply>
  </math>
  <listOfParameters>
    <parameter id="k1" value="500000"/>
  </listOfParameters>
</kineticLaw>
</reaction>
<reaction id="Dissociation" reversible="false">
  <listOfReactants>
    <speciesReference species="P2"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="P" stoichiometry="2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> Cell </ci>
        <ci> k2 </ci>
        <ci> P2 </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="k2" value="0.2"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
</listOfReactions>
</model>
</sbml>

```

A.3.2 Discrete stochastic version

The model below is the SBML for the discrete stochastic version of the dimerisation kinetics model, discussed in Section 7.2.

```

<?xml version="1.0" encoding="UTF-8"?>
<sbml xmlns="http://www.sbml.org/sbml/level2" level="2" version="1">
  <model id="DimerKineticsStoch" name="Dimerisation Kinetics (stochastic)">
    <listOfUnitDefinitions>
      <unitDefinition id="substance">
        <listOfUnits>
          <unit kind="item" multiplier="1" offset="0"/>
        </listOfUnits>
      </unitDefinition>
    </listOfUnitDefinitions>
    <listOfCompartments>
      <compartment id="Cell" size="1e-15"/>
    </listOfCompartments>
    <listOfSpecies>
      <species id="P" compartment="Cell" initialAmount="301" hasOnlySubstanceUnits="true">
      <species id="P2" compartment="Cell" initialAmount="0" hasOnlySubstanceUnits="true">
    </listOfSpecies>
    <listOfReactions>

```

```

<reaction id="Dimerisation" reversible="false">
  <listOfReactants>
    <speciesReference species="P" stoichiometry="2"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="P2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <divide/>
        <apply>
          <times/>
          <ci> c1 </ci>
          <ci> P </ci>
        </apply>
        <minus/>
        <ci> P </ci>
        <cn type="integer"> 1 </cn>
      </apply>
      </apply>
      <cn type="integer"> 2 </cn>
    </math>
    <listOfParameters>
      <parameter id="c1" value="0.00166"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
<reaction id="Dissociation" reversible="false">
  <listOfReactants>
    <speciesReference species="P2"/>
  </listOfReactants>
  <listOfProducts>
    <speciesReference species="P" stoichiometry="2"/>
  </listOfProducts>
  <kineticLaw>
    <math xmlns="http://www.w3.org/1998/Math/MathML">
      <apply>
        <times/>
        <ci> c2 </ci>
        <ci> P2 </ci>
      </apply>
    </math>
    <listOfParameters>
      <parameter id="c2" value="0.2"/>
    </listOfParameters>
  </kineticLaw>
</reaction>
</listOfReactions>
</model>
</sbml>

```



References

- Abramowitz, M. & Stegun, I. A. (1984), *Pocketbook of Mathematical Functions*, Verlag Harri Deutsch, Frankfurt.
- Alfonsi, A., Cances, E., Turinici, G., Di Ventura, B. & Huisinga, W. (2005), 'Adaptive simulation of hybrid stochastic and deterministic models for biochemical systems', *ESAIM: Proceedings* **14**, 1–13.
- Allen, L. J. S. (2003), *Stochastic Processes with Applications to Biology*, Pearson Prentice Hall, Upper Saddle River, New Jersey.
- Ander, M., Beltrao, P., Di Ventura, B., Ferkinghoff-Borg, J., Foglierini, M., Kaplan, A., Lemerle, C., Tomas-Oliveira, I. & Serrano, L. (2004), 'SmartCell, a framework to simulate cellular processes that combines stochastic approximation with diffusion and localisation: analysis of simple networks', *IEE Systems Biology* **1**(1), 129–138.
- Arkin, A., Ross, J. & McAdams, H. H. (1998), 'Stochastic kinetic analysis of developmental pathway bifurcation in phage λ -infected *Escherichia coli* cells', *Genetics* **149**, 1633–1648.
- Bahcall, O. G. (2005), 'Single cell resolution in regulation of gene expression', *Molecular Systems Biology*. doi:10.1038/msb4100020.
- Beskos, A. & Roberts, G. O. (2005), 'Exact simulation of diffusions', *Annals of Applied Probability* **15**(4), 2422–2444.
- Bower, J. M. & Bolouri, H., eds (2000), *Computational Modeling of Genetic and Biochemical Networks*, MIT Press, Cambridge, Massachusetts.
- Boys, R. J., Wilkinson, D. J. & Kirkwood, T. B. L. (2004), Bayesian inference for a discretely observed stochastic kinetic model. In submission.
- Burden, R. & Faires, J. (2000), *Numerical Analysis*, Brooks Cole.
- Chib, S. (1995), 'Marginal likelihood from the Gibbs output', *Journal of the American Statistical Association* **90**(432), 1313–1321.
- Cornish-Bowden, A. (2004), *Fundamentals of Enzyme Kinetics*, third edn, Portland Press.
- Cox, D. R. & Miller, H. D. (1977), *The Theory of Stochastic Processes*, Chapman and Hall, London.
- Devroye, L. (1986), *Non-uniform Random Variate Generation*, Springer Verlag, New York.
- DuCharme, R. (1999), *XML: The Annotated Specification*, Prentice Hall PTR, Upper Saddle River, New Jersey.
- Durham, G. B. & Gallant, R. A. (2002), 'Numerical techniques for maximum likelihood estimation of continuous time diffusion processes', *Journal of Business and Economic Statistics* **20**, 279–316.
- Elf, J. & Ehrenberg, M. (2004), 'Spontaneous separation of bi-stable biochemical systems into spatial domains of opposite phases', *Systems Biology* **1**(2), 230–236.
- Gamerman, D. (1997), *Markov Chain Monte Carlo*, Texts in Statistical Science, Chapman and Hall, New York.
- Gibson, G. J. & Renshaw, E. (1998), 'Estimating parameters in stochastic compartmental models using Markov chain methods', *IMA Journal of Mathematics Applied in Medicine and Biology* **15**, 19–40.

- Gibson, M. A. & Bruck, J. (2000), 'Efficient exact stochastic simulation of chemical systems with many species and many channels', *Journal of Physical Chemistry A* **104**(9), 1876–1889.
- Gillespie, D. T. (1976), 'A general method for numerically simulating the stochastic time evolution of coupled chemical reactions', *Journal of Computational Physics* **22**, 403–434.
- Gillespie, D. T. (1977), 'Exact stochastic simulation of coupled chemical reactions', *Journal of Physical Chemistry* **81**, 2340–2361.
- Gillespie, D. T. (1992a), *Markov Processes: An Introduction for Physical Scientists*, Academic Press, New York.
- Gillespie, D. T. (1992b), 'A rigorous derivation of the chemical master equation', *Physica A* **188**, 404–425.
- Gillespie, D. T. (2000), 'The chemical Langevin equation', *Journal of Chemical Physics* **113**(1), 297–306.
- Gillespie, D. T. (2001), 'Approximate accelerated stochastic simulation of chemically reacting systems', *Journal of Chemical Physics* **115**(4), 1716–1732.
- Gillespie, D. T. & Petzold, L. R. (2003), 'Improved leap-size selection for accelerated stochastic simulation', *Journal of Chemical Physics* **119**(16), 8229–8234.
- Golightly, A. & Wilkinson, D. J. (2005a), 'Bayesian inference for stochastic kinetic models using a diffusion approximation', *Biometrics* **61**(3), 781–788.
- Golightly, A. & Wilkinson, D. J. (2005b), Bayesian sequential inference for nonlinear multivariate diffusions. In submission.
- Golightly, A. & Wilkinson, D. J. (2006), 'Bayesian sequential inference for stochastic kinetic biochemical network models', *Journal of Computational Biology*. In press.
- Golub, G. H. & Van Loan, C. F. (1996), *Matrix Computations*, third edn, Johns Hopkins University Press, Baltimore, Maryland.
- Goss, P. J. E. & Peccoud, J. (1998), 'Quantitative modeling of stochastic systems in molecular biology by using stochastic Petri nets', *Proceedings of the National Academy of Science USA* **95**, 6750–6755.
- Green, P. J. (1995), 'Reversible jump Markov chain Monte Carlo computation and Bayesian model determination', *Biometrika* **82**, 711–732.
- Hardy, S. & Robillard, P. N. (2004), 'Modeling and simulation of molecular biology systems using Petri nets: modeling goal of various approaches', *Journal of Bioinformatics and Computational Biology* **2**(4), 619–637.
- Haseltine, E. L. & Rawlings, J. B. (2002), 'Approximate simulation of coupled fast and slow reactions for stochastic chemical kinetics', *Journal of Chemical Physics* **117**(15), 6959–6969.
- Hastings, W. K. (1970), 'Monte Carlo sampling methods using Markov chains and their applications', *Biometrika* **57**, 97–109.
- Hattne, J., Fange, D. & Elf, J. (2005), 'Stochastic reaction-diffusion simulation with MesoRD', *Bioinformatics* **21**(12), 2923–2924.
- Hucka, M., Finney, A., Sauro, H. M., Bolouri, H., Doyle, J. C., Kitano, H., Arkin, A. P., Bornstein, B. J., Bray, D., Cornish-Bowden, A., Cuellar, A. A., Dronov, S., Gilles, E. D., Ginkel, M., Gor, V., Goryanin, I. I., Hedley, W. J., Hodgman, T. C., Hofmeyr, J.-H., Hunter, P. J., Juty, N. S., Kasberger, J. L., Kremling, A., Kummer, U., Novere, N. L., Loew, L. M., Luccio, D., Mendes, P., Minch, E., Mjolsness, E. D., Nakayama, Y., Nelson, M. R., Nielsen, P. F., Sakurada, T., Schaff, J. C., Shapiro, B. E., Shimizu, T. S., Spence, H. D., Stelling, J., Takahashi, K., Tomita, M., Wagner, J. & Wang, J. (2003), 'The systems biology markup language (SBML): a medium for representation and exchange of biochemical network models', *Bioinformatics* **19**(4), 524–531.

- Johnson, N. L. & Kotz, S. (1969), *Distributions in Statistics*, Vol. 1, Wiley, New York.
- Kennedy, M. C. & O'Hagan, A. (2001), 'Bayesian calibration of computer models (with discussion)', *Journal of the Royal Statistical Society, Series B* **63**, 425–464.
- Kiehl, T. R., Matheyses, R. M. & Simmons, M. K. (2004), 'Hybrid simulation of cellular behavior', *Bioinformatics* **20**(3), 316–322.
- Kirkwood, T. B. L., Boys, R. J., Gillespie, C. S., Proctor, C. J., Shanley, D. P. & Wilkinson, D. J. (2003), 'Towards an e-biology of ageing: integrating theory and data', *Nature Reviews Molecular Cell Biology* **4**(3), 243–249.
- Kitano, H., ed. (2001), *Foundations of Systems Biology*, MIT Press, Cambridge, Massachusetts.
- Klipp, E., Herwig, R., Kowald, A., Wierling, C. & Lehrach, H. (2005), *Systems Biology in Practice: Concepts, Implementation and Application*, Wiley-VCH, New York.
- Kloeden, P. E. & Platen, E. (1992), *Numerical Solution of Stochastic Differential Equations*, Springer Verlag, New York.
- Le Novere, N. & Shimizu, T. S. (2001), 'STOCHSIM: modelling of stochastic biomolecular processes', *Bioinformatics* **17**(6), 575–576.
- Lemerle, C., Di Ventura, B. & Serrano, L. (2005), 'Space as the final frontier in stochastic simulations of biological systems', *FEBS Letters* **579**(8), 1789–1794.
- Lotka, A. J. (1925), *Elements of Physical Biology*, Williams and Wilkins, Baltimore.
- McAdams, H. H. & Arkin, A. (1997), 'Stochastic mechanisms in gene expression', *Proceedings of the National Academy of Science USA* **94**, 814–819.
- McAdams, H. H. & Arkin, A. (1999), 'It's a noisy business: genetic regulation at the nanomolecular scale', *Trends in Genetics* **15**, 65–69.
- McQuarrie, D. A. (1967), 'Stochastic approach to chemical kinetics', *Journal of Applied Probability* **4**, 413–478.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. & Teller, E. (1953), 'Equations of state calculations by fast computing machines', *Journal of Chemical Physics* **21**, 1087–1092.
- Miller, I. & Miller, M. (2004), *John E. Freund's Mathematical Statistics with Applications*, seventh edn, Pearson Prentice Hall, Upper Saddle River, New Jersey.
- Morgan, B. J. T. (1984), *Elements of Simulation*, Chapman & Hall/CRC Press, London.
- Murata, T. (1989), 'Petri nets: properties, analysis and applications', *Proceedings of the IEEE* **77**(4), 541–580.
- O'Hagan, A. & Forster, J. J. (2004), *Bayesian Inference*, Vol. 2B of *Kendall's Advanced Theory of Statistics*, Arnold, London.
- Øksendal, B. (2003), *Stochastic Differential Equations: An Introduction with Applications*, sixth edn, Springer-Verlag, Heidelberg.
- Pinney, J. W., Westhead, D. R. & McConkey, G. A. (2003), 'Petri net representations in systems biology', *Biochemical Society Transactions* **31**(6), 1513–1515.
- Proctor, C. J., Söti, C., Boys, R. J., Gillespie, C. S., Shanley, D. P., Wilkinson, D. J. & Kirkwood, T. B. L. (2005), 'Modelling the action of chaperones and their role in ageing', *Mechanisms of Ageing and Development* **126**(1), 119–131.
- Puchalka, J. & Kierzek, A. M. (2004), 'Bridging the gap between stochastic and deterministic regimes in the kinetic simulations of the biochemical reaction networks', *Biophysical Journal* **86**, 1357–1372.
- R Development Core Team (2005), *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.
URL: <http://www.R-project.org>
- Rao, C. V. & Arkin, A. P. (2003), 'Stochastic chemical kinetics and the quasi-steady-state assumption: application to the Gillespie algorithm', *Journal of Chemical Physics* **118**, 4999–

- 5010.
- Reisig, W. (1985), *Petri Nets: An Introduction*, Monographs on Theoretical Computer Science, Springer-Verlag, New York.
- Rice, J. A. (1993), *Mathematical Statistics and Data Analysis*, second edn, Wadsworth, New York.
- Ripley, B. D. (1987), *Stochastic Simulation*, Wiley, New York.
- Roberts, G. O. & Stramer, O. (2001), 'On inference for non-linear diffusion models using Metropolis-Hastings algorithms', *Biometrika* **88**(3), 603–621.
- Ross, S. M. (1996), *Stochastic Processes*, Wiley, New York.
- Ross, S. M. (2003), *Introduction to Probability Models*, eighth edn, Academic Press, New York.
- Salis, H. & Kaznessis, Y. (2005a), 'Accurate hybrid stochastic simulation of a system of coupled chemical or biochemical reactions', *Journal of Chemical Physics* **122**, 054103.
- Salis, H. & Kaznessis, Y. (2005b), 'Numerical simulation of stochastic gene circuits', *Computers and Chemical Engineering* **29**, 577–588.
- Santner, T. J., Williams, B. J. & Notz, W. I. (2003), *The Design and Analysis of Computer Experiments*, Springer, New York.
- Stryer, L. (1988), *Biochemistry*, third edn, Freeman, New York.
- Van Kampen, N. G. (1992), *Stochastic Processes in Physics and Chemistry*, North-Holland.
- Volterra, V. (1926), 'Fluctuations in the abundance of a species considered mathematically', *Nature* **118**, 558–560.
- W3C (2000), W3c math home.
URL: <http://www.w3.org/Math/>
- Wilkinson, D. J. (2005), Parallel Bayesian computation, in E. J. Kontoghiorghes, ed., 'Handbook of Parallel Computing and Statistics', Marcel Dekker/CRC Press, New York, pp. 481–512.
- Wilkinson, D. J. (2006), Exact sampling of jump times in coupled Markovian jump/Langevin processes, with application to hybrid simulation of stochastic kinetic models. In preparation.
- Wilkinson, D. J. & Yeung, S. K. H. (2002), 'Conditional simulation from highly structured Gaussian systems, with application to blocking-MCMC for the Bayesian analysis of very large linear models', *Statistics and Computing* **12**, 287–300.
- Wilkinson, D. J. & Yeung, S. K. H. (2004), 'A sparse matrix approach to Bayesian computation in large linear models', *Computational Statistics and Data Analysis* **44**, 493–516.

Index

- τ -leap algorithm, 187
- lac* operon, 15

- analysis of simulation output, 165
- AR(1), 117
- associative law, 46
- auto-regulation, 14
 - stochastic, 172

- Bayes Theorem, 54
- Bayesian
 - inference, 197
 - model calibration, 237
 - network (dynamic), 237
- binomial distribution, 64
- bipartite graph, 20
- birth-death process, 3
- Box-Muller method, 96
- Brownian motion, 133

- CDF, 57
 - continuous, 71
- central limit theorem, 85
- Chapman-Kolmogorov equations, 111, 122
- Chebyshev's inequality, 74
- chemical kinetics, 139
- chemical Langevin equation, 188
- chemical master equation, 157
- chemical reactions, 6
- CLE, 188
 - inference, 230
- CLT, 85
- commutative law, 46
- compartments, 33
- complement, 46
- complete data, 220
 - likelihood, 221
- concentration, 139
- conservation law, 29

- continuous probability models, 70
- counting process, 129
- coupled reactions, 19
- cumulative distribution function, 57
 - continuous, 71

- degradation, 13
- DeMorgan's laws, 47
- density function, 70
- detailed balance, 114, 118
- diffusion
 - bridge, 233
 - inference, 230
- diffusion approximations, 134
- diffusion process, 133
- dimensionality reduction, 169
- dimer, 7
- dimerisation kinetics
 - deterministic, 142
 - stochastic, 163
- directed graph (digraph), 20
- discrete probability models, 56
- disjoint, 46
- disjoint union, 47
- distribution
 - binomial, 64
 - exponential, 77
 - gamma, 86
 - Gaussian, 82
 - geometric, 65
 - normal, 82
 - Poisson, 67
 - uniform, 75
- distributive law, 47

- eigenvalue, 110
- eigenvector, 110
- envelope method, 99
- equilibrium, 140

- constant, 142
- equilibrium distribution, 112
- Euler method, 144
- Euler-Maruyama approximation, 134
- event, 45
- exclusive, 46
- expectation, 58
 - continuous, 72
 - linear combination, 62
 - linear transformation, 60
 - product, 62
 - sum, 61
- exponential distribution, 77
 - memoryless property, 79
 - simulation, 94
- first event method, 181
- first reaction method, 182
- Fokker-Planck equation, 188
- gamma
 - distribution, 86
 - function, 86
- Gaussian distribution, 82
 - simulation, 96
- geometric distribution, 65
- Gibbs sampler, 202
 - fixed sweep, 203
 - forward-backward, 205
 - random scan, 205
- Gibson-Bruck algorithm, 182
- Gillespie algorithm, 149
- graph
 - directed, 20
 - reaction, 21
 - theory, 20
- hybrid
 - discrete/CLE, 193
 - discrete/ODE, 191
 - maximal timestep, 192
 - next reaction, 193
 - NRH, 193
 - simulation algorithms, 190
- immigration-death process, 126
- importance sampling, 92
- incidence matrix, 26
- inference
 - network structure, 234
 - parameter, 220
- inhomogeneous Poisson process, 129
- intersection, 46
- intervention analysis, 176
- invariants, 29
- inverse distribution method, 93
- Itô process, 133
- kinetics
 - chemical, 139
 - mass-action, 139
 - molecular, 145
 - stochastic, 145
- Kolmogorov equations, 123
- lac operon
 - stochastic model, 176
- law of large numbers, 73
 - strong, 75
 - weak, 75
- likelihood, 197
- linear congruential generator, 92
- LLN, 73
- lookup method, 97
- Lotka-Volterra dynamics, 141
- lower quartile, 71
- Markov
 - chain, 109
 - order, 109
 - process, 109
 - simulation, 114
- Markov chain Monte Carlo (MCMC), 202
- Markov's inequality, 74
- mass-action kinetics
 - deterministic, 139
 - stochastic, 147
- master equation, 157
- matrix, 24
 - sparse, 24
- maximal timestep method, 192
- MCMC, 202
 - Gibbs sampler, 202
 - hybrid schemes, 216
 - independence sampler, 215

- Metropolis method, 214
- Metropolis-Hastings, 212
- output analysis, 206
- mean
 - sample, 73
- median, 71
- memoryless property, 79
- Metropolis-Hastings, 212
- Michaelis-Menten kinetics, 168
- model calibration, 237
- model reduction, 169
- modelling, 1
- mole, 139
- molecular kinetics, 145
- Monte-Carlo, 91
- multiplication principle, 51

- network inference, 234
- next reaction method, 182
- normal distribution, 82
 - simulation, 96
- NRH, 193
- null space, 29
- numerical integration, 144

- ODEs, 140
 - Euler method, 144
 - numerical integration, 144
 - Runge-Kutta method, 144
- ordinary differential equations (ODEs), 140

- P-invariant, 29
- parameter inference, 220
- PDF, 70
 - 1-1 transformation, 73
 - linear transformation, 72
- Petri nets, 21
 - formalism, 24
 - incidence matrix, 26
 - invariants, 29
 - reachability, 30
- photosynthesis, 1
- PMF, 57
- point process, 129
- Poisson distribution, 67
- Poisson process, 69
 - inhomogeneous, 129
 - simulation, 101
- Poisson timestep method, 186
- population modelling, 16
- probability, 45
 - axioms, 47
 - Bayes Theorem, 54
 - Bayesian, 49
 - classical, 49
 - conditional, 52
 - density function, 70
 - discrete, 56
 - distribution function, 57
 - independence, 53
 - interpretations, 48
 - mass function, 57
 - subjective, 49
- pseudo-random number generators, 92

- quartile, 71

- R statistical programming language, 101
- random number generators, 92
- rate constants
 - conversion, 153
- rate law
 - deterministic, 139
- reachability, 30
- reaction, 7
 - graphs, 21
 - matrix, 26
 - reversible, 141
- rejection sampling, 98
- repression, 11
- reversible Markov process, 113
- reversible reaction, 7, 141
- Runge-Kutta method, 144

- S-invariant, 29
- sample
 - mean, 73
 - space, 45
 - variance, 103
- SBML, 31
 - compartments, 33
 - events, 176
 - full example model, 36
 - kinetic law, 35
 - Level 1, 32
 - Level 2, 32

- parameter, 34
- reaction, 35
- shorthand, 36
- species, 34
- units, 33
- SDE, 133
- sensitivity analysis, 173
- set theory, 46
- SLLN, 75
- spatial stochastic models, 195
- SPN, 150
- stationary distribution, 111
- stochastic
 - kinetics, 145
 - matrix, 110
 - process, 109
 - rate constant conversion, 153
 - rate constants, 147
- stochastic differential equation, 133
- stochastic Petri net (SPN), 150
- stochastic simulation, 91
- stoichiometry, 7
 - matrix, 26
- summarising simulation output, 165
- Systems Biology Markup Language (SBML), 31

- T-invariant, 29
- tau-leap algorithm, 187
- theorem of total probability
 - continuous, 72
- time-discretisation, 186
- time-varying volume, 185
- transcription, 8
- transformation methods, 93
- transition
 - kernel, 116
 - matrix, 110
- translation, 12
- transport (nuclear), 13

- uncertainty analysis, 173
- uniform distribution, 75
 - discrete, 64
 - simulation, 92
- union, 46
- units
 - SBML, 33

- upper quartile, 71

- variance, 59
 - continuous, 72
 - linear combination, 63
 - linear transformation, 62
 - sample, 103
 - sum, 62
- visualisation, 165

- Wiener process, 133
- WLLN, 75

