Graduate Texts in Mathematics



Daniel W. Stroock

An Introduction to Markov Processes

Second Edition



Series Editors:

Sheldon Axler San Francisco State University, San Francisco, CA, USA

Kenneth Ribet University of California, Berkeley, CA, USA

Advisory Board:

Colin Adams, Williams College, Williamstown, MA, USA Alejandro Adem, University of British Columbia, Vancouver, BC, Canada Ruth Charney, Brandeis University, Waltham, MA, USA Irene M. Gamba, The University of Texas at Austin, Austin, TX, USA Roger E. Howe, Yale University, New Haven, CT, USA David Jerison, Massachusetts Institute of Technology, Cambridge, MA, USA Jeffrey C. Lagarias, University of Michigan, Ann Arbor, MI, USA Jill Pipher, Brown University, Providence, RI, USA Fadil Santosa, University of Minnesota, Minneapolis, MN, USA Amie Wilkinson, University of Chicago, Chicago, IL, USA

Graduate Texts in Mathematics bridge the gap between passive study and creative understanding, offering graduate-level introductions to advanced topics in mathematics. The volumes are carefully written as teaching aids and highlight characteristic features of the theory. Although these books are frequently used as textbooks in graduate courses, they are also suitable for individual study.

For further volumes: http://www.springer.com/series/136 Daniel W. Stroock

An Introduction to Markov Processes

Second Edition



Daniel W. Stroock Department of Mathematics Massachusetts Institute of Technology Cambridge, MA, USA

ISSN 0072-5285 ISSN 2197-5612 (electronic) Graduate Texts in Mathematics ISBN 978-3-642-40522-8 ISBN 978-3-642-40523-5 (eBook) DOI 10.1007/978-3-642-40523-5 Springer Heidelberg New York Dordrecht London

Library of Congress Control Number: 2013952463

Mathematics Subject Classification (2010): 60-01, 60J10, 60J27, 60J28, 37L40

© Springer-Verlag Berlin Heidelberg 2005, 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

This book is dedicated to my longtime colleague:

Richard A. Holley

Preface

To some extent, it would be accurate to summarize the contents of this book as an intolerably protracted description of what happens when either one raises a transition probability matrix \mathbf{P} (i.e., all entries $(\mathbf{P})_{ij}$ are non-negative and each row of \mathbf{P} sums to 1) to higher and higher powers or one exponentiates $\mathbf{R}(\mathbf{P} - \mathbf{I})$, where \mathbf{R} is a diagonal matrix with non-negative entries. Indeed, when it comes right down to it, that is all that is done in this book. However, I, and others of my ilk, would take offense at such a dismissive characterization of the theory of Markov chains and processes with values in a countable state space, and a primary goal of mine in writing this book was to convince its readers that our offense would be warranted.

The reason why I, and others of my persuasion, refuse to consider the theory here as no more than a subset of matrix theory is that to do so is to ignore the pervasive role that probability plays throughout. Namely, probability theory provides a model which both motivates and provides a context for what we are doing with these matrices. To wit, even the term "transition probability matrix" lends meaning to an otherwise rather peculiar set of hypotheses to make about a matrix. Specifically, it suggests that we think of the matrix entry $(\mathbf{P})_{ii}$ as giving the probability that, in one step, a system in state i will make a transition to state j. Moreover, if we adopt this interpretation for $(\mathbf{P})_{ii}$, then we must interpret the entry $(\mathbf{P}^n)_{ii}$ of \mathbf{P}^n as the probability of the same transition in n steps. Thus, as $n \to \infty$, \mathbf{P}^n is encoding the long time behavior of a randomly evolving system for which P encodes the one-step behavior, and, as we will see, this interpretation will guide us to an understanding of $\lim_{n\to\infty} (\mathbf{P}^n)_{ii}$. In addition, and perhaps even more important, is the role that probability plays in bridging the chasm between mathematics and the rest of the world. Indeed, it is the probabilistic metaphor which allows one to formulate mathematical models of various phenomena observed in both the natural and social sciences. Without the language of probability, it is hard to imagine how one would go about connecting such phenomena to \mathbf{P}^n .

In spite of the propaganda at the end of the preceding paragraph, this book is written from a mathematician's perspective. Thus, for the most part, the probabilistic metaphor will be used to elucidate mathematical concepts rather than to provide mathematical explanations for non-mathematical phenomena. There are two reasons for my having chosen this perspective. First, and foremost, is my own background. Although I have occasionally tried to help people who are engaged in various sorts of applications, I have not accumulated a large store of examples which are easily translated into terms which are appropriate for a book at this level. In fact, my experience has taught me that people engaged in applications are more than competent to handle the routine problems that they encounter, and that they come to someone like me only as a last resort. As a consequence, the questions which they ask me tend to be quite difficult and the answers to those few which I can solve usually involve material which is well beyond the scope of the present book. The second reason for my writing this book in the way that I have is that I think the material itself is of sufficient interest to stand on its own. In spite of what funding agencies would have us believe, mathematics *qua* mathematics is a worthy intellectual endeavor, and I think there is a place for a modern introduction to stochastic processes which is unabashed about making mathematics its top priority.

I came to this opinion after several semesters during which I taught the introduction to stochastic processes course offered by the M.I.T. department of mathematics. The clientele for that course has been an interesting mix of undergraduate and graduate students, less than half of whom concentrate in mathematics. Nonetheless, most of the students who stay with the course have considerable talent and appreciation for mathematics, even though they lack the formal mathematical training which is requisite for a modern course in stochastic processes, at least as such courses are now taught in mathematics departments to their own graduate students. As a result, I found no ready-made choice of text for the course. On the one hand, the most obvious choice is the classic text A First Course in Stochastic Processes, either the original one by S. Karlin or the updated version [4] by S. Karlin and H. Taylor. Their book gives a no nonsense introduction to stochastic processes, especially Markov processes, on a countable state space, and its consistently honest, if not always easily assimilated, presentation of proofs is complemented by a daunting number of examples and exercises. On the other hand, when I began, I feared that adopting Karlin and Taylor for my course would be a mistake of the same sort as adopting Feller's book for an undergraduate introduction to probability, and this fear prevailed the first two times I taught the course. However, after using, and finding wanting, two derivatives of Karlin's classic, I took the plunge and assigned Karlin and Taylor's book. The result was very much the one which I predicted: I was far more enthusiastic about the text than were my students.

In an attempt to make Karlin and Taylor's book more palatable for the students, I started supplementing their text with notes in which I tried to couch the proofs in terms which I hoped they would find more accessible, and my efforts were rewarded with a quite positive response. In fact, as my notes became more and more extensive and began to diminish the importance of the book, I decided to convert them into what is now this book, although I realize that my decision to do so may have been stupid. For one thing, the market is already close to glutted with books that purport to cover this material. Moreover, some of these books are quite popular, although my experience with them leads me to believe that their popularity is not always correlated with the quality of the mathematics they contain. Having made that pejorative comment, I will not make public which are the books that led me to this conclusion. Instead, I will only mention the books on this topic, besides Karlin and Taylor's, which I very much like. J. Norris's book [5] is an excellent introduction to Markov processes which, at the same time, provides its readers with a good place to exercise their measure-theoretic skills. Of course, Norris's book is only appropriate for students who have measure-theoretic skills to exercise. On the other hand, for students who possess those skills, his book is a place where they can see measure theory put to work in an attractive way. In addition, Norris has included many interesting examples and exercises which illustrate how the subject can be applied. For more advanced students, an excellent treatment of Markov chains on a general state space can be found in the book [6] by D. Revuz.

The present book includes most of the mathematical material contained in [5], but the proofs here demand much less measure theory than his do. In fact, although I have systematically employed measure theoretic terminology (Lebesgue's dominated convergence theorem, the monotone convergence theorem, etc.), which is explained in Chap. 7, I have done so only to familiarize my readers with the jargon that they will encounter if they delve more deeply into the subject. In fact, because the state spaces in this book are countable, the applications which I have made of Lebesgue's theory are, with one notable exception, entirely trivial. The one exception is that I need to know the existence of countably infinite families of mutually independent random variables. In Sect. 7.2 I discuss how one goes about proving their existence, but, as distinguished from the first edition of this text, I do not go into details and instead refer to the treatment in [8]. Be that as it may, the reader who is ready to accept that such families exist has little need to consult Chap. 7 except for terminology and the derivation of a few essentially obvious facts about series.

The organization of this book should be more or less self-evident from the table of contents. In Chap. 1, I give a bare hands treatment of the basic facts, with particular emphasis on recurrence and transience, about nearest neighbor random walks on the square, d-dimensional lattice \mathbb{Z}^d . Chapter 2 introduces the study of ergodic properties, and this becomes the central theme which ties together Chaps. 2 through 6. In Chap. 2, the stochastic processes under consideration are Markov chains (i.e., the time parameter is discrete), and the driving force behind the development there is an idea which was introduced by Doeblin. Restricted as the applicability of Doeblin's idea may be, it has the enormous advantage over the material in Chaps. 4 and 5 that it provides an estimate on the rate at which the chain is converging to its equilibrium distribution. Chapter 3 begins with the classification of states in terms of recurrence and transience and then introduces some computational techniques for computing stationary probabilities. As an application, the final section of Chap. 3 gives a proof that Wilson's algorithm works and that it can be used to derive Kirchhoff's matrix tree theorem. The contents of this section are based on ideas that I learned from S. Sternberg, who learned them from M. Kozdron, who in turn learned them from G. Lawler. It was Kozdron who had the idea of using Wilson's algorithm to derive Kirchhoff's theorem.

In Chap. 4, I study the ergodic properties of Markov chains that do not necessarily satisfy Doeblin's condition. The main result here is the one summarized in equation

(4.1.15). Even though it is completely elementary, the derivation of (4.1.15), is, without doubt, the most demanding piece of analysis in the entire book. So far as I know, every proof of (4.1.15) requires work at some stage. In supposedly "simpler" proofs, the work is hidden elsewhere (either measure theory, as in [5] and [6], or in operator theory, as in [2]). The treatment given here, which is a re-working of the one in [4] based on Feller's renewal theorem, demands nothing more of the reader than a thorough understanding of arguments involving limits superior, limits inferior, and their role in proving that limits exist. In Chap. 5, Markov chains are replaced by continuous-time Markov processes (still on a countable state space). I do this first in the case when the rates are bounded and therefore problems of possible explosion do not arise. Afterwards, I allow for unbounded rates and develop criteria, besides boundedness, which guarantee non-explosion. The remainder of the chapter is devoted to transferring the results obtained for Markov chains in Chaps. 2 and 4 to the continuous-time setting.

Aside from Chap. 7, which is more like an appendix than an integral part of the book, the book ends with Chap. 6. The goal in Chap. 6 is to obtain quantitative results, reminiscent of, if not as strong as, those in Chap. 2, when Doeblin's theory either fails entirely or yields rather poor estimates. The new ingredient in Chap. 6 is the assumption that the chain or process is reversible (i.e., the transition probability is self-adjoint in the L^2 -space of its stationary distribution), and the engine which makes everything go is the associated Dirichlet form. In the final section, the power of the Dirichlet form methodology is tested in an analysis of the Metropolis (a.k.a. as simulated annealing) algorithm. Finally, as I said before, Chap. 7 is an appendix in which the ideas and terminology of Lebesgue's theory of measure and integration are reviewed. Sect. 7.2.1.

I have finally reached the traditional place reserved for thanking those individuals who, either directly or indirectly, contributed to this book. The principal direct contributors are the many students who suffered with various and spontaneously changing versions of this book. I am particularly grateful to Adela Popescu whose careful reading of the first edition brought to light many minor and a few major errors that have been removed and, undoubtedly, replaced by new ones. In addition, I am grateful to Sternberg and Kozdron for introducing me to the ideas in Sect. 3.3.

Thanking, or even identifying, the indirect contributors is trickier. Indeed, they include all the individuals, both dead and alive, from whom I received my education, and I am not about to bore you with even a partial list of who they were or are. Nonetheless, there is one person who, over a period of more than ten years, patiently taught me to appreciate the sort of material treated here. Namely, Richard A. Holley, to whom I have dedicated this book, is a *true probabilist*. To wit, for Dick, intuitive understanding usually precedes his mathematically rigorous comprehension of a probabilistic phenomenon. This statement should lead no one to to doubt Dick's powers as a rigorous mathematician. On the contrary, his intuitive grasp of probability theory not only enhances his own formidable mathematical powers, it has saved me and others from blindly pursuing flawed lines of reasoning. As all who have worked with him know, reconsider what you are saying if ever, during

Preface

some diatribe into which you have launched, Dick quietly says "I don't follow that." In addition to his mathematical prowess, every one of Dick's many students will attest to his wonderful generosity. I was not his student, but I was his colleague, and I can assure you that his generosity is not limited to his students.

Cambridge, MA, USA August 2013 Daniel W. Stroock

Contents

1	Ran	dom W	Valks, a Good Place to Begin	1
	1.1	Neare	st Neighbor Random Walks on \mathbb{Z}	1
		1.1.1	Distribution at Time n	2
		1.1.2	Passage Times via the Reflection Principle	3
		1.1.3	Some Related Computations	4
		1.1.4	Time of First Return	7
		1.1.5	Passage Times via Functional Equations	8
	1.2	Recur	rence Properties of Random Walks	9
		1.2.1	Random Walks on \mathbb{Z}^d	9
		1.2.2	An Elementary Recurrence Criterion	10
		1.2.3	Recurrence of Symmetric Random Walk in \mathbb{Z}^2	12
		1.2.4	Transience in \mathbb{Z}^3	14
	1.3	Exerci	ises	17
2	Doe	blin's T	Theory for Markov Chains	25
	2.1		Generalities	25
		2.1.1	Existence of Markov Chains	26
		2.1.2	Transition Probabilities & Probability Vectors	27
		2.1.3	Transition Probabilities and Functions	28
		2.1.4	The Markov Property	30
	2.2			30
		2.2.1	Doeblin's Basic Theorem	30
		2.2.2	A Couple of Extensions	33
	2.3	Eleme	ents of Ergodic Theory	35
		2.3.1	The Mean Ergodic Theorem	36
		2.3.2	Return Times	37
		2.3.3	Identification of π	41
	2.4	Exerci	ises	43
3	Stat	ionarv	Probabilities	49
	3.1	•	fication of States	49
		3.1.1		50

		3.1.2	Criteria for Recurrence and Transience	52
		3.1.3	Periodicity	56
	3.2	Compu	tation of Stationary Probabilities	58
		3.2.1	Preliminary Results	58
		3.2.2	Computations via Linear Algebra	59
	3.3	Wilson	's Algorithm and Kirchhoff's Formula	64
		3.3.1	Spanning Trees and Wilson Runs	64
		3.3.2	Wilson's Algorithm	65
			Kirchhoff's Matrix Tree Theorem	
	3.4	Exercis	es	69
4	Mor	e About	the Ergodic Properties of Markov Chains	73
	4.1	Ergodic	Theory Without Doeblin	74
			Convergence of Matrices	
		4.1.2	Abel Convergence	75
			Structure of Stationary Distributions	
		4.1.4	A Digression About Moments of Return Times	80
		4.1.5	A Small Improvement	82
		4.1.6	The Mean Ergodic Theorem Again	84
		4.1.7	A Refinement in the Aperiodic Case	85
		4.1.8	Periodic Structure	90
	4.2	Exercis	es	91
5	Mar	kov Pro	cesses in Continuous Time	99
	5 1	р.	Deserves	00
	5.1	Poisson		99
	5.1		Processes	
	5.1	5.1.1	The Simple Poisson Process on \mathbb{Z}^d	99
	5.1	5.1.1 5.1.2	The Simple Poisson Process $\ldots \ldots $	99 102
		5.1.1 5.1.2 Markov	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on $\mathbb{Z}^d \dots \dots \dots \dots \dots \dots$ Processes with Bounded Rates $\dots \dots \dots \dots \dots \dots \dots \dots \dots$	99 102 104
		5.1.1 5.1.2 Markov 5.2.1	The Simple Poisson Process	99 102 104 105
		5.1.1 5.1.2 Markov 5.2.1 5.2.2	The Simple Poisson Process	99 102 104 105 108
		5.1.1 5.1.2 Markov 5.2.1 5.2.2 5.2.3	The Simple Poisson Process	99 102 104 105 108 111
		5.1.1 5.1.2 Markov 5.2.1 5.2.2 5.2.3 5.2.4	The Simple Poisson Process	99 102 104 105 108 111 112
	5.2	5.1.1 5.1.2 Markov 5.2.1 5.2.2 5.2.3 5.2.4 Unbour	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d $\dots \dots \dots$	99 102 104 105 108 111 112 114
	5.2	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d $\dots \dots \dots$	99 102 104 105 108 111 112 114 114
	5.2	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d $\dots \dots \dots \dots \dots \dots \dots \dots$ Processes with Bounded Rates $\dots \dots \dots$	99 102 104 105 108 111 112 114 114 120
	5.2	5.1.1 5.1.2 Markov 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d	99 102 104 105 108 111 112 114 114 120 122
	5.2	5.1.1 5.1.2 Markov 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d $\dots \dots \dots \dots \dots \dots \dots \dots$ Processes with Bounded Rates $\dots \dots \dots$ Basic Construction $\dots \dots \dots$	99 102 104 105 108 111 112 114 114 120 122 122
	5.2	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic 5.4.1	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d $\dots \dots \dots \dots \dots \dots \dots \dots \dots$ Processes with Bounded Rates $\dots \dots \dots$	99 102 104 105 108 111 112 114 114 120 122 122 123
	5.2	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic 5.4.1 5.4.2	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d	99 102 104 105 108 111 112 114 114 120 122 122 123 126
	5.2	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic 5.4.1 5.4.2 5.4.3	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d	99 102 104 105 108 111 112 114 114 120 122 122 123 126 129
6	5.25.35.45.5	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic 5.4.1 5.4.2 5.4.3 Exercis	The Simple Poisson Process $\dots \dots \dots \dots \dots \dots \dots$ Compound Poisson Processes on \mathbb{Z}^d	 99 102 104 105 108 111 112 114 114 120 122 123 126 129 130
6	5.25.35.45.5	5.1.1 5.1.2 Markow 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic 5.4.1 5.4.2 5.4.3 Exercis	The Simple Poisson Process	99 102 104 105 108 111 112 114 114 120 122 123 126 129 130 137
6	 5.2 5.3 5.4 5.5 Rev 	5.1.1 5.1.2 Markov 5.2.1 5.2.2 5.2.3 5.2.4 Unbour 5.3.1 5.3.2 5.3.3 Ergodic 5.4.1 5.4.2 5.4.3 Exercis ersible M Reversi	The Simple Poisson Process \dots \mathbb{Z}^d \dots \mathbb{Z}^d Processes with Bounded Rates \dots Basic Construction \dots Basic Construction \dots	99 102 104 105 108 111 112 114 114 120 122 123 126 129 130 137 138

		6.1.3	The Spectral Gap	141
		6.1.4	Reversibility and Periodicity	143
		6.1.5	Relation to Convergence in Variation	144
	6.2	Dirich	let Forms and Estimation of β	145
		6.2.1	The Dirichlet Form and Poincaré's Inequality	146
		6.2.2	Estimating β_+	148
		6.2.3	Estimating β_{-}	150
	6.3	Revers	sible Markov Processes in Continuous Time	151
		6.3.1	Criterion for Reversibility	151
		6.3.2	Convergence in $L^2(\hat{\pi})$ for Bounded Rates	152
		6.3.3	$L^2(\hat{\pi})$ -Convergence Rate in General	154
		6.3.4	Estimating λ	157
	6.4	Gibbs	States and Glauber Dynamics	157
		6.4.1	Formulation	158
		6.4.2	The Dirichlet Form	159
	6.5	Simul	ated Annealing	162
		6.5.1	The Algorithm	163
		6.5.2	Construction of the Transition Probabilities	164
		6.5.3	Description of the Markov Process	166
		6.5.4	Choosing a Cooling Schedule	166
		6.5.5	Small Improvements	169
	6.6		ses	170
7		Exerci	ses	
7		Exerci Iinimal	ses	170
7	A M	Exerci Iinimal	ses	170 179
7	A M	Exerci linimal A Des	ses Introduction to Measure Theory cription of Lebesgue's Measure Theory	170 179 179
7	A M	Exerci linimal A Des 7.1.1	ses	170 179 179 179
7	A M	Exerci linimal A Des 7.1.1 7.1.2	ses	170 179 179 179 181
7	A M	Exerci linimal A Des 7.1.1 7.1.2 7.1.3	ses Introduction to Measure Theory	170 179 179 179 181 182
7	A M	Exerci linimal A Des 7.1.1 7.1.2 7.1.3 7.1.4	ses	170 179 179 179 181 182 183
7	A M	Exerci Iinimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5	ses Introduction to Measure Theory	170 179 179 179 181 182 183 184
7	A M	Exerci Iinimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6	ses Introduction to Measure Theory	170 179 179 179 181 182 183 184 186
7	A M	Exerci Iinimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8	ses Introduction to Measure Theory	170 179 179 179 181 182 183 184 186 188
7	A M 7.1	Exerci Iinimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8	ses Introduction to Measure Theory	170 179 179 181 182 183 184 186 188 190
7	A M 7.1	Exerci linimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8 Model 7.2.1	ses Introduction to Measure Theory	170 179 179 181 182 183 184 186 188 190 192
7	A M 7.1 7.2	Exerci linimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8 Model 7.2.1	ses Introduction to Measure Theory	170 179 179 181 182 183 184 186 188 190 192 193
7	A M 7.1 7.2	Exerci linimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8 Model 7.2.1 Indeped 7.3.1	ses Introduction to Measure Theory	170 179 179 181 182 183 184 186 188 190 192 193 194 194 196
7	A M 7.1 7.2 7.3	Exerci linimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8 Model 7.2.1 Indeped 7.3.1	ses Introduction to Measure Theory	170 179 179 181 182 183 184 186 188 190 192 193 194 194
-	A M 7.1 7.2 7.3	Exerci linimal A Des 7.1.1 7.1.2 7.1.3 7.1.4 7.1.5 7.1.6 7.1.7 7.1.8 Model 7.2.1 Indepe 7.3.1 Condi 7.4.1	ses Introduction to Measure Theory	170 179 179 181 182 183 184 186 188 190 192 193 194 194 196

List of Notations

$\mathbb{Z} \ \& \ \mathbb{Z}^+$	set of all integers and the subset of positive integers
\mathbb{N}	set of non-negative integers
$\lfloor t \rfloor$	integer part max{ $n \in \mathbb{Z} : n \leq t$ } of $t \in \mathbb{R}$
#S	number of elements in the set S
AC	complement of the set A
1_A	indicator function of the set $A: 1_A(x) = 1$ if $x \in A$ and
	$1_A(x) = 0 \text{ if } x \notin A$
$F \upharpoonright S$	restriction of the function F to the set S
$a \wedge b \& a \vee b$	minimum and the maximum of $a, b \in \mathbb{R}$
$a^{+} \& a^{-}$	positive part $a \lor 0$ & negative part $(-a) \lor 0$ of $a \in \mathbb{R}$
$i \rightarrow j \& i \leftrightarrow j$	state j is accessible from i & communicates with i . See
	Sect. 3.1
$\delta_{i,j}$	Kronecker delta: $\delta_{i,j}$ is 1 or 0 depending on whether <i>i</i> is
	equal or unequal to j
$\boldsymbol{\delta}_i$	point measure at $i \in \mathbb{S}$: $(\delta_i)_j = \delta_{i,j}$ for $j \in \mathbb{S}$
$\mathbb{E}[X, A]$	expected value of X on the event A . See Sect. 6.2
$\mathbb{E}[X \mid A] \& \mathbb{E}[X \mid \Sigma]$	conditional expectatation value of <i>X</i> given the event <i>A</i> &
	the σ -algebra Σ . See Sect. 6.4
$\langle \varphi \rangle_{\pi}$	alternative notations for $\sum_{i \in \mathbb{S}} \varphi(i)(\pi)_i$. See (5.1.4)
$\langle \varphi, \psi \rangle_{\pi} \& \ \varphi\ _{2,\pi}$	inner product and norm in $L^2(\pi)$. See Sect. 5.1.2
Stat(P)	set of stationary distribution for the transition probability
	matrix P . See Sect. 4.1.3
$\ \boldsymbol{\mu}\ _{\mathrm{v}}$	variation norm of the row vector μ . See (2.1.5)
$\ f\ _{\mathrm{u}}$	uniform norm of the function f . See (2.1.10)
$\ \mathbf{M}\ _{u,v}$	uniform-variation norm of the matrix \mathbf{M} . See (4.1.1)
$\operatorname{Var}_{\boldsymbol{\mu}}(f)$	variance of f relative to the probability vector μ

Chapter 1 Random Walks, a Good Place to Begin

The purpose of this chapter is to discuss some examples of Markov processes that can be understood even before the term "Markov process" is. Indeed, anyone who has been introduced to probability theory will recognize that these processes all derive from consideration of elementary "coin tossing."

1.1 Nearest Neighbor Random Walks on \mathbb{Z}

Let *p* be a fixed number from the open interval (0, 1), and suppose that¹ $\{B_n : n \in \mathbb{Z}^+\}$ is a sequence of $\{-1, 1\}$ -valued, identically distributed *Bernoulli random variables*² which are 1 with probability *p*. That is, for any $n \in \mathbb{Z}^+$ and any $E \equiv (\epsilon_1, \ldots, \epsilon_n) \in \{-1, 1\}^n$,

$$\mathbb{P}(B_1 = \epsilon_1, \dots, B_n = \epsilon_n) = p^{N(E)} q^{n-N(E)} \quad \text{where } q \equiv 1-p \quad \text{and}$$
$$N(E) \equiv \#\{m : \epsilon_m = 1\} = \frac{n+S_n(E)}{2} \quad \text{when } S_n(E) \equiv \sum_{1}^{n} \epsilon_m. \tag{1.1.1}$$

Next, set

$$X_0 = 0$$
 and $X_n = \sum_{m=1}^n B_m$ for $n \in \mathbb{Z}^+$. (1.1.2)

The above family of random variables $\{X_n : n \in \mathbb{N}\}$ is often called a *nearest* neighbor random walk on \mathbb{Z} . Nearest neighbor random walks are examples of

 $^{{}^1\}mathbb{Z}$ is used to denote the set of all integers, of which \mathbb{N} and \mathbb{Z}^+ are, respectively, the non-negative and positive members.

 $^{^{2}}$ For historical reasons, mutually independent random variables which take only two values are often said to be Bernoulli random variables. A mathematically rigorous proof that infinitely many exist can be found in Sect. 2.2 of [8].

Markov processes, but the description that I have just given is the one which would be given in elementary probability theory, as opposed to a course, like this one, devoted to stochastic processes. When studying stochastic processes the description should emphasize the dynamic nature of the family. Thus, a stochastic process oriented description might replace (1.1.2) by

$$\mathbb{P}(X_0 = 0) = 1 \text{ and}$$

$$\mathbb{P}(X_n - X_{n-1} = \epsilon \mid X_0, \dots, X_{n-1}) = \begin{cases} p & \text{if } \epsilon = 1 \\ q & \text{if } \epsilon = -1, \end{cases}$$
(1.1.3)

where $\mathbb{P}(X_n - X_{n-1} = \epsilon \mid X_0, ..., X_{n-1})$ denotes the *conditional probability* (cf. Sect. 7.4.1) that $X_n - X_{n-1} = \epsilon$ given $\sigma(\{X_0, ..., X_{n-1}\})$. Certainly (1.1.3) is more dynamic a description than the one in (1.1.2). Specifically, it says that the process starts from 0 at time n = 0 and proceeds so that, at each time $n \in \mathbb{Z}^+$, it moves one step forward with probability p or one step backward with probability q, independent of where it has been before time n.

1.1.1 Distribution at Time n

In this subsection I will present two approaches to computing $\mathbb{P}(X_n = m)$. The first computation is based on the description given in (1.1.2). From (1.1.2) it is clear that $\mathbb{P}(|X_n| \le n) = 1$. In addition, it is obvious that

 $n \text{ odd} \implies \mathbb{P}(X_n \text{ is odd}) = 1 \text{ and } n \text{ even} \implies \mathbb{P}(X_n \text{ is even}) = 1.$

Finally, given $m \in \{-n, ..., n\}$ with the same parity as n and a string $E = (\epsilon_1, ..., \epsilon_n) \in \{-1, 1\}^n$ with (cf. (1.1.1)) $S_n(E) = m$, $N(E) = \frac{n+m}{2}$ and so

$$\mathbb{P}(B_1 = \epsilon_1, \ldots, B_n = \epsilon_n) = p^{\frac{n+m}{2}} q^{\frac{n-m}{2}}.$$

Hence, because, when $\binom{\ell}{k} \equiv \frac{\ell!}{k!(\ell-k)!}$ is the *binomial coefficient* " ℓ choose k," there are $\binom{n}{m+n}$ such strings E, we see that

$$\mathbb{P}(X_n = m) = \binom{n}{\frac{m+m}{2}} p^{\frac{n+m}{2}} q^{\frac{n-m}{2}}$$
(1.1.4)

if $m \in \mathbb{Z}$, $|m| \le n$, and *m* has the same parity as *n* and is 0 otherwise.

Our second computation of the same probability will be based on the more dynamic description given in (1.1.3). To do this, we introduce the notation $(P^n)_m \equiv \mathbb{P}(X_n = m)$. Obviously, $(P^0)_m = \delta_{0,m}$, where $\delta_{k,\ell}$ is the *Kronecker symbol* which is 1 when $k = \ell$ and 0 otherwise. Further, from (1.1.3), we see that $\mathbb{P}(X_n = m)$ equals

$$\mathbb{P}(X_{n-1} = m - 1 \& X_n = m) + \mathbb{P}(X_{n-1} = m + 1 \& X_n = m)$$
$$= p\mathbb{P}(X_{n-1} = m - 1) + q\mathbb{P}(X_{n-1} = m + 1).$$

That is,

$$(P^0)_m = \delta_{0,m}$$
 and $(P^n)_m = p(P^{n-1})_{m-1} + q(P^{n-1})_{m+1}$. (1.1.5)

Obviously, (1.1.5) provides a complete, albeit implicit, prescription for computing the numbers $(P^n)_m$, and one can easily check that the numbers given by (1.1.4) satisfy this prescription. Alternatively, one can use (1.1.5) plus induction on *n* to see that $(P^n)_m = 0$ unless $m = 2\ell - n$ for some $0 \le \ell \le n$ and that $(C^n)_0 = (C^n)_n = 1$ and $(C^n)_\ell = (C^{n-1})_{\ell-1} + (C^{n-1})_\ell$ for $1 \le \ell < n$ where $(C^n)_\ell \equiv p^{-\ell}q^{n-\ell}(P^n)_{2\ell-n}$. In other words, the family $\{(C^n)_\ell : n \in \mathbb{N} \& 0 \le \ell \le n\}$ are given by Pascal's triangle and are therefore the binomial coefficients.

1.1.2 Passage Times via the Reflection Principle

More challenging than the computation in Sect. 1.1.1 is that of finding the distribution of the first passage time to a point $a \in \mathbb{Z}$. That is, given $a \in \mathbb{Z} \setminus \{0\}$, set³

$$\zeta^{\{a\}} = \inf\{n \ge 1 : X_n = a\} \quad (\equiv \infty \text{ when } X_n \neq a \text{ for any } n \ge 1). \tag{1.1.6}$$

Then $\zeta^{\{a\}}$ is the *first passage time* to *a*, and our goal here is to find its distribution. Equivalently, we want an expression for $\mathbb{P}(\zeta^{\{a\}} = n)$, and clearly, by the considerations in Sect. 1.1.1, we need only worry about *n*'s which satisfy $n \ge |a|$ and have the same parity as *a*.

Again I will present two approaches to this problem, here based on (1.1.2) and in Sect. 1.1.5 on (1.1.3). To carry out the one based on (1.1.2), assume that $a \in \mathbb{Z}^+$, suppose that $n \in \mathbb{Z}^+$ has the same parity as a, and observe first that

$$\mathbb{P}(\zeta^{\{a\}} = n) = \mathbb{P}(X_n = a \& \zeta^{\{a\}} > n-1) = p\mathbb{P}(\zeta^{\{a\}} > n-1 \& X_{n-1} = a-1).$$

Hence, it suffices for us to compute $\mathbb{P}(\zeta_a > n - 1 \& X_{n-1} = a - 1)$. For this purpose, note that for any $E \in \{-1, 1\}^{n-1}$ with $S_{n-1}(E) = a - 1$, the event $\{(B_1, \ldots, B_{n-1}) = E\}$ has probability $p^{\frac{n+a}{2}-1}q^{\frac{n-a}{2}}$. Thus,

$$\mathbb{P}(\zeta^{\{a\}} = n) = \mathcal{N}(n, a) p^{\frac{n+a}{2}} q^{\frac{n-a}{2}}$$
(*)

where $\mathcal{N}(n, a)$ is the number of $E \in \{-1, 1\}^{n-1}$ with the properties that $S_m(E) \leq a - 1$ for $0 \leq m \leq n - 1$ and $S_{n-1}(E) = a - 1$. That is, everything comes down to the computation of $\mathcal{N}(n, a)$. Alternatively, since $\mathcal{N}(n, a) = \left(\frac{n-1}{2}\right) - \mathcal{N}'(n, a)$, where $\mathcal{N}'(n, a)$ is the number of $E \in \{-1, 1\}^{n-1}$ such that $S_{n-1}(E) = a - 1$ and $S_m(E) \geq a$ for some $1 \leq m \leq n - 1$, we need only compute $\mathcal{N}'(n, a)$. For this purpose we will use a beautiful argument known as the *reflection*

³As the following indicates, I take the infimum over the empty set to be $+\infty$.

principle. Namely, consider the set P(n, a) of paths $(S_0, \ldots, S_{n-1}) \in \mathbb{Z}^n$ with the properties that $S_0 = 0$, $S_m - S_{m-1} \in \{-1, 1\}$ for $1 \le m \le n-1$, and $S_m \ge a$ for some $1 \le m \le n-1$. Clearly, $\mathcal{N}'(n, a)$ is the number of paths in the set L(n, a) consisting of those $(S_0, \ldots, S_{n-1}) \in P(n, a)$ for which $S_{n-1} = a - 1$, and, as an application of the reflection principle, we will show that the set L(n, a) has the same number of elements as the set U(n, a) whose elements are those paths $(S_0, \ldots, S_{n-1}) \in P(n, a)$ for which $S_{n-1} = a + 1$. Since $(S_0, \ldots, S_{n-1}) \in U(n, a)$ if and only if $S_0 = 0$, $S_m - S_{m-1} \in \{-1, 1\}$ for all $1 \le m \le n - 1$, and $S_{n-1} = a + 1$, the condition that $S_m \ge a$ is redundant. Thus, we already know how to count them: there are $\binom{n-1}{n+a}$ of them. Hence, all that remains is to make the aforementioned application of the reflection principle. To this end, for a given $\mathbf{S} = (S_0, \dots, S_{n-1}) \in P(n, a)$, let $\ell(\mathbf{S})$ be the smallest $0 \le k \le n - 1$ for which $S_k \ge a$, and define the *reflection* $\Re(\mathbf{S}) = (\hat{S}_0, \dots, \hat{S}_{n-1})$ of \mathbf{S} so that $\hat{S}_m = S_m$ if $0 \le m \le \ell(\mathbf{S})$ and $\hat{S}_k = 2a - S_k$ if $\ell(\mathbf{S}) < m \le n-1$. Clearly, \mathfrak{R} maps L(n, a) into U(n, a) and U(n, a) into L(n, a). In addition, \Re is its own inverse: its composition with itself is the identity map. Therefore, as a map from L(n, a) to U(n, a), \mathfrak{R} it must be both one-to-one and onto, and so L(n, a) and U(n, a) have the same numbers of elements.

We have now shown that $\mathcal{N}'(n, a) = {\binom{n-1}{\frac{n+2}{2}}}$ and therefore that

$$\mathcal{N}(n,a) = \binom{n-1}{\frac{n+a}{2}-1} - \binom{n-1}{\frac{n+a}{2}}.$$

Finally, after plugging this into (*), we arrive at

$$\mathbb{P}(\zeta^{\{a\}}=n)=\left[\binom{n-1}{\frac{n+a}{2}-1}-\binom{n-1}{\frac{n+a}{2}}\right]p^{\frac{n+a}{2}}q^{\frac{n-a}{2}},$$

which simplifies to the remarkably simple expression

$$\mathbb{P}\left(\zeta^{\{a\}}=n\right)=\frac{a}{n}\binom{n}{\frac{n+a}{2}}p^{\frac{n+a}{2}}q^{\frac{n-a}{2}}=\frac{a}{n}\mathbb{P}(X_n=a).$$

The computation when a < 0 can be carried out either by repeating the argument just given or, after reversing the roles of p and q, applying the preceding result to -a. However one arrives at it, the general result is that

$$a \neq 0 \quad \Longrightarrow \quad \mathbb{P}\left(\zeta^{\{a\}} = n\right) = \frac{|a|}{n} \binom{n}{\frac{n+a}{2}} p^{\frac{n+a}{2}} q^{\frac{n-a}{2}} = \frac{|a|}{n} \mathbb{P}(X_n = a) \quad (1.1.7)$$

for $n \ge |a|$ with the same parity as a and is 0 otherwise.

1.1.3 Some Related Computations

Although the formula in (1.1.7) is elegant, it is not particularly transparent. For example, it is not at all evident how one can use it to determine whether $\mathbb{P}(\zeta^{\{a\}} < \infty) = 1$. To carry out this computation, let a > 0 be given, and represent $\zeta^{\{a\}}$ as $f_a(B_1, \ldots, B_n, \ldots)$, where f_a is the function which maps $\{-1, 1\}^{\mathbb{Z}^+}$ into $\mathbb{Z}^+ \cup \{\infty\}$ so that, for each $n \in \mathbb{N}$,

$$f_a(\epsilon_1,\ldots,\epsilon_n,\ldots) > n \quad \Longleftrightarrow \quad \sum_{\ell=1}^m \epsilon_\ell < a \quad \text{for } 1 \le m \le n.$$

Because the event $\{\zeta^{\{a\}} = m\}$ depends only on (B_1, \ldots, B_m) and

$$\zeta^{\{a\}} = m \implies \zeta^{\{a+1\}} = m + \zeta^{\{1\}} \circ \Sigma^m$$

where $\zeta_1 \circ \Sigma^m \equiv f_1(B_{m+1}, \dots, B_{m+n}, \dots), (1.1.8)$

 $\{\zeta^{\{a\}} = m \& \zeta^{\{a+1\}} < \infty\} = \{\zeta^{\{a\}} = m\} \cap \{\zeta^{\{1\}} \circ \Sigma^m < \infty\}, \text{ and } \{\zeta^{\{a\}} = m\} \text{ is independent of } \{\zeta^{\{1\}} \circ \Sigma^m < \infty\}.$ In particular, this leads to

$$\mathbb{P}(\zeta^{\{a+1\}} < \infty) = \sum_{m=1}^{\infty} \mathbb{P}(\zeta^{\{a\}} = m \& \zeta_{a+1} < \infty)$$
$$= \sum_{m=1}^{\infty} \mathbb{P}(\zeta^{\{a\}} = m) \mathbb{P}(\zeta^{\{1\}} \circ \Sigma^m < \infty)$$
$$= \mathbb{P}(\zeta^{\{1\}} < \infty) \sum_{m=1}^{\infty} \mathbb{P}(\zeta^{\{a\}} = m) = \mathbb{P}(\zeta^{\{1\}} < \infty) \mathbb{P}(\zeta^{\{a\}} < \infty),$$

since $(B_{m+1}, \ldots, B_{m+n}, \ldots)$ and $(B_1, \ldots, B_n, \ldots)$ have the same distribution and therefore so do $\zeta^{\{1\}} \circ \Sigma^m$ and $\zeta^{\{1\}}$. The same reasoning applies equally well when a < 0, only now with -1 playing the role of 1. In other words, we have proved that

$$\mathbb{P}\big(\zeta^{\{a\}} < \infty\big) = \mathbb{P}\big(\zeta^{\{\operatorname{sgn}(a)\}} < \infty\big)^{|a|} \quad \text{for } a \in \mathbb{Z} \setminus \{0\},$$
(1.1.9)

where sgn(*a*), the *signum* of *a*, is 1 or -1 according to whether a > 0 or a < 0. In particular, this shows that $\mathbb{P}(\zeta^{\{1\}} < \infty) = 1 \implies \mathbb{P}(\zeta^{\{a\}} < \infty) = 1$ and $\mathbb{P}(\zeta^{\{-1\}} < \infty) = 1 \implies \mathbb{P}(\zeta^{\{-a\}} < \infty) = 1$ for all $a \in \mathbb{Z}^+$.

In view of the preceding, we need only look at $\mathbb{P}(\zeta^{\{1\}} < \infty)$. Moreover, by the monotone convergence theorem, Theorem 7.1.9,

$$\mathbb{P}(\zeta^{\{1\}} < \infty) = \lim_{s \neq 1} \mathbb{E}[s^{\zeta^{\{1\}}}] = \lim_{s \neq 1} \sum_{n=1}^{\infty} s^{2n-1} \mathbb{P}(\zeta^{\{1\}} = 2n-1).$$

Applying (1.1.7) with a = 1, we know that

$$\mathbb{P}(\zeta^{\{1\}} = 2n - 1) = \frac{1}{2n - 1} \binom{2n - 1}{n} p^n q^{n - 1}.$$

Next, note that

$$\frac{1}{2n-1} \binom{2n-1}{n} = \frac{(2(n-1))!}{n!(n-1)!} = \frac{2^{n-1}}{n!} \prod_{m=1}^{n-1} (2m-1)$$
$$= \frac{4^{n-1}}{n!} \prod_{m=1}^{n-1} (m-\frac{1}{2}) = (-1)^{n-1} \frac{4^n}{2} \binom{\frac{1}{2}}{n},$$

where,⁴ for any $\alpha \in \mathbb{R}$,

$$\binom{\alpha}{n} \equiv \begin{cases} 1 & \text{if } n = 0\\ \frac{1}{n!} \prod_{m=0}^{n-1} (\alpha - m) & \text{if } n \in \mathbb{Z}^+ \end{cases}$$

is the *generalized binomial coefficient* which gives the coefficient of x^n in the Taylor's expansion of $(1 + x)^{\alpha}$ around x = 0. Hence, since $4pq \le 1$,

$$\sum_{n=1}^{\infty} s^{2n-1} \mathbb{P}(\zeta^{\{1\}} = 2n-1) = -\frac{1}{2qs} \sum_{n=1}^{\infty} {\binom{\frac{1}{2}}{n}} (-4pqs^2)^n = \frac{1-\sqrt{1-4pqs^2}}{2qs},$$

and so

$$\mathbb{E}\left[s^{\zeta^{\{1\}}}\right] = \frac{1 - \sqrt{1 - 4pqs^2}}{2qs} \quad \text{for } |s| < 1.$$
(1.1.10)

Obviously, by symmetry, one can reverse the roles of p and q to obtain

$$\mathbb{E}\left[s^{\zeta^{\{-1\}}}\right] = \frac{1 - \sqrt{1 - 4pqs^2}}{2ps} \quad \text{for } |s| < 1.$$
(1.1.11)

By letting $s \nearrow 1$ in (1.1.10) and noting that $1 - 4pq = (p+q)^2 - 4pq = (p-q)^2$, we see that⁵

$$\lim_{s \neq 1} \mathbb{E}[s^{\zeta^{(1)}}] = \frac{1 - |p - q|}{2q} = \frac{p \wedge q}{q}$$

and so

$$\mathbb{P}(\zeta^{\{1\}} < \infty) = \begin{cases} 1 & \text{if } p \ge q \\ \frac{p}{q} & \text{if } p < q \end{cases}$$

⁴In the preceding, I have adopted the convention that $\prod_{j=k}^{\ell} a_j = 1$ if $\ell < k$.

⁵I use $a \wedge b$ to denote the minimum min $\{a, b\}$ of $a, b \in \mathbb{R}$.

Of course, $\mathbb{P}(\zeta^{\{-1\}} < \infty)$ is given by the same formula, only with the roles of *p* and *q* reversed. Thus, by (1.1.12),

$$\mathbb{P}\left(\zeta^{\{a\}} < \infty\right) = \begin{cases} 1 & \text{if } a \in \mathbb{Z}^+ \& p \ge q \text{ or } -a \in \mathbb{Z}^+ \& p \le q \\ \left(\frac{p}{q}\right)^a & \text{if } a \in \mathbb{Z}^+ \& p < q \text{ or } -a \in \mathbb{Z}^+ \& p > q. \end{cases}$$
(1.1.12)

1.1.4 Time of First Return

Having gone to so much trouble to arrive at (1.1.12), it is only reasonable to draw from it a famous conclusion about the *recurrence* properties of nearest neighbor random walks on \mathbb{Z} . Namely, let

$$\rho_0 \equiv \inf\{n \ge 1 : X_n = 0\} \quad (\equiv \infty \text{ if } X_n \neq 0 \text{ for all } n \ge 1)$$

be the *time of first return* to 0. Then, by precisely the same sort of reasoning which allowed us to arrive at (1.1.9), we see that $\mathbb{P}(X_1 = 1 \& \rho_0 < \infty) = p\mathbb{P}(\zeta_{-1} < \infty)$ and $\mathbb{P}(X_1 = -1 \& \rho_0 < \infty) = q\mathbb{P}(\zeta^{\{1\}} < \infty)$, and so, by (1.1.12),

$$\mathbb{P}(\rho_0 < \infty) = 2(p \land q). \tag{1.1.13}$$

In other words, the random walk $\{X_n : n \ge 0\}$ will return to 0 with probability 1 if and only if it is symmetric in the sense that $p = \frac{1}{2}$.

By sharpening the preceding a little, one sees that $\mathbb{P}(X_1 = 1 \& \rho_0 = 2n) = p\mathbb{P}(\zeta^{\{1\}} = 2n - 1)$ and $\mathbb{P}(X_1 = -1 \& \rho_0 = 2n) = q\mathbb{P}(\zeta^{\{1\}} = 2n - 1)$, and so, by (1.1.10) and (1.1.11),

$$\mathbb{E}[s^{\rho_0}] = 1 - \sqrt{1 - 4pqs^2} \quad \text{for } |s| < 1.$$
 (1.1.14)

Hence,

$$\mathbb{E}\left[\rho_0 s^{\rho_0}\right] = s \frac{d}{ds} \mathbb{E}\left[s^{\rho_0}\right] = \frac{4pqs^2}{\sqrt{1 - 4pqs^2}} \quad \text{for } |s| < 1,$$

and therefore, since $\mathbb{E}[\rho_0 s^{\rho_0}] \nearrow \mathbb{E}[\rho_0, \rho_0 < \infty]$ as $s \nearrow 1$,

$$\mathbb{E}[\rho_0, \rho_0 < \infty] = \frac{4pq}{|p-q|},$$

which, in conjunction with (1.1.13), means that⁷

$$\mathbb{E}[\rho_0 \mid \rho_0 < \infty] = \frac{2(p \lor q)}{|p - q|} = 1 + \frac{1}{|p - q|}.$$
(1.1.15)

⁶When X is a random variable and A is an event, we will often use $\mathbb{E}[X, A]$ to denote $\mathbb{E}[X\mathbf{1}_A]$. ⁷ $a \lor b$ is used to denote the maximum max $\{a, b\}$ of $a, b \in \mathbb{R}$.

The conclusions drawn in the preceding provide significant insight into the behavior of nearest neighbor random walks on \mathbb{Z} . In the first place, they say that when the random walk is symmetric, it returns to 0 with probability 1 but the expected amount of time it takes to do so is infinite. Secondly, when the random walk is not symmetric, it will, with positive probability, fail to return. On the other hand, in the non-symmetric case, the behavior of the trajectories is interesting. Namely, (1.1.13) in combination with (1.1.15) say that either they fail to return at all or they return relatively quickly.

1.1.5 Passage Times via Functional Equations

I close this discussion of passage times for nearest neighbor random walks with a less computational derivation of (1.1.10). For this purpose, set $u_a(s) = \mathbb{E}[s^{\zeta^{\{a\}}}]$ for $a \in \mathbb{Z} \setminus \{0\}$ and $s \in (-1, 1)$. Given $a \in \mathbb{Z}^+$, we use the ideas in Sect. 1.1.3, especially (1.1.8), to arrive at

$$u_{a+1}(s) = \sum_{m=1}^{\infty} s^m \mathbb{E} \left[s^{\zeta^{\{1\}} \circ \Sigma^m}, \, \zeta^{\{a\}} = m \right] = \sum_{m=1}^{\infty} s^m \mathbb{P} \left(\zeta^{\{a\}} = m \right) \mathbb{E} \left[s^{\zeta^{\{1\}} \circ \Sigma^m} \right]$$
$$= \sum_{m=1}^{\infty} s^m \mathbb{P} (\zeta_a = m) u_1(s) = u_a(s) u_1(s).$$

Similarly, if $-a \in \mathbb{Z}^+$, then $u_{a-1}(s) = u_a(s)u_{-1}(s)$. Hence

$$u_a(s) = u_{\text{sgn}(a)}(s)^{|a|}$$
 for $a \in \mathbb{Z} \setminus \{0\}$ and $|s| < 1$. (1.1.16)

Continuing with the same line of reasoning, we also have

$$u_1(s) = \mathbb{E}[s^{\zeta^{(1)}}, X_1 = 1] + \mathbb{E}[s^{\zeta^{(1)}}, X_1 = -1]$$

= $ps + qs\mathbb{E}[s^{\zeta_2 \circ \Sigma^1}] = ps + qsu_2(s) = ps + qsu_1(s)^2.$

Hence, by the quadratic formula,

$$u_1(s) = \frac{1 \pm \sqrt{1 - 4pqs^2}}{2qs}.$$

Because $\mathbb{P}(\zeta^{\{1\}} \text{ is odd}) = 1$, $u_1(-s) = -u_1(s)$. At the same time,

$$s \in (0,1) \quad \Longrightarrow \quad \frac{1 + \sqrt{1 - 4pqs^2}}{2qs} > \frac{1 + \sqrt{1 - 4pq}}{2q} = \frac{p \lor q}{q} \ge 1.$$

Hence, since $s \in (0, 1) \implies u_1(s) < 1$, we can eliminate the "+" solution and thereby arrive at a second derivation of (1.1.10). In fact, after combining this with

1.2 Recurrence Properties of Random Walks

(1.1.16), we have shown that

$$\mathbb{E}\left[s^{\zeta_{a}}\right] = \begin{cases} (\frac{1-\sqrt{1-4pqs^{2}}}{2qs})^{a} & \text{if } a \in \mathbb{Z}^{+} \\ (\frac{1-\sqrt{1-4pqs^{2}}}{2ps})^{-a} & \text{if } -a \in \mathbb{Z}^{+} \end{cases}$$
 (1.1.17)

In theory one can recover (1.1.7) by differentiating (1.1.17) *n* times with respect to *s* at 0, but the computation is tedious.

1.2 Recurrence Properties of Random Walks

In Sect. 1.1.4, we studied the time ρ_0 of first return of a nearest neighbor random walk to 0. As we will see in Chaps. 2 and 4, times of first return are critical (cf. Sect. 2.3.2) for an understanding of the long time behavior of random walks and related processes. Indeed, when the random walk returns to 0, it starts all over again. Thus, if it returns with probability 1, then the entire history of the walk will consist of epochs, each epoch being a sojourn which begins and ends at 0. Because it marks the time at which one epoch ends and a second, identically distributed, one begins, a time of first return is often called a *renewal time*, and the walk is said to be *recurrent* if $\mathbb{P}(\rho_0 < \infty) = 1$. Walks which are not recurrent are said to be *transient*.

In this section, we will discuss the recurrence properties of nearest neighbor random walks. Of course, we already know (cf. (1.1.13)) that a nearest neighbor random on \mathbb{Z} is recurrent if and only if it is symmetric. Thus, our interest here will be in higher dimensional analogs. In particular, in the hope that it will be convincing evidence that recurrence is subtle, we will show that the recurrence of the nearest neighbor, symmetric random walk on \mathbb{Z} persists when \mathbb{Z} is replaced by \mathbb{Z}^2 but disappears in \mathbb{Z}^3 .

1.2.1 Random Walks on \mathbb{Z}^d

To describe the analog on \mathbb{Z}^d of a nearest neighbor random walk on \mathbb{Z} , we begin by thinking of $\mathbf{N}_1 \equiv \{-1, 1\}$ as the set of *nearest neighbors* in \mathbb{Z} of 0. It should then be clear why the set \mathbf{N}_d of nearest neighbors of the origin in \mathbb{Z}^d consists of the 2*d* points in \mathbb{Z}^d for which (d - 1) coordinates are 0 and the remaining coordinate is in \mathbf{N}_1 . Next, we replace the \mathbf{N}_1 -valued Bernoulli random variables in Sect. 1.1 by their *d*-dimensional analogs, that is: independent, identically distributed \mathbf{N}_d -valued random variables $\mathbf{B}_1, \ldots, \mathbf{B}_n, \ldots$.⁸ Finally, a nearest neighbor random walk on \mathbb{Z}^d

⁸The existence of the \mathbf{B}_n 's can be seen as a consequence of Theorem 7.3.2. Namely, let $\{U_n : n \in \mathbb{Z}^+\}$ be a family of mutually independent random variables which are uniformly distributed on [0, 1). Next, let $(\mathbf{k}_1, \ldots, \mathbf{k}_{2d})$ be an ordering of the elements of \mathbf{N}_d , set $\beta_0 = 0$ and $\beta_m = \sum_{\ell=1}^m \mathbb{P}(\mathbf{B}_1 = \mathbf{k}_\ell)$ for $1 \le m \le 2d$, define $F : [0, 1) \longrightarrow \mathbf{N}_d$ so that $F \upharpoonright [\beta_{m-1}, \beta_m) = \mathbf{k}_m$, and set $\mathbf{B}_n = F(U_n)$.

is a family $\{\mathbf{X}_n : n \ge 0\}$ of the form

$$\mathbf{X}_0 = \mathbf{0}$$
 and $\mathbf{X}_n = \sum_{m=1}^n \mathbf{B}_m$ for $n \ge 1$.

The equivalent, stochastic process oriented description of $\{X_n : n \ge 0\}$ is

$$\mathbb{P}(\mathbf{X}_0 = \mathbf{0}) = 1 \quad \text{and}, \quad \text{for } n \ge 1 \text{ and } \boldsymbol{\epsilon} \in \mathbf{N}_d,$$

$$\mathbb{P}(\mathbf{X}_n - \mathbf{X}_{n-1} = \boldsymbol{\epsilon} \mid \mathbf{X}_0, \dots, \mathbf{X}_{n-1}) = p_{\boldsymbol{\epsilon}},$$

(1.2.1)

where $p_{\epsilon} \equiv \mathbb{P}(\mathbf{B}_1 = \epsilon)$. When \mathbf{B}_1 is uniformly distributed on \mathbf{N}_d , the random walk is said to be *symmetric*.

In keeping with the notation and terminology introduced above, we define the time ρ_0 of first return to the origin to be *n* if $n \ge 1$, $\mathbf{X}_n = \mathbf{0}$, and $\mathbf{X}_m \neq \mathbf{0}$ for $1 \le m < n$, and we take $\rho_0 = \infty$ if no such $n \ge 1$ exists. Also, we will say that the walk is *recurrent* or *transient* according to whether $\mathbb{P}(\rho_0 < \infty)$ is 1 or strictly less than 1.

1.2.2 An Elementary Recurrence Criterion

Given $n \ge 1$, let $\rho_{\mathbf{0}}^{(n)}$ be the time of the *n*th return to **0**. That is, $\rho_{\mathbf{0}}^{(1)} = \rho_{\mathbf{0}}$ and, for $n \ge 2$,

$$\rho_{\mathbf{0}}^{(n-1)} < \infty \quad \Longrightarrow \quad \rho_{\mathbf{0}}^{(n)} = \inf\{m > \rho^{(n-1)} : \mathbf{X}_m = \mathbf{0}\}$$

and $\rho_{\mathbf{0}}^{(n-1)} = \infty \implies \rho_{\mathbf{0}}^{(n)} = \infty$. Equivalently, if $g : (\mathbf{N}_d)^{\mathbb{Z}^+} \longrightarrow \mathbb{Z}^+ \cup \{\infty\}$ is determined so that

$$g(\boldsymbol{\epsilon}_1,\ldots\boldsymbol{\epsilon}_\ell,\ldots) > n$$
 if and only if $\sum_{\ell=1}^m \boldsymbol{\epsilon}_\ell \neq \mathbf{0}$ for $1 \le m \le n$,

then $\rho_0 = g(\mathbf{B}_1, \dots, \mathbf{B}_{\ell}, \dots)$, and $\rho_0^{(n)} = m \implies \rho_0^{(n+1)} = m + \rho_0 \circ \Sigma^m$ where $\rho_0 \circ \Sigma^m$ is equal to $g(\mathbf{B}_{m+1}, \dots, \mathbf{B}_{m+\ell}, \dots)$. In particular, this leads to

$$\begin{split} \mathbb{P}\big(\rho_{\mathbf{0}}^{(n+1)} < \infty\big) &= \sum_{m=1}^{\infty} \mathbb{P}\big(\rho_{\mathbf{0}}^{(n)} = m \And \rho_{\mathbf{0}} \circ \Sigma^{m} < \infty\big) \\ &= \mathbb{P}\big(\rho_{\mathbf{0}}^{(n)} < \infty\big) \mathbb{P}(\rho_{\mathbf{0}} < \infty), \end{split}$$

since $\{\rho_0^{(n)} = m\}$ depends only on $(\mathbf{B}_1, \dots, \mathbf{B}_m)$, and is therefore independent of $\rho_0 \circ \Sigma^m$, and the distribution of $\rho_0 \circ \Sigma^m$ is the same as that of ρ_0 . Thus, we have proved that

$$\mathbb{P}(\rho_{\mathbf{0}}^{(n)} < \infty) = \mathbb{P}(\rho_{\mathbf{0}} < \infty)^{n} \quad \text{for } n \ge 1.$$
(1.2.2)

One dividend of (1.2.2) is that it supports the epochal picture described above about the structure of recurrent walks. Namely, it says that *if the walk returns once* to **0** with probability 1, then, with probability 1, it will do so infinitely often. This observation has many applications. For example, it shows that if the mean value of *i*th coordinate of **B**₁ is different from 0, then {**X**_n : $n \ge 0$ } must be transient. To see this, use Y_n to denote the *i*th coordinate of **B**_n, and observe that { $Y_n - Y_{n-1} : n \ge 1$ } is a sequence of mutually independent, identically distributed {-1, 0, 1}-valued random variables with mean value $\mu \ne 0$. By the strong law of large numbers (cf. Exercise 1.3.4 below), this means that $\frac{Y_n}{n} \longrightarrow \mu \ne 0$ with probability 1, which is possible only if $|\mathbf{X}_n| \ge |Y_n| \longrightarrow \infty$ with probability 1, and clearly this eliminates the possibility that, even with positive probability, $\mathbf{X}_n = \mathbf{0}$ infinitely often.

A second dividend of (1.2.2) is the following. Define

$$T_{\mathbf{0}} = \sum_{n=0}^{\infty} \mathbf{1}_{\{\mathbf{0}\}}(\mathbf{X}_n)$$

to be the total time that $\{\mathbf{X}_n : n \ge 0\}$ spends at the origin. Since $\mathbf{X}_0 = \mathbf{0}, T_{\mathbf{0}} \ge 1$. Moreover, for $n \ge 1, T_{\mathbf{0}} > n \iff \rho_{\mathbf{0}}^{(n)} < \infty$. Hence, by (1.2.2),

$$\mathbb{E}[T_{\mathbf{0}}] = \sum_{n=0}^{\infty} \mathbb{P}(T_{\mathbf{0}} > n) = 1 + \sum_{n=1}^{\infty} \mathbb{P}(\rho_{\mathbf{0}}^{(n)} < \infty) = 1 + \sum_{n=1}^{\infty} \mathbb{P}(\rho_{\mathbf{0}} < \infty)^{n},$$

and so

$$\mathbb{E}[T_0] = \frac{1}{1 - \mathbb{P}(\rho_0 < \infty)} = \frac{1}{\mathbb{P}(\rho_0 = \infty)}.$$
 (1.2.3)

Before applying (1.2.3) to the problem of recurrence, it is interesting to note that T_0 is a random variable for which the following peculiar dichotomy holds:

$$\mathbb{P}(T_0 < \infty) > 0 \implies \mathbb{E}[T_0] < \infty$$

$$\mathbb{E}[T_0] = \infty \implies \mathbb{P}(T_0 = \infty) = 1.$$
 (1.2.4)

Indeed, if $\mathbb{P}(T_0 < \infty) > 0$, then, with positive probability, \mathbf{X}_n cannot be **0** infinitely often, and so, by (1.2.2), $\mathbb{P}(\rho_0 < \infty) < 1$, which, by (1.2.3), means that $\mathbb{E}[T_0] < \infty$. On the other hand, if $\mathbb{E}[T_0] = \infty$, then (1.2.3) implies that $\mathbb{P}(\rho_0 < \infty) = 1$ and therefore, by (1.2.2), that $\mathbb{P}(T_0 > n) = \mathbb{P}(\rho_0^{(n)} < \infty) = 1$ for all $n \ge 1$. Hence (cf. (7.1.3)), $\mathbb{P}(T_0 = \infty) = 1$.

1.2.3 Recurrence of Symmetric Random Walk in \mathbb{Z}^2

The most frequent way that (1.2.3) gets applied to determine recurrence is in conjunction with the formula

$$\mathbb{E}[T_0] = \sum_{n=0}^{\infty} \mathbb{P}(\mathbf{X}_n = \mathbf{0}).$$
(1.2.5)

Although the proof of (1.2.5) is essentially trivial (cf. Theorem 7.1.15):

$$\mathbb{E}[T_0] = \mathbb{E}\left[\sum_{n=0}^{\infty} \mathbf{1}_{\{0\}}(\mathbf{X}_n)\right] = \sum_{n=0}^{\infty} \mathbb{E}\left[\mathbf{1}_{\{0\}}(\mathbf{X}_n)\right] = \sum_{n=0}^{\infty} \mathbb{P}(\mathbf{X}_n = \mathbf{0}),$$

in conjunction with (1.2.3) it becomes powerful. Namely, it says that

$$\{\mathbf{X}_n : n \ge 0\}$$
 is recurrent if and only if $\sum_{n=0}^{\infty} \mathbb{P}(\mathbf{X}_n = \mathbf{0}) = \infty,$ (1.2.6)

and, since $\mathbb{P}(\mathbf{X}_n = \mathbf{0})$ is more amenable to estimation than quantities which involve knowing the trajectory at more than one time, this is valuable information.

In order to apply (1.2.6) to symmetric random walks, it is important to know that when the walk is symmetric, then **0** is the most likely place for the walk to be at any even time. To verify this, note that if $\mathbf{k} \in \mathbb{Z}^d$, then

$$\mathbb{P}(\mathbf{X}_{2n} = \mathbf{k}) = \sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell \& \mathbf{X}_{2n} - \mathbf{X}_n = \mathbf{k} - \ell)$$

$$= \sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell) \mathbb{P}(\mathbf{X}_{2n} - \mathbf{X}_n = \mathbf{k} - \ell)$$

$$= \sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell) \mathbb{P}(\mathbf{X}_n = \mathbf{k} - \ell)$$

$$\leq \left(\sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell)^2\right)^{\frac{1}{2}} \left(\sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \mathbf{k} - \ell)^2\right)^{\frac{1}{2}} = \sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell)^2,$$

where, in the passage to the last line, we have applied Schwarz's inequality (cf. Exercise 1.3.1 below). Up to this point we have not used symmetry. However, if the walk is symmetric, then $\mathbb{P}(\mathbf{X}_n = \ell) = \mathbb{P}(\mathbf{X}_n = -\ell)$, and so the last line of the preceding can be continued as

$$\sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell) \mathbb{P}(\mathbf{X}_n = -\ell)$$
$$= \sum_{\ell \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_n = \ell) \mathbb{P}(\mathbf{X}_{2n} - \mathbf{X}_n = -\ell) = \mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}).$$

1.2 Recurrence Properties of Random Walks

Thus,

$$\{\mathbf{X}_n : n \ge 0\} \text{ symmetric} \implies \mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) = \max_{\mathbf{k} \in \mathbb{Z}^d} \mathbb{P}(\mathbf{X}_{2n} = \mathbf{k}).$$
(1.2.7)

To develop a feeling for how these considerations get applied, we begin by using them to give a second derivation of the recurrence of the nearest neighbor, symmetric random walk on \mathbb{Z} . For this purpose, note that, because $\mathbb{P}(|X_n| \le n) = 1$, (1.2.7) implies that

$$1 = \sum_{\ell=-2n}^{2n} \mathbb{P}(X_{2n} = \ell) \le (4n+1)\mathbb{P}(X_{2n} = 0),$$

and therefore, since the harmonic series diverges, that $\sum_{n=0}^{\infty} \mathbb{P}(X_n = 0) = \infty$.

The analysis for the symmetric, nearest neighbor random walk in \mathbb{Z}^2 requires an additional ingredient. Namely, the *d*-dimensional analog of the preceding line of reasoning would lead to $\mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) \ge (4n + 1)^{-d}$, which is inconclusive except when d = 1. In order to do better, we need to use the fact that

$$\{\mathbf{X}_n : n \ge 0\}$$
 symmetric $\implies \mathbb{E}[|\mathbf{X}_n|^2] = n.$ (1.2.8)

To prove (1.2.8), note that each coordinate of \mathbf{B}_n is a random variable with mean value 0 and variance $\frac{1}{d}$. Hence, because the \mathbf{B}_n 's are mutually independent, the second moment of each coordinate of \mathbf{X}_n is $\frac{n}{d}$.

Knowing (1.2.8), Markov's inequality (7.1.12) says that

$$\mathbb{P}(|\mathbf{X}_{2n}| \geq 2\sqrt{n}) \leq \frac{1}{4n} \mathbb{E}[|\mathbf{X}_{2n}|^2] = \frac{1}{2},$$

which allows us to sharpen the preceding argument to give9

$$\frac{1}{2} \leq \mathbb{P}(|\mathbf{X}_{2n}| < 2\sqrt{n}) = \sum_{|\boldsymbol{\ell}| < 2\sqrt{n}} \mathbb{P}(\mathbf{X}_{2n} = \boldsymbol{\ell})$$
$$\leq (4\sqrt{n} + 1)^d \mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) \leq 2^{d-1} (4n^{\frac{d}{2}} + 1) \mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}).$$

That is, we have now shown that

$$\mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) \ge 2^{-d} \left(4n^{\frac{d}{2}} + 1\right)^{-1}$$
(1.2.9)

for the symmetric, nearest neighbor random walk on \mathbb{Z}^d . In particular, when d = 2, this proves that *the symmetric, nearest neighbor random walk on* \mathbb{Z}^2 *is recurrent.*

⁹For any $a, b \in [0, \infty)$ and $p \in [1, \infty)$, $(a + b)^p \le 2^{p-1}(a^p + b^p)$. This can be seen as an application of Jensen's inequality (cf. Exercise 6.6.2), which, in this case, is simply the statement that $x \in [0, \infty) \mapsto x^p$ is convex.

1.2.4 Transience in \mathbb{Z}^3

Although (1.2.9) was sufficient to prove recurrence for the symmetric, nearest neighbor random walk in \mathbb{Z}^2 , it leaves open the possibility of transience in \mathbb{Z}^d for $d \ge 3$. Thus, in order to nail down the question when $d \ge 3$, we will need to see how good an estimate (1.2.9) really is. In particular, it would suffice to prove that there is an upper bound of the same form.

To get an upper bound which complements the lower bound in (1.2.9), we first do so in the case when d = 1. For this purpose, let $0 \le \ell \le n$ be given, and observe that

$$\frac{\mathbb{P}(X_{2n} = 2\ell)}{\mathbb{P}(X_{2n} = 0)} = \frac{(n!)^2}{(n+\ell)!(n-\ell)!} = \frac{n(n-1)\cdots(n-\ell+1)}{(n+\ell)(n+\ell-1)\cdots(n+1)}$$
$$= \prod_{k=0}^{\ell-1} \left(1 - \frac{\ell}{n+\ell-k}\right) \ge \left(1 - \frac{\ell}{n+1}\right)^{\ell}.$$

Now recall that

$$\log(1-x) = -\sum_{m=1}^{\infty} \frac{x^m}{m} \quad \text{for } |x| < 1 \tag{1.2.10}$$

and therefore that $\log(1-x) \ge -\frac{3x}{2}$ for $0 \le x \le \frac{1}{2}$. Hence, the preceding shows that

$$\frac{\mathbb{P}(X_{2n}=2\ell)}{\mathbb{P}(X_{2n}=0)} \ge \exp\left(\ell \log \frac{\ell}{n+1}\right) \ge e^{-\frac{3\ell^2}{2(n+1)}}$$

as long as $0 \le \ell \le \frac{n+1}{2}$. Because $\mathbb{P}(X_{2n} = -2\ell) = \mathbb{P}(X_{2n} = 2\ell)$, we can now say that

$$\mathbb{P}(X_{2n} = 0) \le e^{\frac{3}{2}} \mathbb{P}(X_{2n} = 2\ell) \quad \text{for } |\ell| \le \sqrt{n}$$

But, because $\sum_{\ell \in \mathbb{Z}} \mathbb{P}(X_{2n} = 2\ell) = 1$, this means that $(2\sqrt{n} - 1)\mathbb{P}(X_{2n} = 0) \le e^{\frac{3}{2}}$, and so

$$\mathbb{P}(X_{2n} = 0) \le e^{\frac{3}{2}} (2\sqrt{n} - 1)^{-1}, \quad n \ge 1,$$
(1.2.11)

when $\{X_n : n \ge 0\}$ is the symmetric, nearest neighbor random walk on \mathbb{Z} .

If, as they most definitely are not, the coordinates of the symmetric, nearest neighbor random walk were independent, then (1.2.11) would yield the sort of upper bound for which we are looking. Thus it is reasonable to examine to what extent we can relate the symmetric, nearest neighbor random walk on \mathbb{Z}^d to d mutually independent symmetric, nearest neighbor random walks $\{X_{i,n} : n \ge 0\}, 1 \le i \le d$, on \mathbb{Z} . To this end, refer to (1.2.1) with $p_{\epsilon} \equiv \frac{1}{2d}$, and think of choosing $\mathbf{X}_n - \mathbf{X}_{n-1}$ in two steps: first choose the coordinate which is to be non-zero and then choose whether it is to be +1 or -1. With this in mind, let $\{I_n : n \ge 0\}$ be a sequence of $\{1, \ldots, d\}$ -valued, mutually independent, uniformly distributed random variables

which are independent of $\{X_{i,n} : 1 \le i \le d \& n \ge 0\}$, set, for $1 \le i \le d$, $N_{i,0} = 0$ and $N_{i,n} = \sum_{m=1}^{n} \mathbf{1}_{\{i\}}(I_m)$ when $n \ge 1$, and consider the sequence $\{\mathbf{Y}_n : n \ge 0\}$ given by

$$\mathbf{Y}_n = (X_{1,N_{1,n}}, \dots, X_{d,N_{d,n}}). \tag{1.2.12}$$

Without too much effort, one can check that $\{\mathbf{Y}_n : n \ge 0\}$ satisfies the conditions in (1.2.1) for the symmetric, nearest neighbor random walk on \mathbb{Z}^d and therefore has the same distribution as $\{\mathbf{X}_n : n \ge 0\}$. In particular, by (1.2.11),

$$\mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) = \sum_{\mathbf{m} \in \mathbb{N}^d} \mathbb{P}(X_{i,2m_i} = 0 \& N_{i,2n} = 2m_i \text{ for } 1 \le i \le d)$$

$$= \sum_{\substack{\mathbf{m} \in \mathbb{N}^d \\ m_1 \land \dots \land m_d \ge \frac{n}{d}}} \left(\prod_{i=1}^d \mathbb{P}(X_{i,2m_i} = 0) \right) \mathbb{P}(N_{i,2n} = 2m_i \text{ for } 1 \le i \le d)$$

$$+ \sum_{\substack{\mathbf{m} \in \mathbb{N}^d \\ m_1 \land \dots \land m_d < \frac{n}{2d}}} \left(\prod_{i=1}^d \mathbb{P}(X_{i,2m_i} = 0) \right) \mathbb{P}(N_{i,2n} = 2m_i \text{ for } 1 \le i \le d)$$

$$\le e^{\frac{3d}{2}} \left(2\sqrt{\frac{n}{d}} - 1 \right)^{-d} + \mathbb{P}\left(N_{i,2n} \le \frac{n}{d} \text{ for some } 1 \le i \le d \right).$$

Thus, we will have proved that there is a constant $A(d) < \infty$ such that

$$\mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) \le A(d)n^{-\frac{d}{2}}, \quad n \ge 1,$$
(1.2.13)

once we show that there is a constant $B(d) < \infty$ such that

$$\mathbb{P}(N_{i,2n} \le \frac{n}{d} \text{ for some } 1 \le i \le d) \le B(d)n^{-\frac{d}{2}}, \quad n \ge 1.$$
(1.2.14)

In particular, this will complete the proof that

$$d \ge 3 \implies \sum_{n=0}^{\infty} \mathbb{P}(\mathbf{X}_{2n} = \mathbf{0}) < \infty$$

and therefore that the symmetric, nearest neighbor random walk in \mathbb{Z}^d is transient when $d \ge 3$.

To prove (1.2.14), first note that

$$\mathbb{P}\left(N_{i,2n} \leq \frac{n}{d} \text{ for some } 1 \leq i \leq d\right) \leq d\mathbb{P}\left(N_{1,2n} \leq \frac{n}{d}\right).$$

Next, set $Z_m = \mathbf{1}_{\{1\}}(I_m)$, write $N_{1,n} = \sum_{1}^{n} Z_m$, and observe that $\{Z_m : m \ge 1\}$ is a sequence of $\{0, 1\}$ -valued Bernoulli random variables such that $\mathbb{P}(Z_m = 1) =$

1 Random Walks, a Good Place to Begin

 $p \equiv \frac{1}{d}$. In particular, for any $\lambda \in \mathbb{R}$,

$$\mathbb{E}\left[\exp\left(\lambda\sum_{1}^{n}Z_{m}\right)\right] = \left(pe^{\lambda} + q\right)^{n},$$

and so

$$\mathbb{E}\left[\exp\left(\lambda\left(np-\sum_{1}^{n}Z_{m}\right)\right)\right]=e^{n\psi(\lambda)}\quad\text{where }\psi(\lambda)\equiv\log\left(pe^{-\lambda q}+qe^{\lambda p}\right).$$

Since $\psi(0) = \psi'(0) = 0$, and

$$\psi''(\lambda) = \frac{pqe^{\lambda(p-q)}}{(qe^{\lambda p} + pe^{-\lambda q})^2} = \frac{pq}{(qx_{\lambda} + px_{\lambda}^{-1})^2} \le \frac{1}{4}$$

where $x_{\lambda} \equiv e^{\frac{1}{2}\lambda(p+q)}$, Taylor's formula allows us to conclude that

$$\mathbb{E}\left[\exp\left(\lambda\left(np-\sum_{1}^{n}Z_{m}\right)\right)\right] \le e^{\frac{n\lambda^{2}}{8}}, \quad \lambda \in \mathbb{R}.$$
(1.2.15)

Starting from (1.2.15), there are many ways to arrive at (1.2.14). For example, for any $\lambda > 0$ and R > 0, Markov's inequality (7.1.12) plus (1.2.15) say that

$$\mathbb{P}\left(\sum_{1}^{n} Z_{m} \le np - nR\right) = \mathbb{P}\left(\exp\left(\lambda\left(np - \sum_{1}^{n} Z_{m}\right)\right) \ge e^{n\lambda R}\right)$$
$$\le e^{-\lambda nR + \frac{n\lambda^{2}}{8}},$$

which, when $\lambda = 4nR$, gives

$$\mathbb{P}\left(\sum_{1}^{n} Z_m \le np - nR\right) \le e^{-2nR^2}.$$
(1.2.16)

Returning to the notation used earlier and using the remark with which our discussion of (1.2.14) began, one sees from (1.2.16) that

$$\mathbb{P}\left(N_{i,2n} \leq \frac{n}{d} \text{ for some } 1 \leq i \leq d\right) \leq de^{-\frac{4n}{d^2}},$$

which is obviously far more than is required by (1.2.14).

1.3 Exercises

Exercise 1.3.1 Schwarz's inequality comes in many forms, the most elementary of which is the statement that, for any $\{a_n : n \in \mathbb{Z}\} \subseteq \mathbb{R}$ and $\{b_n : n \in \mathbb{Z}\} \subseteq \mathbb{R}$,

$$\sum_{n\in\mathbb{Z}}|a_nb_n|\leq \sqrt{\sum_{n\in\mathbb{Z}}a_n^2}\sqrt{\sum_{n\in\mathbb{Z}}b_n^2}.$$

Moreover, when the right hand side is finite, then

$$\left|\sum_{n\in\mathbb{Z}}a_nb_n\right| = \sqrt{\sum_{n\in\mathbb{Z}}a_n^2}\sqrt{\sum_{n\in\mathbb{Z}}b_n^2}$$

if and only if there is an $\alpha \in \mathbb{R}$ for which either $b_n = \alpha a_n$, $n \in \mathbb{Z}$, or $a_n = \alpha b_n$, $n \in \mathbb{Z}$. Here is an outline of one proof of these statements.

- (a) Given a real, quadratic polynomial $P(x) = Ax^2 + 2Bx + C$, use the quadratic formula to see that $P \ge 0$ everywhere if and only if $C \ge 0$ and $B^2 \le AC$. Similarly, show that P > 0 everywhere if and only if C > 0 and $B^2 < AC$.
- (b) Begin by noting that there is nothing to do unless $\sum_{n \in \mathbb{Z}} (a_n^2 + b_n^2) < \infty$, in which case $\sum_{n \in \mathbb{Z}} |a_n b_n| < \infty$ as well. Assuming that $\sum_{n \in \mathbb{Z}} (a_n^2 + b_n^2) < \infty$, show that

$$P(x) \equiv \sum_{n \in \mathbb{Z}} (a_n x + b_n)^2 = Ax^2 + 2Bx + C$$

where $A = \sum_{n \in \mathbb{Z}} a_n^2$, $B = \sum_{n \in \mathbb{Z}} b_n^2$, and $C = \sum_{n \in \mathbb{Z}} a_n b_n$

and apply (a).

Exercise 1.3.2 Let $\{Y_n : n \ge 1\}$ be a sequence of mutually independent, identically distributed random variables satisfying $\mathbb{E}[|Y_1|] < \infty$. Set $X_n = \sum_{m=1}^n Y_m$ for $n \ge 1$. *The weak law of large numbers* says that

$$\mathbb{P}\left(\left|\frac{X_n}{n} - \mathbb{E}[Y_1]\right| \ge \epsilon\right) \longrightarrow 0 \quad \text{for all } \epsilon > 0.$$

In fact,

$$\lim_{n \to \infty} \mathbb{E}\left[\left| \frac{X_n}{n} - \mathbb{E}[Y_1] \right| \right] = 0, \qquad (1.3.3)$$

from which the above follows as an application of Markov's inequality. Here are steps which lead to (1.3.3).

(a) First reduce to the case when $\mathbb{E}[Y_1] = 0$. Next, assume that $\mathbb{E}[Y_1^2] < \infty$, and show that

$$\mathbb{E}\left[\left|\frac{X_n}{n}\right|\right]^2 \le \mathbb{E}\left[\left|\frac{X_n}{n}\right|^2\right] = \frac{\mathbb{E}[Y_1^2]}{n}$$

Hence the result is proved when Y_1 has a finite second moment.

(b) Given R > 0, set $Y_n^{(R)} = Y_n \mathbf{1}_{[0,R)}(|Y_n|) - \mathbb{E}[Y_n, |Y_n| < R]$ and $X_n^{(R)} = \sum_{m=1}^n Y_m^{(R)}$. Note that, for any R > 0,

$$\mathbb{E}\left[\left|\frac{X_n}{n}\right|\right] \le \mathbb{E}\left[\left|\frac{X_n^{(R)}}{n}\right|\right] + \mathbb{E}\left[\left|\frac{X_n - X_n^{(R)}}{n}\right|\right]$$
$$\le \sqrt{\mathbb{E}\left[\left(\frac{X_n^{(R)}}{n}\right)^2\right]} + 2\mathbb{E}\left[|Y_1|, |Y_1| \ge R\right] \le \frac{R}{n^{\frac{1}{2}}} + 2\mathbb{E}\left[|Y_1|, |Y_1| \ge R\right].$$

and use this, together with the monotone convergence theorem, to complete the proof of (1.3.3).

Exercise 1.3.4 Refer to Exercise 1.3.2. *The strong law of large numbers* says that the statement in the weak law can be improved to the statement that $\frac{X_n}{n} \longrightarrow \mathbb{E}[Y_1]$ with probability 1. The proof of the strong law when one assumes only that $\mathbb{E}[|Y_1|] < \infty$ is a bit tricky. However, if one is willing to assume that $\mathbb{E}[Y_1^4] < \infty$, then a proof can be based on the same type argument that leads to the weak law.

Let $\{Y_n\}_1^\infty$ be a sequence of mutually independent random variables with the properties that $M = \sup_n \mathbb{E}[|Y_n|^4] < \infty$, and prove that, with probability 1, $\lim_{n\to\infty} \frac{1}{n} \sum_{m=1}^n (Y_m - \mathbb{E}[Y_m]) = 0$. Note that we have not assumed yet that they are identically distributed, but when we add this assumption we get $\lim_{n\to\infty} \frac{1}{n} \sum_{m=1}^n Y_m = \mathbb{E}[Y_1]$ with probability 1.

Here is an outline.

- (a) Begin by reducing to the case when $\mathbb{E}[Y_n] = 0$ for all $n \in \mathbb{Z}^+$.
- (b) After writing

$$\mathbb{E}\left[\left(\sum_{1}^{n} Y_{k}\right)^{4}\right] = \sum_{k_{1},\dots,k_{4}=1}^{n} \mathbb{E}[Y_{k_{1}}\cdots Y_{k_{4}}]$$

and noting that the only terms which do not vanish are those for which each index is equal to at least one other index, conclude that

$$\mathbb{E}\left[\left(\sum_{1}^{n} Y_{k}\right)^{4}\right] = \sum_{k=1}^{n} \mathbb{E}[Y_{k}^{4}] + 6 \sum_{1 \le k < \ell \le n} \mathbb{E}[Y_{k}^{2}]\mathbb{E}[Y_{\ell}^{2}].$$

Hence, since $\mathbb{E}[Y_k^4] - \mathbb{E}[Y_k^2]^2 = \operatorname{Var}(Y_k^2) \ge 0$,

$$\mathbb{E}\left[\left(\sum_{1}^{n} Y_{k}\right)^{4}\right] \leq 3Mn^{2}.$$
(*)

1.3 Exercises

(c) Starting from (*), show that

$$\mathbb{P}\left(\left|\frac{\sum_{1}^{n} Y_{k}}{n}\right| \ge \epsilon\right) \le \frac{1}{\epsilon^{4}} \mathbb{E}\left[\left|\frac{\sum_{1}^{n} Y_{k}}{n}\right|^{4}\right] \le \frac{3M}{\epsilon^{4} n^{2}} \longrightarrow 0$$

for all $\epsilon > 0$. This is the weak law of large numbers for independent random variables with bounded fourth moments. Of course, the use of four moments here is somewhat ridiculous since the argument using only two moments is easier.

(d) Starting again from (*) and using (7.1.4), show that

$$\mathbb{P}\left(\sup_{n>m} \left|\frac{\sum_{1}^{n} Y_{k}}{n}\right| \ge \epsilon\right) \le \sum_{n=m+1}^{\infty} \mathbb{P}\left(\left|\sum_{1}^{n} Y_{k}\right| \ge n\epsilon\right)$$
$$\le \frac{4M}{\epsilon^{4}} \sum_{n=m+1}^{\infty} \frac{1}{n^{2}}$$
$$\le \frac{4M}{\epsilon^{4}m} \longrightarrow 0 \quad \text{as } m \to \infty \text{ for all } \epsilon > 0$$

(e) Use the definition of convergence plus (7.1.4) to show that

$$\mathbb{P}\left(\frac{\sum_{1}^{n} Y_{k}}{n} \nrightarrow 0\right) = \mathbb{P}\left(\bigcup_{N=1}^{\infty} \bigcap_{m=1}^{\infty} \bigcup_{n>m} \left\{ \left|\frac{\sum_{1}^{n} Y_{k}}{n}\right| \ge \frac{1}{N} \right\} \right)$$
$$\leq \sum_{N=1}^{\infty} \mathbb{P}\left(\bigcap_{m=1}^{\infty} \bigcup_{n>m} \left\{ \left|\frac{\sum_{1}^{n} Y_{k}}{n}\right| \ge \frac{1}{N} \right\} \right).$$

Finally, apply the second line of (7.1.3) plus (d) above to justify

$$\mathbb{P}\left(\bigcap_{m=1}^{\infty}\bigcup_{n>m}\left\{\left|\frac{\sum_{1}^{n}Y_{k}}{n}\right|\geq\frac{1}{N}\right\}\right)=\lim_{m\to\infty}\mathbb{P}\left(\bigcup_{n>m}\left\{\left|\frac{\sum_{1}^{n}Y_{k}}{n}\right|\geq\frac{1}{N}\right\}\right)=0$$

for each $N \in \mathbb{Z}^+$. Hence, with probability $1, \frac{1}{n} \sum_{1}^{n} Y_k \longrightarrow 0$, which is the strong law of large numbers for independent random variables with bounded fourth moments.

Exercise 1.3.5 Readers who know DeMoivre's proof of the central limit theorem will have realized that the estimate in (1.2.11) is a poor man's substitute for what one can get as a consequence of *Stirling's formula*

$$n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \quad \text{as } n \to \infty,$$
 (1.3.6)

whose meaning is that the ratio of the quantities on the two sides of " \sim " tends to 1. Indeed, given (1.3.6), show that

$$\mathbb{P}(X_{2n}=0)\sim \sqrt{\frac{2}{\pi n}}.$$

Next, give a proof of (1.3.6) based on the following line of reasoning.

(a) Let τ_1, \ldots, τ_n be a mutually independent, unit exponential random variables,¹⁰ and show that for any $0 < R \le \sqrt{n}$

$$1 - \frac{1}{R^2} \le \mathbb{P}\left(-R \le \frac{\tau_1 + \dots + \tau_n - n}{\sqrt{n}} \le R\right) = \frac{1}{(n-1)!} \int_{-\sqrt{n}R+n}^{\sqrt{n}R+n} t^{n-1} e^{-t} dt.$$

(b) Make a change of variables followed by elementary manipulations to show that

$$\int_{-\sqrt{nR+n}}^{\sqrt{nR+n}} t^{n-1} e^{-t} dt = \sqrt{n} e^{-n} \int_{-R}^{R} (n + \sqrt{n}\sigma)^{n-1} e^{-\sqrt{n}\sigma} d\sigma$$
$$= n^{n-\frac{1}{2}} e^{-n} \int_{-R}^{R} \left(1 + \frac{\sigma}{\sqrt{n}}\right)^{n-1} e^{-\sqrt{n}\sigma} d\sigma$$
$$= n^{n-\frac{1}{2}} e^{-n} \int_{-R}^{R} \exp\left(-\frac{\sigma^2}{2} + E_n(\sigma)\right) d\sigma,$$

where

$$E_n(\sigma) \equiv (n-1)\log\left(1+\frac{\sigma}{\sqrt{n}}\right) - \sqrt{n\sigma} + \frac{\sigma^2}{2}$$

(c) As an application of the Taylor's series for $\log(1 + x)$ (cf. (1.2.10)), show that $E_n(\sigma) \longrightarrow 0$ uniformly for $|\sigma| \le R$ when $n \to \infty$, and combine this with the results in (a) and (b) to arrive at

$$\overline{\lim_{n \to \infty}} \frac{n^{n+\frac{1}{2}}e^{-n}}{n!} \int_{-R}^{R} e^{-\frac{\sigma^2}{2}} d\sigma \le 1$$

and

$$\lim_{n \to \infty} \frac{n^{n+\frac{1}{2}}e^{-n}}{n!} \int_{-R}^{R} e^{-\frac{\sigma^2}{2}} d\sigma \ge 1 - \frac{1}{R^2}.$$

Because $\int_{-\infty}^{\infty} e^{-\frac{\sigma^2}{2}} d\sigma = \sqrt{2\pi}$, it is clear that (1.3.6) follows after one lets $R \nearrow \infty$.

¹⁰A unit exponential random variable is a random variable τ for which $\mathbb{P}(\tau > t) = e^{-t \vee 0}$.

1.3 Exercises

Exercise 1.3.7 The argument in Sect. 1.2.3 is quite robust. Indeed, let $\{\mathbf{X}_n : n \ge 0\}$ be any symmetric random walk on \mathbb{Z}^2 whose jumps have finite second moment. That is, $\mathbf{X}_0 = \mathbf{0}$, $\{\mathbf{X}_n - \mathbf{X}_{n-1} : n \ge 1\}$ are mutually independent, identically distributed, symmetric (\mathbf{X}_1 has the same distribution as $-\mathbf{X}_1$), \mathbb{Z}^2 -valued random variables with finite second moment. Show that $\{\mathbf{X}_n : n \ge 0\}$ is recurrent in the sense that $\mathbb{P}(\exists n \ge 1 \mathbf{X}_n = \mathbf{0}) = 1$.

Exercise 1.3.8 Let $\{\mathbf{X}_n : n \ge 0\}$ be a random walk on $\mathbb{Z}^d : \mathbf{X}_0 = \mathbf{0}, \{\mathbf{X}_n - \mathbf{X}_{n-1} : n \ge 1\}$ are mutually independent, identically distributed, \mathbb{Z}^d -valued random variables. Further, for each $1 \le i \le d$, let $(\mathbf{X}_n)_i$ be the *i*th coordinate of \mathbf{X}_n , and assume that

$$\min_{1\leq i\leq d} \mathbb{P}((\mathbf{X}_1)_i\neq 0) > 0 \quad \text{but} \quad \mathbb{P}(\exists i\neq j \ (\mathbf{X}_1)_i (\mathbf{X}_1)_j\neq 0) = 0.$$

If, for some $C < \infty$ and $(\alpha_1, \ldots, \alpha_d) \in [0, \infty)^d$ with $\sum_{i=1}^d \alpha_i > 1$,

$$\mathbb{P}((\mathbf{X}_n)_i = 0) \le Cn^{-\alpha_i} \quad \text{for all } n \ge 1,$$

show that $\{\mathbf{X}_n : n \ge 0\}$ is transient in the sense that $\mathbb{P}(\exists n \ge 1 | \mathbf{X}_n = \mathbf{0}) < 1$.

Exercise 1.3.9 As in the preceding, let $\{\mathbf{X}_n : n \ge 0\}$ be a random walk on \mathbb{Z}^d . Given $\mathbf{k} \in \mathbb{Z}^d$, set

$$T_{\mathbf{k}} = \sum_{n=0}^{\infty} \mathbf{1}_{\{\mathbf{k}\}}(\mathbf{X}_n) \quad \text{and} \quad \zeta^{\{\mathbf{k}\}} = \inf\{n \ge 0 : \mathbf{X}_n = \mathbf{k}\}.$$

Show that

$$\mathbb{E}[T_{\mathbf{k}}] = \mathbb{P}(\zeta^{\{\mathbf{k}\}} < \infty) \mathbb{E}[T_{\mathbf{0}}] = \frac{\mathbb{P}(\zeta^{\{\mathbf{k}\}} < \infty)}{\mathbb{P}(\rho_{\mathbf{0}} = \infty)}, \qquad (1.3.10)$$

where $\rho_0 = \inf\{n \ge 1 : \mathbf{X}_n = \mathbf{0}\}$ is the time of first return to **0**. In particular, if $\{\mathbf{X}_n : n \ge 0\}$ is transient in the sense described in the preceding exercise, show that

$$\mathbb{E}\left[\sum_{m=0}^{\infty} \mathbf{1}_{B(r)}(\mathbf{X}_n)\right] < \infty \quad \text{for all } r \in (0,\infty),$$

where $B(r) = \{\mathbf{k} : |\mathbf{k}| \le r\}$; and from this conclude that $|\mathbf{X}_n| \longrightarrow \infty$ with probability 1. On the other hand, if $\{\mathbf{X}_n : n \ge 0\}$ is recurrent, show that $\mathbf{X}_n = \mathbf{0}$ infinitely often with probability 1. Hence, either $\{\mathbf{X}_n : n \ge 0\}$ is recurrent and $\mathbf{X}_n = \mathbf{0}$ infinitely often with probability 1 or it is transient and $|\mathbf{X}_n| \longrightarrow \infty$ with probability 1.

Exercise 1.3.11 Take d = 1 in the preceding, $X_0 = 0$, and $\{X_n - X_{n-1} : n \ge 1\}$ to be mutually independent, identically distributed random variables for which $0 < \mathbb{E}[|X_1|] < \infty$ and $\mathbb{E}[X_1] = 0$. By a slight variation on the argument given in Sect. 1.2.1, we will show here that this random walk is not only recurrent but that

$$\overline{\lim_{n \to \infty}} X_n = \infty \quad \text{and} \quad \underline{\lim_{n \to \infty}} X_n = -\infty \quad \text{with probability 1.}$$

(a) First show that it suffices to prove that $\sup_n X_n = \infty$ and that $\inf_n X_n = -\infty$. Next, use (1.3.3) to show that

$$\lim_{n\to\infty}\max_{1\le m\le n}\frac{\mathbb{E}[|X_m|]}{n}=0.$$

(b) For $n \ge 1$, set $T_k^{(n)} = \sum_{m=0}^{n-1} \mathbf{1}_{\{k\}}(X_m)$, show that $\mathbb{E}[T_k^{(n)}] \le \mathbb{E}[T_0^{(n)}]$ for all $k \in \mathbb{Z}$, and use this to arrive at

$$(4\mu(n)+1)\mathbb{E}[T_0^{(n)}] \ge \frac{n}{2}$$
 where $\mu(n) \equiv \max_{0 \le m \le n-1} \mathbb{E}[|X_m|]$.

Finally, apply part (a) to see that $\mathbb{E}[T_0] = \infty$ and (1.3.10) to get $\mathbb{P}(\rho_0 < \infty) = 1$, which means that $\{X_n : n \ge 0\}$ is recurrent.

(c) To complete the program, proceed as in the derivation of (1.2.2) to pass from (b) to

$$\mathbb{P}(\rho_0^{(m)} < \infty) = 1 \quad \text{for all } m \ge 1, \tag{(*)}$$

where $\rho_0^{(m)}$ is the time of the *m*th return to 0. Next, set $\eta_r = \inf\{n \ge 0 : X_n \ge r\}$ for $r \in \mathbb{Z}^+$, show that $\epsilon \equiv \mathbb{P}(\eta_1 > \rho_0) < 1$, and conclude that $\mathbb{P}(\eta_1 > \rho_0^{(m)}) \le \epsilon^m$. Now, combine this with (*) to get $\mathbb{P}(\eta_1 < \infty) = 1$. Finally, argue that

$$\mathbb{P}(\eta_{r+1} < \infty) \ge \mathbb{P}(\eta_r < \infty) \mathbb{P}(\eta_1 < \infty)$$

and therefore that $\mathbb{P}(\eta_r < \infty) = 1$ for all $r \ge 1$. Since this means that, with probability 1, $\sup_n X_n \ge r$ for all $r \ge 1$, it follows that $\sup_n X_n = \infty$ with probability 1. To prove that $\inf_n X_n = -\infty$ with probability 1, simply replace $\{X_n : n \ge 0\}$ by $\{-X_n : n \ge 0\}$.

Exercise 1.3.12¹¹ Here is an interesting application of one dimensional random walks to elementary *queuing theory*. Queuing theory deals with the distribution of the number of people waiting to be served (i.e., the length of the queue) when, during each time interval, the number of people who arrive and the number of people who are served are random. The queuing model which we will consider here is among the simplest. Namely, we will assume that the queue is initially empty and that during each time interval [n - 1, n) the number of people who arrive minus the number who can be served is given by a \mathbb{Z} -valued random variable B_n . Further, we assume that the B_n 's are mutually independent and identically distributed random variables satisfying $0 < \mathbb{E}[|B_1|] < \infty$. The associated queue is, apart from the fact that there are never a negative number of people waiting, the random walk $\{X_n : n \ge 0\}$ determined by the B_n 's: $X_0 = 0$ and $X_n = \sum_{m=1}^n B_m$. To take into account the prohibition against having a queue of negative length, the queuing model $\{Q_n : n \ge 0\}$ is given by the prescription

$$Q_0 = 0$$
 and $Q_n = (Q_{n-1} + B_n)^+$ for $n \ge 1$.

¹¹So far as I know, this example was invented by Wm. Feller.

1.3 Exercises

(a) Show that

$$Q_n = X_n - \min_{0 \le m \le n} X_m = \max_{0 \le m \le n} (X_n - X_m),$$

and conclude that, for each $n \ge 0$, the distribution of Q_n is the same as that of $M_n \equiv \max_{0 \le m \le n} X_m$.

(b) Set $M_{\infty} \equiv \lim_{n \to \infty} M_n \in \mathbb{N} \cup \{\infty\}$, and, as a consequence of (a), arrive at

$$\lim_{n \to \infty} \mathbb{P}(Q_n = j) = \mathbb{P}(M_\infty = j) \quad \text{for } j \in \mathbb{N}.$$

- (c) Set $\mu \equiv \mathbb{E}[B_1]$. The weak law of large numbers says that, for each $\epsilon > 0$, $\mathbb{P}(|X_n - n\mu| \ge n\epsilon) \longrightarrow 0$ as $n \to \infty$. In particular, when $\mu > 0$, show that $\mathbb{P}(M_{\infty} = \infty) = 1$. When $\mu = 0$, use Exercise 1.3.11 to reach the same conclusion. Hence, when $\mathbb{E}[B_1] \ge 0$, $\mathbb{P}(Q_n = j) \longrightarrow 0$ for all $j \in \mathbb{N}$. That is, when the expected difference between the number arrivals and the number of people served is non-negative, then $\mathbb{P}(Q_n \ge M) \longrightarrow 1$ for every $M \in \mathbb{Z}^+$.
- (d) Now assume that $\mu \equiv \mathbb{E}[B_1] < 0$. Then the strong law of large numbers (cf. Exercise 1.3.4 for the case when B_1 has a finite fourth moment and Theorem 3.3.10 in [8] for the general case) says that $\frac{X_n}{n} \longrightarrow \mu$ with probability 1. In particular, conclude that $M_{\infty} < \infty$ with probability 1 and therefore that $\sum_{j \in \mathbb{N}} v_j = 1$ when $v_j \equiv \lim_{n \to \infty} \mathbb{P}(Q_n = j) = \mathbb{P}(M_{\infty} = j)$.
- (e) Specialize to the case when the B_m 's are $\{-1, 1\}$ -valued Bernoulli random variables with $p \equiv \mathbb{P}(B_1 = 1) \in (0, 1)$, and set q = 1 p. Use the calculations in (1.1.12) to show that

$$\lim_{n \to \infty} \mathbb{P}(Q_n = j) = \begin{cases} 0 & \text{if } p \ge q \\ \frac{q-p}{q} (\frac{p}{q})^j & \text{if } p < q. \end{cases}$$

(f) Generalize (e) to the case when $B_m \in \{-1, 0, 1\}$, $p = \mathbb{P}(B_1 = 1)$, and $q = \mathbb{P}(B_1 = -1)$. The idea is that M_{∞} in this case has the same distribution as $\sup_n Y_n$, where $\{Y_n : n \ge 0\}$ is the random walk corresponding to $\{-1, 1\}$ -valued Bernoulli random variables which are 1 with probability $\frac{p}{p+q}$.

Chapter 2 Doeblin's Theory for Markov Chains

In this chapter we begin in earnest our study of Markov processes. Like the random walks in Chap. 1, the processes with which we will be dealing here take only countably many values and have a discrete (as opposed to continuous) time parameter. In fact, in many ways, these processes are the simplest generalizations of random walks. To be precise, random walks proceed in such a way that the distribution of their increments are independent of everything that has happened before the increment takes place. The processes at which we will be looking now proceed in such a way that the distribution of their new position depends on where they are at the time when they move but not on where they were in the past. A process with this sort of dependence property is said to have the *Markov property* and is called a *Markov chain*.¹

The set S in which a process takes its values is called its *state space*, and, as I said, our processes will have state spaces which are either finite or countably infinite. Thus, at least for theoretical purposes, there is no reason for us not to think of S as the set $\{1, \ldots, N\}$ or \mathbb{Z}^+ , depending on whether S is finite or countably infinite. On the other hand, always taking S to be one of these has the disadvantage that it may mask important properties. For example, it would have been a great mistake to describe the nearest neighbor random walk on \mathbb{Z}^2 after mapping \mathbb{Z}^2 isomorphically onto \mathbb{Z}^+ .

2.1 Some Generalities

Before getting started, there are a few general facts that we will need to know about Markov chains.

A *Markov chain* on a finite or countably infinite state space S is a family of S-valued random variables $\{X_n : n \ge 0\}$ with the property that, for all $n \ge 0$ and

¹The term "chain" is commonly applied to processes with a time discrete parameter.

 $(i_0,\ldots,i_n,j)\in\mathbb{S}^{n+2},$

$$\mathbb{P}(X_{n+1} = j \mid X_0 = i_0, \dots, X_n = i_n) = (\mathbf{P})_{i_n j},$$
(2.1.1)

where **P** is a matrix all of whose entries are non-negative and each of whose rows sums to 1. Equivalently (cf. Sect. 7.4.1)

$$\mathbb{P}(X_{n+1} = j \mid X_0, \dots, X_n) = (\mathbf{P})_{X_n j}.$$
(2.1.2)

It should be clear that (2.1.2) is a mathematically precise expression of the idea that, when a Markov chain jumps, the distribution of where it lands depends only on where it was at the time when it jumped and not on where it was in the past.

2.1.1 Existence of Markov Chains

For obvious reasons, a matrix whose entries are non-negative and each of whose rows sum to 1 is called a *transition probability matrix*: it gives the probability that the Markov chain will move to the state *j* at time n + 1 given that it is at state *i* at time *n*, independent of where it was prior to time *n*. Further, it is clear that only a transition probability matrix could appear on the right of (2.1.1). What may not be so immediate is that one can go in the opposite direction. Namely, let μ be a *probability vector*² and **P** a transition probability matrix. Then there exists a Markov chain { $X_n : n \ge 0$ } with *initial distribution* μ and transition probability matrix **P**. That is, $\mathbb{P}(X_0 = i) = (\mu)_i$ and (2.1.1) holds.

To prove the preceding existence statement, one can proceed as follows. Begin by assuming, without loss in generality, that \mathbb{S} is either $\{1, \ldots, N\}$ or \mathbb{Z}^+ . Next, given $i \in \mathbb{S}$, set $\beta(i, 0) = 0$ and $\beta(i, j) = \sum_{k=1}^{j} (\mathbf{P})_{ik}$ for $j \ge 1$, and define $F : \mathbb{S} \times [0, 1) \longrightarrow \mathbb{S}$ so that F(i, u) = j if $\beta(i, j-1) \le u < \beta(i, j)$. In addition, set $\alpha(0) = 0$ and $\alpha(i) = \sum_{k=1}^{i} (\mu)_k$ for $i \ge 1$, and define $f : [0, 1) \longrightarrow \mathbb{S}$ so that f(u) = i if $\alpha(i-1) \le u < \alpha(i)$. Finally, let $\{U_n : n \ge 0\}$ be a sequence of mutually independent random variables (cf. Theorem 7.3.2) which are uniformly distributed on [0, 1), and set

$$X_n = \begin{cases} f(U_0) & \text{if } n = 0\\ F(X_{n-1}, U_n) & \text{if } n \ge 1. \end{cases}$$
(2.1.3)

We will now show that the sequence $\{X_n : n \ge 0\}$ in (2.1.3) is a Markov chain with the required properties. For this purpose, suppose that $(i_0, \ldots, i_n) \in \mathbb{S}^{n+1}$, and

²A probability vector is a row vector whose coordinates are non-negative and sum to 1.

observe that

$$\mathbb{P}(X_0 = i_0, \dots, X_n = i_n)$$

= $\mathbb{P}(U_0 \in [\alpha(i_0 - 1), \alpha(i_0))$
& $U_m \in [\beta(i_{m-1}, i_m - 1), \beta(i_{m-1}, i_m)) \text{ for } 1 \le m \le n)$
= $\mu_{i_0}(\mathbf{P})_{i_0i_1} \cdots (\mathbf{P})_{i_{n-1}i_n} = \mathbb{P}(X_0 = i_0, \dots, X_{n-1} = i_{n-1})\mathbf{P}_{i_{n-1},j}.$

2.1.2 Transition Probabilities & Probability Vectors

Notice that the use of matrix notation here is clever. To wit, if μ is the row vector with *i*th entry $(\mu)_i = \mathbb{P}(X_0 = i)$, then μ is called the *initial distribution* of the chain and

$$(\boldsymbol{\mu}\mathbf{P}^n)_j = \mathbb{P}(X_n = j), \quad n \ge 0 \text{ and } j \in \mathbb{S},$$
 (2.1.4)

where we have adopted the convention that \mathbf{P}^0 is the identity matrix and $\mathbf{P}^n = \mathbf{P}\mathbf{P}^{n-1}$ for $n \ge 1$.³ To check (2.1.4), let $n \ge 1$ be given, and note that, by (2.1.1) and induction,

$$\mathbb{P}(X_0 = i_0, \dots, X_{n-1} = i_{n-1}, X_n = j) = (\boldsymbol{\mu})_{i_0}(\mathbf{P})_{i_0 i_1} \cdots (\mathbf{P})_{i_{n-1} j}$$

Hence (2.1.4) results after one sums with respect to (i_0, \ldots, i_{n-1}) . Obviously, (2.1.4) is the statement that the row vector $\mu \mathbf{P}^n$ is the distribution of the Markov chain at time *n* if μ is its initial distribution (i.e., its distribution at time 0). Alternatively, \mathbf{P}^n is the *n*-step transition probability matrix: $(\mathbf{P}^n)_{ij}$ is the conditional probability that $X_{m+n} = j$ given that $X_m = i$.

For future reference, we will introduce here an appropriate way in which to measure the length of row vectors when they are being used to represent measures. Namely, given a row vector ρ , we set

$$\|\boldsymbol{\rho}\|_{\mathbf{v}} = \sum_{i \in \mathbb{S}} |(\boldsymbol{\rho})_i|, \qquad (2.1.5)$$

where the subscript "v" is used in recognition that this is the notion of length which corresponds to the *variation norm* on the space of measures. The basic reason for our making this choice of norm is that

$$\|\boldsymbol{\rho}\mathbf{P}\|_{\mathsf{v}} \le \|\boldsymbol{\rho}\|_{\mathsf{v}},\tag{2.1.6}$$

since, by Theorem 7.1.15,

$$\|\boldsymbol{\rho}\mathbf{P}\|_{\mathbf{v}} = \sum_{j\in\mathbb{S}} \left| \sum_{i\in\mathbb{S}} (\boldsymbol{\rho})_i (\mathbf{P})_{ij} \right| \leq \sum_{i\in\mathbb{S}} \left(\sum_{j\in\mathbb{S}} |(\boldsymbol{\rho})_i| (\mathbf{P})_{ij} \right) = \|\boldsymbol{\rho}\|_{\mathbf{v}}.$$

³The reader should check for itself that \mathbf{P}^n is again a transition probability matrix for all $n \in \mathbb{N}$: all entries are non-negative and each row sums to 1.

Notice that this is a quite different way of measuring the length from the way Euclid would have: he would have used

$$\|\boldsymbol{\rho}\|_2 = \left(\sum_{i\in\mathbb{S}} (\boldsymbol{\rho})_i^2\right)^{\frac{1}{2}}.$$

When S is finite, these two norms are comparable. Namely,

$$\|\boldsymbol{\rho}\|_2 \le \|\boldsymbol{\rho}\|_v \le \sqrt{\#\mathbb{S}} \|\boldsymbol{\rho}\|_2$$
, where $\#\mathbb{S}$ denotes the cardinality of \mathbb{S} .

The first inequality is easily seen by squaring both sides, and the second is an application of Schwarz's inequality (cf. Exercise 1.3.1). However, when S is infinite, they are not comparable. Nonetheless, $\|\cdot\|_{v}$ is a good *norm* (i.e., measure of length) in the sense that $\|\rho\|_{v} = 0$ if and only if $\rho = 0$ and that it satisfies the *triangle inequality*: $\|\rho + \rho'\|_{v} \le \|\rho\|_{v} + \|\rho'\|_{v}$. Furthermore, Cauchy's convergence criterion holds for $\|\cdot\|_{v}$. That is, if $\{\rho_{n}\}_{1}^{\infty}$ is a sequence in $\mathbb{R}^{\mathbb{S}}$, then there exists $\rho \in \mathbb{R}^{\mathbb{S}}$ for which $\|\rho_{n} - \rho\|_{v} \longrightarrow 0$ if and only $\{\rho_{n}\}_{1}^{\infty}$ is *Cauchy convergent* in the sense that

$$\lim_{m\to\infty}\sup_{n>m}\|\boldsymbol{\rho}_n-\boldsymbol{\rho}_m\|_{\mathrm{v}}=0.$$

As usual, the "only if" direction is an easy application of the triangle inequality:

$$\|\boldsymbol{\rho}_n - \boldsymbol{\rho}_m\|_{\mathsf{v}} \leq \|\boldsymbol{\rho}_n - \boldsymbol{\rho}\|_{\mathsf{v}} + \|\boldsymbol{\rho} - \boldsymbol{\rho}_m\|_{\mathsf{v}}.$$

To go the other direction, suppose that $\{\rho_n\}_1^\infty$ is Cauchy convergent, and observe that each coordinate of $\{\rho_n\}_1^\infty$ must be Cauchy convergent as real numbers. Hence, by Cauchy's criterion for real numbers, there exists a ρ to which $\{\rho_n\}_1^\infty$ converges in the sense that each coordinate of the ρ_n 's tends to the corresponding coordinate of ρ . Thus, by Fatou's lemma, Theorem 7.1.10, as $m \to \infty$,

$$\|\boldsymbol{\rho}-\boldsymbol{\rho}_m\|_{\mathbf{v}}=\sum_{i\in\mathbb{S}}\left|(\boldsymbol{\rho})_i-(\boldsymbol{\rho}_m)_i\right|\leq \lim_{n\to\infty}\sum_{i\in\mathbb{S}}\left|(\boldsymbol{\rho}_n)_i-(\boldsymbol{\rho}_m)_i\right|\longrightarrow 0.$$

2.1.3 Transition Probabilities and Functions

As we saw in Sect. 2.1.2, the representation of the transition probability as a matrix and the initial distributions as a row vector facilitates the representation of the distribution at later times. In order to understand how to get the analogous benefit when computing expectation values of functions, think of a function f on the state space S as the column vector **f** whose *j*th coordinate is the value of the function f at *j*. Clearly, if μ is the row vector which represents the probability measure μ on S and **f** is the column vector which represents a function f which is either non-negative or bounded, then $\mu \mathbf{f} = \sum_{i \in S} f(i)\mu(\{i\})$ is the expected value of f with respect to μ . Similarly, the column vector $\mathbf{P}^n \mathbf{f}$ represents that function whose value at *i* is the conditional expectation value of $f(X_n)$ given that $X_0 = i$. Indeed,

$$\mathbb{E}\left[f(X_n) \mid X_0 = i\right] = \sum_{j \in \mathbb{S}} f(j)\mathbb{P}(X_n = j \mid X_0 = i)$$
$$= \sum_{j \in \mathbb{S}} (\mathbf{P}^n)_{ij} (\mathbf{f})_j = (\mathbf{P}^n \mathbf{f})_i.$$

More generally, if f is either a non-negative or bounded function on S and **f** is the column vector which it determines, then, for $0 \le m \le n$,

$$\mathbb{E}[f(X_n) | X_0 = i_0, \dots, X_m = i_m] = (\mathbf{P}^{n-m} \mathbf{f})_{i_m}, \quad \text{or, equivalently,}$$
$$\mathbb{E}[f(X_n) | X_0, \dots, X_m] = (\mathbf{P}^{n-m} \mathbf{f})_{X_m}$$
(2.1.7)

since

$$\mathbb{E}\Big[f(X_n) \mid X_0 = i_0, \dots, X_m = i_m\Big]$$

= $\sum_{j \in \mathbb{S}} f(j) \mathbb{P}(X_n = j \mid X_0 = i_0, \dots, X_m = i_m)$
= $\sum_{j \in \mathbb{S}} f(j) (\mathbf{P}^{n-m})_{i_m j} = (\mathbf{P}^{n-m} \mathbf{f})_{i_m}.$

In particular, if μ is the initial distribution of $\{X_n : n \ge 0\}$, then

$$\mathbb{E}\big[f(X_n)\big] = \boldsymbol{\mu} \mathbf{P}^n \mathbf{f},$$

since $\mathbb{E}[f(X_n)] = \sum_i (\boldsymbol{\mu})_i \mathbb{E}[f(X_n)|X_0 = i].$

Just as $\|\cdot\|_v$ was the appropriate way to measure the length of row vectors when we were using them to represent measures, the appropriate way to measure the length of column vectors which represent functions is with the *uniform norm* $\|\cdot\|_u$:

$$\|\mathbf{f}\|_{\mathbf{u}} = \sup_{j \in \mathbb{S}} |(\mathbf{f})_j|.$$
(2.1.8)

The reason why $\|\cdot\|_{u}$ is the norm of choice here is that $\|\boldsymbol{\mu}\mathbf{f}\| \le \|\boldsymbol{\mu}\|_{v} \|\mathbf{f}\|_{u}$, since

$$|\boldsymbol{\mu}\mathbf{f}| \leq \sum_{i\in\mathbb{S}} |(\boldsymbol{\mu})_i| |(\mathbf{f})_i| \leq \|\mathbf{f}\|_{\mathrm{u}} \sum_{i\in\mathbb{S}} |(\boldsymbol{\mu})_i|.$$

In particular, we have the complement to (2.1.6):

$$\|\mathbf{Pf}\|_{u} \le \|\mathbf{f}\|_{u}. \tag{2.1.9}$$

2.1.4 The Markov Property

By definition, if μ is the initial distribution of $\{X_n : n \ge 0\}$, then

$$\mathbb{P}(X_0 = i_0, \dots, X_n = i_n) = (\boldsymbol{\mu})_{i_0} (\mathbf{P})_{i_0 i_1} \cdots (\mathbf{P})_{i_{n-1} i_n}.$$
(2.1.10)

Hence, if $m, n \ge 1$ and $F : \mathbb{S}^{n+1} \longrightarrow \mathbb{R}$ is either bounded or non-negative, then

$$\mathbb{E} \Big[F(X_m, \dots, X_{m+n}), X_0 = i_0, \dots, X_m = i_m \Big]$$

= $\sum_{j_1, \dots, j_n \in \mathbb{S}} F(i_m, j_1, \dots, j_n) \mu_{i_0}(\mathbf{P})_{i_0 i_1} \cdots (\mathbf{P})_{i_{m-1} i_m} (\mathbf{P})_{i_m j_1} \cdots (\mathbf{P})_{j_{1n-1} j_n}$
= $\mathbb{E} \Big[F(X_0, \dots, X_n) \, \big| \, X_0 = i_m \Big] \mathbb{P}(X_0 = i_0, \dots, X_m = i_m).$

Equivalently, we have now proved the Markov property in the form

$$\mathbb{E} \Big[F(X_m, \dots, X_{m+n}) \, \big| \, X_0 = i_0, \dots, X_m = i_m \Big] \\= \mathbb{E} \Big[F(X_0, \dots, X_n) \, \big| \, X_0 = i_m \Big].$$
(2.1.11)

2.2 Doeblin's Theory

In this section I will introduce an elementary but basic technique, due to Doeblin, which will allow us to study the long time distribution of a Markov chain, particularly ones on a finite state space.

2.2.1 Doeblin's Basic Theorem

For many purposes, what one wants to know about a Markov chain is its distribution after a long time, and, at least when the state space is finite, it is reasonable to think that the distribution of the chain will stabilize. To be more precise, if one is dealing with a chain which can go in a single step from any state *i* to some state *j* with positive probability, then that state *j* is going to visited again and again, and so, after a while, the chain's initial distribution is going to get "forgotten." In other words, we are predicting for such a chain that $\mu \mathbf{P}^n$ will, for sufficiently large *n*, be nearly independent of μ . In particular, this would mean that $\mu \mathbf{P}^n = (\mu \mathbf{P}^{n-m})\mathbf{P}^m$ is very nearly equal to $\mu \mathbf{P}^m$ when *m* is large and therefore, by Cauchy's convergence criterion, that $\pi = \lim_{n\to\infty} \mu \mathbf{P}^n$ exists. In addition, if this were the case, then we would have that $\pi = \lim_{n\to\infty} \mu \mathbf{P}^{n+1} = \lim_{n\to\infty} (\mu \mathbf{P}^n)\mathbf{P} = \pi \mathbf{P}$. That is, π would have to be a left eigenvector for **P** with eigenvalue 1. A probability vector π is, for obvious reasons, called a *stationary probability* for the transition probability matrix **P** if $\pi = \pi \mathbf{P}$.

2.2 Doeblin's Theory

Although a state j of the sort in the preceding discussion is most likely to exist when the state space is finite, there are situations in which these musings apply even to infinite state spaces. That is, if, no matter where the chain starts, it has a positive probability of immediately visiting some fixed state, then, as the following theorem shows, it will stabilize.

Theorem 2.2.1 (Doeblin's Theorem) Let **P** be a transition probability matrix with the property that, for some state $j_0 \in \mathbb{S}$ and $\epsilon > 0$, $(\mathbf{P})_{ij_0} \ge \epsilon$ for all $i \in \mathbb{S}$. Then **P** has a unique stationary probability vector $\boldsymbol{\pi}$, $(\boldsymbol{\pi})_{j_0} \ge \epsilon$, and, for all initial distributions $\boldsymbol{\mu}$,

$$\left\|\boldsymbol{\mu}\mathbf{P}^n - \boldsymbol{\pi}\right\|_{\mathbf{v}} \le (1-\epsilon)^n \|\boldsymbol{\mu} - \boldsymbol{\pi}\|_{\mathbf{v}} \le 2(1-\epsilon)^n, \quad n \ge 0$$

Proof The key to the proof lies in the observations that if $\rho \in \mathbb{R}^{\mathbb{S}}$ is a row vector with $\|\rho\|_{V} < \infty$, then

$$\sum_{j \in \mathbb{S}} (\boldsymbol{\rho} \mathbf{P})_j = \sum_{i \in \mathbb{S}} (\boldsymbol{\rho})_i,$$

and
$$\sum_{i \in \mathbb{S}} (\boldsymbol{\rho})_i = 0 \implies \|\boldsymbol{\rho} \mathbf{P}^n\|_{\mathbf{v}} \le (1 - \epsilon)^n \|\boldsymbol{\rho}\|_{\mathbf{v}} \text{ for } n \ge 1.$$

$$(2.2.2)$$

The first of these is trivial, because, by Theorem 7.1.15,

$$\sum_{j\in\mathbb{S}} (\boldsymbol{\rho}\mathbf{P})_j = \sum_{j\in\mathbb{S}} \left(\sum_{i\in\mathbb{S}} (\boldsymbol{\rho})_i (\mathbf{P})_{ij} \right) = \sum_{i\in\mathbb{S}} \left(\sum_{j\in\mathbb{S}} (\boldsymbol{\rho})_i (\mathbf{P})_{ij} \right) = \sum_{i\in\mathbb{S}} (\boldsymbol{\rho})_i.$$

As for the second, note that, by an easy induction argument, it suffices to check it when n = 1. Next, suppose that $\sum_{i} (\rho)_i = 0$, and observe that

$$\left| (\boldsymbol{\rho} \mathbf{P})_{j} \right| = \left| \sum_{i \in \mathbb{S}} (\boldsymbol{\rho})_{i} (\mathbf{P})_{ij} \right|$$
$$= \left| \sum_{i \in \mathbb{S}} (\boldsymbol{\rho})_{i} ((\mathbf{P})_{ij} - \epsilon \delta_{j,j_{0}}) \right| \le \sum_{i \in \mathbb{S}} |(\boldsymbol{\rho})_{i}| ((\mathbf{P})_{ij} - \epsilon \delta_{j,j_{0}}),$$

and therefore that

$$\|\boldsymbol{\rho}\mathbf{P}\|_{\mathbf{v}} \leq \sum_{j\in\mathbb{S}} \left(\sum_{i\in\mathbb{S}} |(\boldsymbol{\rho})_i| \left((\mathbf{P})_{ij} - \epsilon \delta_{j,j_0} \right) \right)$$
$$= \sum_{i\in\mathbb{S}} |(\boldsymbol{\rho})_i| \left(\sum_{j\in\mathbb{S}} \left((\mathbf{P})_{ij} - \epsilon \delta_{j,j_0} \right) \right) = (1-\epsilon) \|\boldsymbol{\rho}\|_{\mathbf{v}}.$$

Now let $\boldsymbol{\mu}$ be a probability vector, and set $\boldsymbol{\mu}_n = \boldsymbol{\mu} \mathbf{P}^n$. Then, because $\boldsymbol{\mu}_n = \boldsymbol{\mu}_{n-m} \mathbf{P}^m$ and $\sum_i ((\boldsymbol{\mu}_{n-m})_i - \boldsymbol{\mu}_i) = 1 - 1 = 0$,

$$\|\boldsymbol{\mu}_n - \boldsymbol{\mu}_m\|_{\mathsf{v}} \leq (1-\epsilon)^m \|\boldsymbol{\mu}_{n-m} - \boldsymbol{\mu}\|_{\mathsf{v}} \leq 2(1-\epsilon)^m$$

for $1 \le m < n$. Hence $\{\mu_n\}_1^\infty$ is Cauchy convergent, and therefore there exists a π for which $\|\mu_n - \pi\|_v \longrightarrow 0$. Since each μ_n is a probability vector, it is clear that π must also be a probability vector. In addition, $\pi = \lim_{n\to\infty} \mu \mathbf{P}^{n+1} = \lim_{n\to\infty} (\mu \mathbf{P}^n) \mathbf{P} = \pi \mathbf{P}$, and so π is stationary. In particular,

$$(\boldsymbol{\pi})_{j_0} = \sum_{i \in \mathbb{S}} (\boldsymbol{\pi})_i (\mathbf{P})_{ij_0} \ge \epsilon \sum_{i \in \mathbb{S}} (\boldsymbol{\pi})_i = \epsilon.$$

Finally, if v is any probability vector, then

$$\left\|\boldsymbol{\nu}\mathbf{P}^m-\boldsymbol{\pi}\right\|_{\mathbf{v}}=\left\|(\boldsymbol{\nu}-\boldsymbol{\pi})\mathbf{P}^m\right\|_{\mathbf{v}}\leq(1-\epsilon)^m\|\boldsymbol{\mu}-\boldsymbol{\pi}\|_{\mathbf{v}}\leq2(1-\epsilon)^m,$$

which, of course, proves both the stated convergence result and the uniqueness of π as the only stationary probability vector for **P**.

The condition in Doeblin's Theorem is called *Doeblin's condition*, and it is instructive to understand what his theorem says in the language of *spectral theory*. Namely, as an operator on the space of bounded functions (a.k.a. column vectors with finite uniform norm), **P** has the function **1** as a right eigenfunction with eigenvalue 1: **P1** = **1**. Thus, at least if S is finite, general principles say that there must exist a row vector which is a left eigenvector of **P** with eigenvalue 1. Moreover, because 1 and the entries of **P** are real, this left eigenvector can be taken to have real components. Thus, from the spectral point of view, it is no surprise that there is a non-zero row vector $\mu \in \mathbb{R}^{\mathbb{S}}$ with the property that $\mu \mathbf{P} = \mu$. On the other hand, standard spectral theory would not predict that μ can be chosen to have non-negative components, and this is the first place where Doeblin's theorem gives information which is not readily available from standard spectral theory, even when S is finite. To interpret the estimate in Doeblin's Theorem, let $M_1(\mathbb{S}; \mathbb{C})$ denote the space of row vectors $\mathbf{v} \in \mathbb{C}^{\mathbb{S}}$ with $\|\mathbf{v}\|_{v} = 1$. Then,

$$\|\mathbf{v}\mathbf{P}\|_{\mathbf{v}} \leq 1$$
 for all $\mathbf{v} \in M_1(\mathbb{S}; \mathbb{C})$,

and so

$$\sup\{|\alpha|: \alpha \in \mathbb{C} \& \exists v \in M_1(\mathbb{S}; \mathbb{C}) v\mathbf{P} = \alpha v\} \le 1.$$

Moreover, if $\mathbf{vP} = \alpha \mathbf{v}$ for some $\alpha \neq 1$, then $\mathbf{v1} = \mathbf{v}(\mathbf{P1}) = (\mathbf{vP})\mathbf{1} = \alpha \mathbf{v1}$, and therefore $\mathbf{v1} = 0$. Thus, the estimate in (2.2.2) says that all eigenvalues of **P** which are different from 1 have absolute value dominated by $1 - \epsilon$. That is, the entire spectrum of **P** lies in the complex unit disk, 1 is a simple eigenvalue, and all the other eigenvalues lie in the disk of radius $1 - \epsilon$. Finally, although general spectral theory fails to predict Doeblin's Theorem, it should be said that there is a spectral theory, the one initiated by Frobenius and developed further by Kakutani, that does cover Doeblin's results. The interested reader should consult Chap. VIII in [2].

2.2.2 A Couple of Extensions

An essentially trivial extension of Theorem 2.2.1 is provided by the observation that, for any $M \ge 1$ and $\epsilon > 0$,⁴

$$\sup_{j} \inf_{i} (\mathbf{P}^{M})_{ij} \ge \epsilon$$

$$\implies \| \boldsymbol{\mu} \mathbf{P}^{n} - \boldsymbol{\pi} \|_{v} \le (1 - \epsilon)^{\lfloor \frac{n}{M} \rfloor} \| \boldsymbol{\mu} - \boldsymbol{\pi} \|_{v} \le 2(1 - \epsilon)^{\lfloor \frac{n}{M} \rfloor}$$
(2.2.3)

for all probability vectors $\boldsymbol{\mu}$ and a unique stationary probability vector $\boldsymbol{\pi}$. To see this, let $\boldsymbol{\pi}$ be the stationary probability vector for \mathbf{P}^{M} , the one guaranteed by Theorem 2.2.1, and note that, for any probability vector $\boldsymbol{\mu}$, any $m \in \mathbb{N}$, and any $0 \leq r < M$,

$$\left\|\boldsymbol{\mu}\mathbf{P}^{mM+r}-\boldsymbol{\pi}\right\|_{v}=\left\|\left(\boldsymbol{\mu}\mathbf{P}^{r}-\boldsymbol{\pi}\mathbf{P}^{r}\right)\mathbf{P}^{mM}\right\|_{v}\leq(1-\epsilon)^{m}\|\boldsymbol{\mu}-\boldsymbol{\pi}\|_{v}\leq2(1-\epsilon)^{m}.$$

Thus (2.2.3) has been proved, and from (2.2.3) the argument needed to show that π is the one and only stationary measure for **P** is the same as the one given in the proof of Theorem 2.2.1.

The next extension is a little less trivial. In order to appreciate the point that it is addressing, one should keep in mind the following example. Consider the transition probability matrix

$$\mathbf{P} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \text{on } \{1, 2\}.$$

Obviously, this two state chain goes in a single step from one state to the other. Thus, it certainly visits all its states. On the other hand, it does not satisfy the hypothesis in (2.2.3): $(\mathbf{P}^n)_{ij} = 0$ if either i = j and n is odd or if $i \neq j$ and n is even. Thus, it should not be surprising that the conclusion in (2.2.3) fails to hold for this **P**. Indeed, it is easy to check that although $(\frac{1}{2}, \frac{1}{2})$ is the one and only stationary probability vector for **P**, $||(1, 0)\mathbf{P}^n - (\frac{1}{2}, \frac{1}{2})||_v = 1$ for all $n \ge 0$. As we will see later (cf. Sect. 3.1.3), the problems encountered here stem from the fact that $(\mathbf{P}^n)_{ii} > 0$ only if n is even.

In spite of the problems raised by the preceding example, one should expect that the chain corresponding to this \mathbf{P} does equilibrate in some sense. To describe what I have in mind, set

$$\mathbf{A}_{n} = \frac{1}{n} \sum_{m=0}^{n-1} \mathbf{P}^{m}.$$
 (2.2.4)

Although the matrix A_n is again a transition probability matrix, it is not describing transitions but instead it is giving the average amount of time that the chain will visit

⁴Here and elsewhere, we use $\lfloor s \rfloor$ to denote the *integer part of s* of $s \in \mathbb{R}$. That is, $\lfloor s \rfloor$ is the largest integer dominated by *s*.

states. To be precise, because

$$(\mathbf{A}_n)_{ij} = \frac{1}{n} \sum_{m=0}^{n-1} \mathbb{P}(X_m = j \mid X_0 = i) = \mathbb{E}\left[\frac{1}{n} \sum_{m=0}^{n-1} \mathbf{1}_{\{j\}}(X_m) \mid X_0 = i\right],$$

 $(\mathbf{A}_n)_{ij}$ is the expected value of the average time spent at state *j* during the time interval [0, n-1] given that *i* was the state from which the chain started. Experience teaches us that data becomes much more forgiving when it is averaged, and the present situation is no exception. Indeed, continuing with the example given above, observe that, for any probability vector $\boldsymbol{\mu}$,

$$\left\| \boldsymbol{\mu} \mathbf{A}_n - \left(\frac{1}{2}, \frac{1}{2}\right) \right\|_{\mathbf{v}} \leq \frac{1}{n} \quad \text{for } n \geq 1.$$

What follows is a statement which shows that this sort of conclusion is quite general.

Theorem 2.2.5 Suppose that **P** is a transition probability matrix on S. If for some $M \in \mathbb{Z}^+$, $j_0 \in \mathbb{S}$, and $\epsilon > 0$, $(\mathbf{A}_M)_{ij_0} \ge \epsilon$ for all $i \in \mathbb{S}$, then there is precisely one stationary probability vector $\boldsymbol{\pi}$ for **P**, $(\boldsymbol{\pi})_{i_0} \ge \epsilon$, and

$$\|\boldsymbol{\mu}\mathbf{A}_n - \boldsymbol{\pi}\|_{\mathsf{v}} \leq \frac{M-1}{n\epsilon}$$

for any probability vector $\boldsymbol{\mu}$.

To get started, let π be the unique stationary probability that Theorem 2.2.3 guarantees for \mathbf{A}_M . Then, because any μ which is stationary for \mathbf{P} is certainly stationary for \mathbf{A}_M , it is clear that π is the only candidate for \mathbf{P} -stationarity. Moreover, to see that π is \mathbf{P} -stationary, observe that, because \mathbf{P} commutes with \mathbf{A}_M , $(\mathbf{\pi}\mathbf{P})\mathbf{A}_M = (\mathbf{\pi}\mathbf{A}_M)\mathbf{P} = \mathbf{\pi}\mathbf{P}$. Hence, $\mathbf{\pi}\mathbf{P}$ is stationary for \mathbf{A}_M and therefore, by uniqueness, must be equal to π . That is, $\pi = \pi\mathbf{P}$.

In order to prove the asserted convergence result, we will need an elementary property of averaging procedures. Namely, for any probability vector μ ,

$$\|\boldsymbol{\mu}\mathbf{A}_{n}\mathbf{A}_{m}-\boldsymbol{\mu}\mathbf{A}_{n}\|_{\mathrm{v}} \leq \frac{m-1}{n} \quad \text{for all } m, n \geq 1.$$
 (2.2.6)

To check this, first note that, by the triangle inequality,

$$\|\boldsymbol{\mu}\mathbf{A}_{n}\mathbf{A}_{m}-\boldsymbol{\mu}\mathbf{A}_{n}\|_{\mathrm{v}}=\frac{1}{m}\left\|\sum_{k=0}^{m-1}(\boldsymbol{\mu}\mathbf{A}_{n}\mathbf{P}^{k}-\boldsymbol{\mu}\mathbf{A}_{n})\right\|_{\mathrm{v}}$$
$$\leq\frac{1}{m}\sum_{k=0}^{m-1}\|\boldsymbol{\mu}\mathbf{A}_{n}\mathbf{P}^{k}-\boldsymbol{\mu}\mathbf{A}_{n}\|_{\mathrm{v}}.$$

Second, if $k \ge n$ then $\|\boldsymbol{\mu} \mathbf{A}_n \mathbf{P}^k - \boldsymbol{\mu} \mathbf{A}_n\|_{\mathbf{v}} \le 2 \le \frac{2k}{n}$, and if $0 \le k < n$ then

$$\boldsymbol{\mu}\mathbf{A}_{n}\mathbf{P}^{k}-\boldsymbol{\mu}\mathbf{A}_{n}=\frac{1}{n}\sum_{\ell=0}^{n-1}(\boldsymbol{\mu}\mathbf{P}^{\ell+k}-\boldsymbol{\mu}\mathbf{P}^{\ell})=\frac{1}{n}\left(\sum_{\ell=k}^{n+k-1}\boldsymbol{\mu}\mathbf{P}^{\ell}-\sum_{\ell=0}^{n-1}\boldsymbol{\mu}\mathbf{P}^{\ell}\right),$$

and so $\|\boldsymbol{\mu}\mathbf{P}^k\mathbf{A}_n - \boldsymbol{\mu}\mathbf{A}_n\|_{v} \le \frac{2k}{n}$ for all $n \ge 1$. Hence, after combining this with the first observation, we are lead to

$$\|\boldsymbol{\mu}\mathbf{A}_{n}\mathbf{A}_{m}-\boldsymbol{\mu}\mathbf{A}_{n}\|_{\mathrm{v}}\leq\frac{2}{mn}\sum_{k=0}^{m-1}k=\frac{m-1}{n},$$

which is what we wanted.

To complete the proof of Theorem 2.2.5 from here, assume that $(\mathbf{A}_M)_{ij_0} \ge \epsilon$ for all *i*, and, as above, let π be the unique stationary probability vector for \mathbf{P} . Then, π is also the unique stationary probability vector for \mathbf{A}_M , and so, by the estimate in the second line of (2.2.2) applied to \mathbf{A}_M , $\|\boldsymbol{\mu}\mathbf{A}_n\mathbf{A}_M - \pi\|_{v} = \|(\boldsymbol{\mu}\mathbf{A}_n - \pi)\mathbf{A}_M\|_{v} \le (1-\epsilon)\|\boldsymbol{\mu}\mathbf{A}_n - \pi\|_{v}$, which, in conjunction with (2.2.6), leads to

$$\|\boldsymbol{\mu}\mathbf{A}_n - \boldsymbol{\pi}\|_{\mathsf{v}} \leq \|\boldsymbol{\mu}\mathbf{A}_n - \boldsymbol{\mu}\mathbf{A}_n\mathbf{A}_M\|_{\mathsf{v}} + \|\boldsymbol{\mu}\mathbf{A}_n\mathbf{A}_M - \boldsymbol{\pi}\|_{\mathsf{v}}$$
$$\leq \frac{M-1}{n} + (1-\epsilon)\|\boldsymbol{\mu}\mathbf{A}_n - \boldsymbol{\pi}\|_{\mathsf{v}}.$$

Finally, after elementary rearrangement, this gives the required result.

2.3 Elements of Ergodic Theory

In the preceding section we saw that, under suitable conditions, either $\mu \mathbf{P}^n$ or $\mu \mathbf{A}_n$ converge and that the limit is the unique stationary probability vector $\boldsymbol{\pi}$ for **P**. In the present section, we will provide a more probabilistically oriented interpretation of these results. In particular, we will give a probabilistic interpretation of $\boldsymbol{\pi}$. This will be done again, by entirely different methods, in Chap. 4.

Before going further, it will be useful to have summarized our earlier results in the form (cf. (2.2.3) and remember that $|\mu \mathbf{f}| \le ||\mu||_v ||f||_u)^5$

$$\sup_{j} \inf_{i} \left(\mathbf{P}^{M} \right)_{ij} \ge \epsilon \implies \left\| \mathbf{P}^{n} \mathbf{f} - \boldsymbol{\pi} \mathbf{f} \right\|_{u} \le 2(1-\epsilon)^{\lfloor \frac{n}{M} \rfloor} \| \mathbf{f} \|_{u} \qquad (2.3.1)$$

and (cf. Theorem 2.2.5)

$$\sup_{j} \inf_{i} (\mathbf{A}_{M})_{ij} \ge \epsilon \implies \|\mathbf{A}_{n}\mathbf{f} - \boldsymbol{\pi}\mathbf{f}\|_{u} \le \frac{M-1}{n\epsilon} \|\mathbf{f}\|_{u}$$
(2.3.2)

when **f** is a bounded column vector.

⁵Here, and elsewhere, I abuse notation by using a constant to stand for the associated constant function.

2.3.1 The Mean Ergodic Theorem

Let {**X**_{*n*} : $n \ge 0$ } be a Markov chain with transition probability **P**. Obviously,

$$\overline{T}_{j}^{(n)} \equiv \frac{1}{n} \sum_{m=0}^{n-1} \mathbf{1}_{\{j\}}(X_{m})$$
(2.3.3)

is the average amount of time that the chain spends at *j* before time *n*. Thus, if $\boldsymbol{\mu}$ is the initial distribution of the chain (i.e., $(\boldsymbol{\mu})_i = \mathbb{P}(X_0 = i)$), then $(\boldsymbol{\mu} \mathbf{A}_n)_j = \mathbb{E}[\bar{T}_j^{(n)}]$, and so, when it applies, Theorem 2.2.5 implies that $\mathbb{E}[\bar{T}_j^{(n)}] \longrightarrow (\boldsymbol{\pi})_j$ as $n \to \infty$. Here we will be proving that the random variables $\bar{T}_j^{(n)}$ themselves, not just their expected values, tend to $(\boldsymbol{\pi})_j$ as $n \to \infty$. Such results come under the heading of *ergodic theory*. Ergodic theory is the mathematics of the principle, first enunciated by the physicist J.W. Gibbs in connection with the kinetic theory of gases, which asserts that the time-average over a particular trajectory of a dynamical system will approximate the equilibrium state of that system. Unfortunately, in spite of results, like those given here, confirming this principle, even now, nearly 150 years after Gibbs, there are embarrassingly few physically realistic situations in which Gibbs's principle has been mathematically confirmed.

Theorem 2.3.4 (Mean Ergodic Theorem) Under the hypotheses in Theorem 2.2.5,

$$\sup_{j\in\mathbb{S}} \mathbb{E}\left[\left(\overline{T}_{j}^{(n)}-(\boldsymbol{\pi})_{j}\right)^{2}\right] \leq \frac{2(M-1)}{n\epsilon} \quad for \ all \ n\geq 1.$$

(See (2.3.10) below for a more refined, less quantitative version.) More generally, for any bounded function f on S and all $n \ge 1$:

$$\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}f(X_m)-\boldsymbol{\pi}\mathbf{f}\right)^2\right] \leq \frac{4(M-1)\|\mathbf{f}\|_{\mathbf{u}}^2}{n\epsilon},$$

where \mathbf{f} denotes the column vector determined by f.

Proof Let $\mathbf{\bar{f}}$ be the column vector determined by the function $f = f - \pi \mathbf{f}$. Obviously,

$$\frac{1}{n}\sum_{m=0}^{n-1}f(X_m) - \pi \mathbf{f} = \frac{1}{n}\sum_{m=0}^{n-1}\bar{f}(X_m),$$

and so

$$\left(\frac{1}{n}\sum_{m=0}^{n-1}f(X_m) - \pi \mathbf{f}\right)^2 = \frac{1}{n^2} \left(\sum_{m=0}^{n-1}\bar{f}(X_m)\right)^2 = \frac{1}{n^2}\sum_{k,\ell=0}^{n-1}\bar{f}(X_k)\bar{f}(X_\ell) = \frac{2}{n^2}\sum_{0 \le k \le \ell < n}\bar{f}(X_k)\bar{f}(X_\ell) - \frac{1}{n^2}\sum_{k=0}^{n-1}\bar{f}(X_k)^2 \le \frac{2}{n^2}\sum_{0 \le k \le \ell < n}\bar{f}(X_k)\bar{f}(X_\ell).$$

Hence,

$$\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}f(X_m) - \pi \mathbf{f}\right)^2\right] \le \frac{2}{n^2}\sum_{k=0}^{n-1}\mathbb{E}\left[\bar{f}(X_k)\sum_{\ell=0}^{n-k-1}\bar{f}(X_{k+\ell})\right]$$
$$= \frac{2}{n^2}\sum_{k=0}^{n-1}\mathbb{E}\left[\bar{f}(X_k)\sum_{\ell=0}^{n-k-1}\left(\mathbf{P}^\ell \bar{\mathbf{f}}\right)_{X_k}\right]$$
$$= \frac{2}{n^2}\sum_{k=0}^{n-1}(n-k)\mathbb{E}\left[\bar{f}(X_k)(\mathbf{A}_{n-k}\bar{\mathbf{f}})_{X_k}\right].$$

But, by (2.3.2), $\|\mathbf{A}_{n-k}\bar{\mathbf{f}}\|_{u} \leq \frac{M-1}{(n-k)\epsilon}\|\bar{\mathbf{f}}\|_{u}$, and so, since $\|\bar{\mathbf{f}}\|_{u} \leq 2\|\mathbf{f}\|_{u}$,

$$(n-k)\mathbb{E}\left[\overline{f}(X_k)(\mathbf{A}_{n-k}\overline{\mathbf{f}})_{X_k}\right] \leq \frac{2(M-1)\|\mathbf{f}\|_{\mathbf{u}}^2}{\epsilon}$$

After plugging this into the preceding, we get the second result. To get the first, simply take $f = \mathbf{1}_{\{j\}}$ and observe that, in this case, $\|\bar{\mathbf{f}}\|_{u} \leq 1$.

2.3.2 Return Times

As the contents of Sects. 1.1 and 1.2 already indicate, return times ought to play an important role in the analysis of the long time behavior of Markov chains. In particular, if $\rho_j^{(0)} \equiv 0$ and, for $m \ge 1$, the *time of mth return* to *j* is defined so that $\rho_j^{(m)} = \infty$ if $\rho_j^{(m-1)} = \infty$ or $X_n \ne j$ for every $n > \rho^{(m-1)}$ and $\rho_j^{(m)} = \inf\{n > \rho_j^{(m-1)} : X_n = j\}$ otherwise, then we say that *j* is *recurrent* if $\mathbb{P}(\rho_j^{(1)} < \infty | X_0 = j) = 1$ and that it is *transient* if $\mathbb{P}(\rho_j^{(1)} < \infty | X_0 = j) < 1$; and we can hope that when *j* is recurrent, then the history of the chain breaks into epochs which are punctuated by the successive returns to *j*. In this subsection we will provide evidence which bolsters that hope.

Notice that $\rho_j \equiv \rho_j^{(1)} \ge 1$ and, for $n \ge 1$,

$$\mathbf{1}_{(n,\infty]}(\rho_j) = F_{n,j}(X_0, \dots, X_n) \quad \text{where}$$

$$F_{n,j}(i_0, \dots, i_n) = \begin{cases} 1 & \text{if } i_m \neq j \text{ for } 1 \leq m \leq n, \\ 0 & \text{otherwise.} \end{cases}$$
(2.3.5)

In particular, this shows that the event $\{\rho_j > n\}$ is a measurable function of (X_0, \ldots, X_n) . More generally, because

$$\mathbf{1}_{(n,\infty]}(\rho_j^{(m+1)}) = \mathbf{1}_{[n,\infty]}(\rho_j^{(m)}) + \sum_{\ell=1}^{n-1} \mathbf{1}_{\{\ell\}}(\rho_j^{(m)}) F_{n-\ell,j}(X_\ell,\ldots,X_n),$$

an easy inductive argument shows that, for each $m \in \mathbb{N}$ and $n \in \mathbb{N}$, $\{\rho_j^{(m)} > n\}$ is a measurable function of (X_0, \ldots, X_n) .

Theorem 2.3.6 For all $m \in \mathbb{Z}^+$ and $(i, j) \in \mathbb{S}^2$,

$$\mathbb{P}\left(\rho_j^{(m)} < \infty \mid X_0 = i\right) = \mathbb{P}(\rho_j < \infty \mid X_0 = i)\mathbb{P}(\rho_j < \infty \mid X_0 = j)^{m-1}.$$

In particular, if j is recurrent, then $\mathbb{P}(\rho_j^{(m)} < \infty | X_0 = j) = 1$ for all $m \in \mathbb{N}$. In fact, if j is recurrent, then, conditional on $X_0 = j$, $\{\rho_j^{(m)} - \rho_j^{(m-1)} : m \ge 1\}$ is a sequence of mutually independent random variables each of which has the same distribution as ρ_j .

Proof To prove the first statement, we apply (2.1.11) and the monotone convergence theorem, Theorem 7.1.9, to justify

$$\mathbb{P}(\rho_{j}^{(m)} < \infty \mid X_{0} = i)$$

$$= \sum_{n=1}^{\infty} \mathbb{P}(\rho_{j}^{(m-1)} = n \& \rho_{j}^{(m)} < \infty \mid X_{0} = i)$$

$$= \lim_{N \to \infty} \sum_{n=1}^{\infty} \mathbb{P}(\rho_{j}^{(m-1)} = n \& \rho_{j}^{(m)} \le N \mid X_{0} = i)$$

$$= \sum_{n=1}^{\infty} \lim_{N \to \infty} \mathbb{E}[1 - F_{N,j}(X_{n}, \dots, X_{n+N}), \rho_{j}^{(m-1)} = n \mid X_{0} = i]$$

$$= \sum_{n=1}^{\infty} \lim_{N \to \infty} \mathbb{E}[1 - F_{N,j}(X_{0}, \dots, X_{N}) \mid X_{0} = j] \mathbb{P}(\rho_{j}^{(m-1)} = n \mid X_{0} = i)$$

$$= \sum_{n=1}^{\infty} \lim_{N \to \infty} \mathbb{P}(\rho_{j} \le N \mid X_{0} = j) \mathbb{P}(\rho_{j}^{(m-1)} = n \mid X_{0} = i)$$

$$= \mathbb{P}(\rho_{j} < \infty \mid X_{0} = j) \mathbb{P}(\rho_{j}^{(m-1)} < \infty \mid X_{0} = i).$$

Turning to the second statement, note that it suffices for us prove that

$$\mathbb{P}(\rho_j^{(m+1)} > n + n_m \mid X_0 = j, \rho_j^{(1)} = n_1, \dots, \rho_j^{(m)} = n_m)$$

= $\mathbb{P}(\rho_j > n \mid X_0 = j).$

But, again by (2.1.11), the expression on the left is equal to

$$\mathbb{E} \Big[F_{n,j}(X_{n_m}, \dots, X_{n_m+n}) \, \big| \, X_0 = j, \, \rho_j^{(1)} = n_1, \dots, \, \rho_j^{(m)} = n_m \Big] \\ = \mathbb{E} \Big[F_{n,j}(X_0, \dots, X_n) \, \big| \, X_0 = j \Big] = \mathbb{P}(\rho_j > n \, | \, X_0 = j). \qquad \Box$$

Reasoning as we did in Sect. 1.2.2, we can derive from the first part of Theorem 2.3.6:

$$\mathbb{E}[T_j \mid X_0 = i] = \delta_{i,j} + \frac{\mathbb{P}(\rho_j < \infty \mid X_0 = i)}{\mathbb{P}(\rho_j = \infty \mid X_0 = j)},$$

$$\mathbb{E}[T_j \mid X_0 = j] = \infty \quad \Longleftrightarrow \quad \mathbb{P}(T_j = \infty \mid X_0 = j) = 1,$$

$$\mathbb{E}[T_j \mid X_0 = j] < \infty \quad \Longleftrightarrow \quad \mathbb{P}(T_j < \infty \mid X_0 = j) = 1,$$
(2.3.7)

where $T_j = \sum_{m=0}^{\infty} \mathbf{1}_{\{j\}}(X_m)$ is the total time the chain spends in the state *j*. Indeed, because

$$\mathbb{P}(T_j > m \mid X_0 = i) = \begin{cases} \mathbb{P}(\rho_j^{(m)} < \infty \mid X_0 = j) & \text{if } i = j \\ \mathbb{P}(\rho_j^{(m+1)} < \infty \mid X_0 = i) & \text{if } i \neq j, \end{cases}$$

all three parts of (2.3.7) follow immediately from the first part of Theorem 2.3.6.

Of course, from (2.3.7) we know that

j is recurrent if and only if $\mathbb{E}[T_j|X_0 = j] = \infty$.

In particular, under the conditions in Theorem 2.2.5, this means that j_0 is recurrent since $(\mathbf{A}_n)_{j_0j_0} \longrightarrow (\boldsymbol{\pi})_{j_0} > 0$ and therefore

$$\mathbb{E}[T_{j_0} | X_0 = j_0] = \sum_{m=0}^{\infty} (\mathbf{P}^m)_{j_0 j_0} = \lim_{n \to \infty} n(\mathbf{A}_n)_{j_0 j_0} = \infty.$$

To facilitate the statement of the next result, we will say that *j* is *accessible* from *i* and will write $i \rightarrow j$ if $(\mathbf{P}^n)_{ij} > 0$ for some $n \ge 0$. Equivalently, $i \rightarrow j$ if and only if i = j or $i \ne j$ and $\mathbb{P}(\rho_j < \infty | X_0 = i) > 0$.

Theorem 2.3.8 Assume that $\inf_k(\mathbf{A}_M)_{ki} \ge \epsilon$ for some $M \ge 1$, $i \in \mathbb{S}$, and $\epsilon > 0$. Then *j* is recurrent if and only if $i \rightarrow j$, in which case $\inf_k(\mathbf{A}_{M'})_{kj} > 0$ for some M' > 0. Moreover, if *k* is recurrent, then $\sup_{j \in \mathbb{S}} \mathbb{E}[\rho_k^p | X_0 = j] < \infty$ for all $p \in (0, \infty)$. *Proof* First suppose that $i \nleftrightarrow j$. Equivalently, $\mathbb{P}(\rho_j = \infty | X_0 = i) = 1$. At the same time, because $(\mathbf{A}_M)_{ji} \ge \epsilon$, there exists an $1 \le m < M$ such that $(\mathbf{P}^m)_{ji} > 0$, and so (cf. (2.3.5))

$$\mathbb{P}(\rho_j = \infty \mid X_0 = j) \ge \mathbb{P}(X_n \neq j \text{ for } n \ge m \& X_m = i \mid X_0 = j)$$

$$\ge \lim_{N \to \infty} \mathbb{E} \Big[F_{N,j}(X_m, \dots, X_{m+N}), \ X_m = i \mid X_0 = j \Big]$$

$$= \lim_{N \to \infty} \mathbb{E} \Big[F_{N,j}(X_0, \dots, X_N) \mid X_0 = i \Big] \mathbb{P}(X_m = i \mid X_0 = j)$$

$$= \mathbb{P}(\rho_j = \infty \mid X_0 = i) \Big(\mathbf{P}^m \Big)_{ji} > 0.$$

Therefore j cannot be recurrent.

We next show that

$$i \rightarrow j \implies \inf_k (\mathbf{A}_{M'})_{kj} > 0 \text{ for some } M' \ge 1.$$
 (*)

To this end, choose $m \in \mathbb{N}$ so that $(\mathbf{P}^m)_{ij} > 0$. Then, for all $k \in \mathbb{S}$,

$$(\mathbf{A}_{m+M})_{kj} = \frac{1}{m+M} \sum_{\ell=0}^{M+m-1} (\mathbf{P}^{\ell})_{kj} \ge \frac{1}{m+M} \sum_{\ell=0}^{M-1} (\mathbf{P}^{\ell})_{ki} (\mathbf{P}^{m})_{ij}$$
$$= \frac{M}{m+M} (\mathbf{A}_{M})_{ki} (\mathbf{P}^{m})_{ij} \ge \frac{M\epsilon}{m+M} (\mathbf{P}^{m})_{ij} > 0.$$

In view of (*) and what we have already shown, it suffices to show that $\mathbb{E}[\rho_i^p | X_0 = j] < \infty$ for all $j \in \mathbb{S}$. For this purpose, set

$$u(n,k) = \mathbb{P}(\rho_i > nM \mid X_0 = k) \text{ for } n \in \mathbb{Z}^+ \text{ and } k \in \mathbb{S}.$$

Then, by (2.1.11),

$$u(n+1,k) = \sum_{j \in \mathbb{S}} \mathbb{P}(\rho_i > (n+1)M \& X_{nM} = j \mid X_0 = k)$$

= $\sum_{j \in \mathbb{S}} \mathbb{E}[F_{M,i}(X_{nM}, \dots, X_{(n+1)M}), \rho_i > nM \& X_{nM} = j \mid X_0 = k]$
= $\sum_{j \in \mathbb{S}} \mathbb{P}(\rho_i > M \mid X_0 = j)\mathbb{P}(\rho_i > nM \& X_{nM} = j \mid X_0 = k)$
= $\sum_{j \in \mathbb{S}} u(1, j)\mathbb{P}(\rho_i > nM \& X_{nM} = j \mid X_0 = k).$

Hence, $u(n + 1, k) \le Uu(n, k)$ where $U \equiv \max_{j \in \mathbb{S}} u(1, j)$. Finally, since $u(1, j) = 1 - \mathbb{P}(\rho_i \le M | X_0 = j)$ and

$$\mathbb{P}(\rho_i \le M \mid X_0 = j) \ge \max_{0 \le m < M} \left(\mathbf{P}^m\right)_{ji} \ge (\mathbf{A}_M)_{ji} \ge \epsilon,$$

 $U \le 1 - \epsilon$. In particular, this means that $u(n+1, k) \le (1-\epsilon)u(n, k)$, and therefore that $\mathbb{P}(\rho_i > nM | X_0 = k) \le (1-\epsilon)^n$, from which

$$\mathbb{E}\left[\rho_i^p \mid X_0 = k\right] = \sum_{n=1}^{\infty} n^p \mathbb{P}(\rho_i = n \mid X_0 = k)$$

$$\leq \sum_{m=1}^{\infty} (mM)^p \sum_{n=(m-1)M+1}^{mM} \mathbb{P}(\rho_i = n \mid X_0 = k)$$

$$\leq M^p \sum_{m=1}^{\infty} m^p \mathbb{P}\left(\rho_i > (m-1)M \mid X_0 = k\right)$$

$$\leq M^p \sum_{m=1}^{\infty} m^p (1-\epsilon)^{m-1} < \infty$$

follows immediately.

2.3.3 Identification of π

Under the conditions in Theorem 2.2.5, we know that there is precisely one **P**-stationary probability vector $\boldsymbol{\pi}$. In this section, we will give a probabilistic interpretation of $(\boldsymbol{\pi})_i$. Namely, we will show that

$$\sup_{M \ge 1} \sup_{j \in \mathbb{S}} \inf(\mathbf{A}_M)_{ij} > 0$$

$$\implies \quad (\boldsymbol{\pi})_j = \frac{1}{\mathbb{E}[\rho_j | X_0 = j]} \quad (\equiv 0 \text{ if } j \text{ is transient}). \tag{2.3.9}$$

The idea for the proof of (2.3.9) is that, on the one hand, (cf. (2.3.3))

$$\mathbb{E}\left[\overline{T}_{j}^{(n)} \mid X_{0} = j\right] = (\mathbf{A}_{n})_{jj} \longrightarrow (\boldsymbol{\pi})_{j},$$

while, on the other hand,

$$X_0 = j \implies \overline{T}_j^{(\rho_j^{(m)})} = \frac{1}{\rho_j^{(m)}} \sum_{\ell=0}^{\rho_j^{(m)}-1} \mathbf{1}_{\{j\}}(X_\ell) = \frac{m}{\rho_j^{(m)}}$$

Thus, since, at least when j is recurrent, Theorem 2.3.6 says that $\rho_j^{(m)}$ is the sum of m mutually independent copies of ρ_j , the preceding combined with the weak law of large numbers should lead

$$(\boldsymbol{\pi})_j = \lim_{m \to \infty} \mathbb{E}\left[\overline{T}_j^{(\rho_j^{(m)})} \, \big| \, X_0 = j\right] = \frac{1}{\mathbb{E}[\rho_j | X_0 = j]}.$$

To carry out the program suggested above, we will actually prove a stronger result. Namely, we will show that, for each $j \in \mathbb{S}^{6}$,

$$\mathbb{P}\left(\lim_{n \to \infty} \overline{T}_{j}^{(n)} = \frac{1}{\mathbb{E}[\rho_{j} \mid X_{0} = j]} \mid X_{0} = j\right) = 1.$$
(2.3.10)

In particular, because $0 \le \overline{T}^{(n)} \le 1$, Lebesgue's dominated convergence theorem, Theorem 7.1.11, says that

$$(\boldsymbol{\pi})_j = \lim_{n \to \infty} (\mathbf{A}_n)_{jj} = \lim_{n \to \infty} \mathbb{E}\left[\overline{T}_j^{(n)} \mid X_0 = j\right] = \frac{1}{\mathbb{E}[\rho_j \mid X_0 = j]}$$

follows from (2.3.10). Thus, we need only prove (2.3.10). To this end, choose j_0 , M, and $\epsilon > 0$ so that $(\mathbf{A}_M)_{ij_0} \ge \epsilon$ for all i. If $j_0 \nleftrightarrow j$, then, by Theorem 2.3.8, j is transient, and so, by (2.3.7), $\mathbb{P}(T_j < \infty | X_0 = j) = 1$. Hence, conditional on $X_0 = j$, $\overline{T}_j^{(n)} \le \frac{1}{n}T_j \longrightarrow 0$ with probability 1. At the same time, because j is transient, $\mathbb{P}(\rho_j = \infty | X_0 = j) > 0$, and so $\mathbb{E}[\rho_j | X_0 = j] = \infty$. Hence, we have proved (2.3.10) in the case when $j_0 \nleftrightarrow j$.

Next assume that $j_0 \rightarrow j$. Then, again by Theorem 2.3.8, $\mathbb{E}[\rho_j^4|X_0 = j] < \infty$ and, conditional on $X_0 = j$, Theorem 2.3.6 says that the random variables $\rho_j^{(m)} - \rho_j^{(m-1)}$ are mutually independent random variables and have the same distribution as ρ_j . In particular, by the strong law of large numbers (cf. Exercise 1.3.4)

$$\mathbb{P}\left(\lim_{m \to \infty} \frac{\rho_j^{(m)}}{m} = r_j \mid X_0 = j\right) = 1 \quad \text{where } r_j \equiv \mathbb{E}[\rho_j | X_0 = j].$$

On the other hand, for any $m \ge 1$,

$$\left|\overline{T}_{j}^{(n)}-\frac{1}{r_{j}}\right| \leq \left|\overline{T}_{j}^{(n)}-\overline{T}_{j}^{(\rho_{j}^{(m)})}\right| + \left|\overline{T}_{j}^{(\rho_{j}^{(m)})}-\frac{1}{r_{j}}\right|,$$

and

$$\begin{aligned} |\overline{T}_{j}^{(n)} - \overline{T}_{j}^{(\rho_{j}^{(m)})}| &\leq \frac{|T_{j}^{(n)} - T_{j}^{(\rho_{j}^{(m)})}|}{n} + \left|1 - \frac{\rho_{j}^{(m)}}{n}\right| \overline{T}_{j}^{(\rho_{j}^{(m)})} \\ &\leq 2\left|1 - \frac{\rho_{j}^{(m)}}{n}\right| \leq 2\left|1 - \frac{mr_{j}}{n}\right| + \frac{2m}{n}\left|\frac{\rho_{j}^{(m)}}{m} - r_{j}\right| \end{aligned}$$

⁶Statements like the one which follows are called *individual ergodic theorems* because they, as distinguished from the first part of Theorem 2.3.4, are about convergence with probability 1 as opposed to convergence in square mean. See Exercise 4.2.10 below for more information.

while, since $\overline{T}_{j}^{(\rho_{j}^{(m)})} = \frac{m}{\rho_{j}^{(m)}} \leq 1$,

$$\left|\overline{T}_{j}^{(\rho_{j}^{(m)})}-\frac{1}{r_{j}}\right|\leq\frac{1}{r_{j}}\left|\frac{\rho_{j}^{(m)}}{m}-r_{j}\right|.$$

Hence,

$$\left|\overline{T}_{j}^{(n)}-\frac{1}{r_{j}}\right| \leq 2\left|1-\frac{mr_{j}}{n}\right| + \left(\frac{2m}{n}+\frac{1}{r_{j}}\right)\left|\frac{\rho_{j}^{(m)}}{m}-r_{j}\right|.$$

Finally, by taking $m_n = \lfloor \frac{n}{r_i} \rfloor$ we get

$$\left|\overline{T}_{j}^{(n)}-\frac{1}{r_{j}}\right| \leq \frac{2r_{j}}{n} + \frac{3}{r_{j}}\left|\frac{\rho_{j}^{(m_{n})}}{m_{n}}-r_{j}\right| \longrightarrow 0 \quad \text{as } n \to \infty.$$

Notice that (2.3.10) is precisely the sort of statement for which Gibbs was looking. That is, it says that, with probability 1, when one observes an individual path, the average time that it spends in each state tends, as one observes for a longer and longer time, to the probability that the equilibrium (i.e., stationary) distribution assigns to that state.

2.4 Exercises

Exercise 2.4.1 In this exercise we will give a probabilistic interpretation of the *adjoint* of a transition probability matrix with respect to a stationary probability. To be precise, suppose that the transition probability matrix **P** admits a stationary distribution $\boldsymbol{\pi}$, assume $(\boldsymbol{\pi})_i > 0$ for each $i \in \mathbb{S}$, and determine the matrix \mathbf{P}^\top by $(\mathbf{P}^\top)_{ij} = \frac{(\boldsymbol{\pi})_j}{(\boldsymbol{\pi})} (\mathbf{P})_{ji}$.

- (a) Show that \mathbf{P}^{\top} is a transition probability matrix for which π is again a stationary probability.
- (b) Use P and P^T to denote probabilities computed for the chains determined, respectively, by P and P^T with initial distribution π, and show that these chains are the *reverse* of one another in the sense that, for each n ≥ 0 the distribution of (X₀,..., X_n) under P^T is the same as the distribution of (X_n,..., X₀) under P. That is,

$$\mathbb{P}^+(X_0 = i_0, \dots, X_n = i_n) = \mathbb{P}(X_n = i_0, \dots, X_0 = i_n)$$

for all $n \ge 0$ and $(i_0, \ldots, i_n) \in \mathbb{S}^{n+1}$.

Exercise 2.4.2 The Doeblin theory applies particularly well to chains on a finite state. For example, suppose that **P** is a transition probability matrix on an *N* element state space \mathbb{S} , and show that there exists an $\epsilon > 0$ such that $(\mathbf{A}_N)_{ij_0} \ge \epsilon$ for all $i \in \mathbb{S}$

if and only if $i \rightarrow j_0$ for all $i \in S$. In particular, if such a j_0 exists, conclude that, for all probability vectors μ ,

$$\|\boldsymbol{\mu}\mathbf{A}_n - \boldsymbol{\pi}\|_{\mathsf{v}} \leq \frac{2(N-1)}{n\epsilon}, \quad n \geq 1,$$

where π is the unique stationary probability vector for **P**.

Exercise 2.4.3 Here is a version of Doeblin's theorem that sometimes gives a slightly better estimate. Namely, assume that $(\mathbf{P})_{ij} \ge \epsilon_j$ for all (i, j), and set $\epsilon = \sum_j \epsilon_j$. If $\epsilon > 0$, show that the conclusion of Theorem 2.2.1 holds with this ϵ and that $(\boldsymbol{\pi})_i \ge \epsilon_i$ for each $i \in \mathbb{S}$.

Exercise 2.4.4 Assume that **P** is a transition probability matrix on the finite state space S, and show that

$$j \in \mathbb{S}$$
 is recurrent if and only if $\mathbb{E}[\rho_j | X_0 = j] < \infty$.

Of course, the "if" part is trivial and has nothing to do with the finiteness of the state space.

Exercise 2.4.5 Again assume that **P** is a transition probability matrix on the finite state space S. In addition, assume that **P** is *doubly stochastic* in the sense that each of its columns as well as each of its rows sums to 1. Under the condition that every state is accessible from every other state, show that $\mathbb{E}[\rho_j|X_0 = j] = \#S$ for each $j \in S$.

Exercise 2.4.6 In order to test how good Doeblin's theorem is, consider the case when $S = \{1, 2\}$ and

$$\mathbf{P} = \begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix} \text{ for some } (\alpha, \beta) \in (0, 1).$$

Show that $\pi = (\alpha + \beta)^{-1}(\beta, \alpha)$ is a stationary probability for **P**, that

 $\|\mathbf{v}\mathbf{P} - \boldsymbol{\pi}\|_{v} = |1 - \alpha - \beta| \|\boldsymbol{\mu} - \boldsymbol{\pi}\|_{v}$ for all probability vectors $\boldsymbol{\mu}$.

Hence, in this case, Doeblin's theorem gives the optimal result.

Exercise 2.4.7 One of the earliest examples of Markov processes are the *branching processes* introduced, around the end of the nineteenth century, by Galton and Watson to model demographics. In this model, $\mathbb{S} = \mathbb{N}$, the state $i \in \mathbb{N}$ representing the number of members in the population, and the process evolves so that, at each stage, every individual, independently of all other members of the population, dies and is replaced by a random number of offspring. Thus, 0 is an absorbing state, and, given that there are $i \geq 1$ individuals alive at a given time *n*, the number of

individuals alive at time n + 1 will be distributed like the sum of i mutually independent, \mathbb{N} -valued, identically distributed random variables. To be more precise, if $\boldsymbol{\mu} = (\mu_0, \dots, \mu_k, \dots)$ is the probability vector giving the number of offspring each individual produces, define the *m*-fold convolution power $\mu^{\star m}$ so that $(\mu^{\star 0})_i = \delta_{0,i}$ and, for $m \ge 1$,

$$(\boldsymbol{\mu}^{\star m})_{j} = \sum_{i=0}^{j} (\boldsymbol{\mu}^{\star (m-1)})_{j-i} \mu_{i}.$$

Then the transition probability matrix **P** is given by $(\mathbf{P})_{ij} = (\boldsymbol{\mu}^{\star i})_j$.

The first interesting question which one should ask about this model is what it predicts will be the probability of eventual *extinction*. That is, what is $\lim_{n\to\infty} \mathbb{P}(X_n = 0)$? A naïve guess is that eventual extinction should occur or should not occur depending on whether the expected number $\gamma \equiv \sum_{k=0}^{\infty} k \mu_k$ of progeny is strictly less or strictly greater than 1, with the case when the expected number is precisely 1 being more ambiguous. In order to verify this guess and remove trivial special cases, we make the assumptions that $\mu_0 > 0$, $\mu_0 + \mu_1 < 1$, and $\gamma \equiv \sum_{k=0}^{\infty} k \mu_k < \infty.$

(a) Set $f(s) = \sum_{k=0}^{\infty} s^k \mu_k$ for $s \in [0, 1]$, and define $f^{\circ n}(s)$ inductively so that $f^{\circ 0}(s) = s$ and $f^{\circ n} = f \circ f^{\circ (n-1)}$ for $n \ge 1$. Show that $\gamma = f'(1)$ and that

$$f^{\circ n}(s)^{i} = \mathbb{E}[s^{X_{n}} | X_{0} = i] = \sum_{j=0}^{\infty} s^{j} (\mathbf{P}^{n})_{ij} \text{ for } s \in [0, 1] \text{ and } i \ge 0.$$

- **Hint**: Begin by showing that $f(s)^i = \sum_{j=0}^{\infty} s^j (\boldsymbol{\mu}^{\star i})_j$. (b) Observe that $s \in [0, 1] \mapsto f(s) s$ is a continuous function which is positive at s = 0, zero at s = 1, and smooth and strictly convex (i.e., f'' > 0) on (0, 1). Conclude that either $\gamma \leq 1$ and f(s) > s for all $s \in [0, 1)$ or $\gamma > 1$ and there is exactly one $\alpha \in (0, 1)$ at which $f(\alpha) = \alpha$.
- (c) Referring to the preceding, show that

$$\gamma \le 1 \implies \lim_{n \to \infty} \mathbb{E}[s^{X_n} \mid X_0 = i] = 1 \text{ for all } s \in (0, 1]$$

and that

$$\gamma > 1 \implies \lim_{n \to \infty} \mathbb{E}[s^{X_n} | X_0 = i] = \alpha^i \text{ for all } s \in (0, 1).$$

(d) Based on (c), conclude that $\gamma \leq 1 \implies \mathbb{P}(X_n = 0 | X_0 = i) \longrightarrow 1$ and that $\gamma > 1 \implies \lim_{n \to \infty} \mathbb{P}(X_n = 0 | X_0 = i) = \alpha^i$ and

$$\lim_{n \to \infty} \mathbb{P}(1 \le X_n \le L | X_0 = i) = 0 \quad \text{for all } L \ge 1.$$

The last conclusion has the ominous implication that, when the expected number of progeny is larger than 1, then the population either becomes extinct or, what may be worse, grows indefinitely.

Exercise 2.4.8 Continue with the setting and notation in Exercise 2.4.7. We will to show in this exercise that there are significant differences between the cases when $\gamma < 1$ and $\gamma = 1$.

- (a) Show that E[X_n | X₀ = i] = iγⁿ. Hence, when γ < 1, the expected size of the population goes to 0 at an exponential rate. On the other hand, when γ = 1, the expected size remains constant, this in spite of the fact that as n → ∞ P(X_n = 0|X₀ = i) → 1. Thus, when γ = 1, we have a typical situation of the sort which demonstrates why Lebesgue had to make the hypotheses he did in his dominated convergence theorem, Theorem 7.1.11. In the present case, the explanation is simple: as n → ∞, with large probability X_n = 0 but, nonetheless, with positive probability X_n is enormous.
- (b) Let ρ_0 be the time of first return to 0. Show that

$$\mathbb{P}(\rho_0 \le n | X_0 = i) = \mathbb{P}(X_n = 0 | X_0 = i) = \left(f^{\circ (n-1)}(\mu_0) \right)^i,$$

and use this to get the estimate

$$\mathbb{P}(\rho_0 > n | X_0 = i) \le i \gamma^{n-1} (1 - \mu_0)$$

In particular, this shows that $\mathbb{E}[\rho_0^k | X_0 = i] < \infty$ for all $k \in \mathbb{Z}^+$ when $\gamma < 1$. (c) Now assume that $\gamma = 1$. Under the additional condition that $\beta \equiv f''(1) = \sum_{k \ge 2} k(k-1)\mu_k < \infty$, start from $\mathbb{P}(\rho_0 \le n | X_0 = 1) = f^{\circ (n-1)}(\mu_0)$, and show

that $\mathbb{E}[\rho_0|X_0=i] = \infty$ for all $i \ge 1$.

Hint: Begin by showing that

$$1 - f^{\circ n}(\mu_0) \ge \left(\prod_{\ell=m}^{n-1} \left(1 - \beta \left(1 - f^{\circ \ell}(\mu_0) \right) \right) \right) \left(1 - f^{\circ m}(\mu_0) \right)$$

for n > m. Next, use this to show that

$$\infty > \mathbb{E}[\rho_0|X_0 = 1] = 1 + \sum_{0}^{\infty} (1 - f^{\circ n}(\mu_0))$$

would lead to a contradiction.

(d) Here we show that the conclusion in (c) will, in general, be false without the finiteness condition on the second derivative. To see this, let $\theta \in (0, 1)$ be given, and check that $f(s) \equiv s + \frac{(1-s)^{1+\theta}}{1+\theta} = \sum_{k=0}^{\infty} s^k \mu_k$, for some probability vector $\boldsymbol{\mu} = (\mu_0, \dots, \mu_k, \dots)$ with $\mu_k > 0$ unless k = 1. Now use this choice of $\boldsymbol{\mu}$ to see that, when the second derivative condition in (c) fails, $\mathbb{E}[\rho_0|X_0 = 1]$ can be finite even though $\gamma = 1$.

Hint: Set $a_n = 1 - f^{\circ n}(\mu_0)$, note that $a_n - a_{n+1} = \mu_0 a_n^{1+\theta}$, and use this first to see that $\frac{a_{n+1}}{a_n} \longrightarrow 1$ and then that there exist $0 < c_2 < c_2 < \infty$ such that $c_1 \le a_{n+1}^{-\theta} - a_n^{-\theta} \le c_2$ for all $n \ge 1$. Conclude that $\mathbb{P}(\rho_0 > n | X_0 = 1)$ tends to 0 like $n^{-\frac{1}{\theta}}$.

Exercise 2.4.9 The idea underlying this exercise was introduced by J.L. Doob and is called⁷ *Doob's h-transformation*. Let **P** is a transition probability matrix on the state space S. Next, let $\emptyset \neq \Gamma \subsetneq S$ be given, set

$$\rho_{\Gamma} = \inf\{n \ge 1 : X_n \in \Gamma\},\$$

and assume that

$$h(i) \equiv \mathbb{P}(\rho_{\Gamma} = \infty \mid X_0 = i) > 0 \quad \text{for all } i \in \hat{\mathbb{S}} \equiv \mathbb{S} \setminus \Gamma.$$

- (a) Show that $h(i) = \sum_{j \in \hat{\mathbb{S}}} (\mathbf{P})_{ij} h(j)$ for all $i \in \hat{\mathbb{S}}$, and conclude that the matrix $\hat{\mathbf{P}}$ given by $(\hat{\mathbf{P}})_{ij} = \frac{1}{h(i)} (\mathbf{P})_{ij} h(j)$ for $(i, j) \in (\hat{\mathbb{S}})^2$ is a transition probability matrix on $\hat{\mathbb{S}}$.
- (b) For all $n \in \mathbb{N}$ and $(j_0, \ldots, j_n) \in (\hat{\mathbb{S}})^{n+1}$, show that, for each $i \in \hat{\mathbb{S}}$,

$$\hat{\mathbb{P}}(X_0 = j_0, \dots, X_n = j_n | X_0 = i) = \mathbb{P}(X_0 = j_0, \dots, X_n = j_n | \rho_{\Gamma} = \infty \& X_0 = i),$$

where $\hat{\mathbb{P}}$ is used here to denote probabilities computed for the Markov chain on $\hat{\mathbb{S}}$ whose transition probability matrix is $\hat{\mathbf{P}}$. That is, the Markov chain determined by $\hat{\mathbf{P}}$ is the Markov chain determined by \mathbf{P} conditioned to never hit Γ .

Exercise 2.4.10 Here is another example of an *h*-transform. Assume that $j_0 \in \mathbb{S}$ is transient but that $i \rightarrow j_0$ for all $i \in \mathbb{S}$.⁸ Set

$$h(j_0) = 1$$
 and $h(i) = \mathbb{P}(\rho_{j_0} < \infty | X_0 = i)$ for $i \neq j_0$.

(a) After checking that h(i) > 0 for all $i \in S$, define $\hat{\mathbf{P}}$ so that

$$(\hat{\mathbf{P}})_{ij} = \begin{cases} (\mathbf{P})_{j_0 j} & \text{if } i = j_0 \\ h(i)^{-1}(\mathbf{P})_{ij}h(j) & \text{if } i \neq j_0. \end{cases}$$

Show that $\hat{\mathbf{P}}$ is again a transition probability matrix.

(b) Using $\hat{\mathbb{P}}$ to denote probabilities computed relative to the chain determined by $\hat{\mathbf{P}}$, show that

$$\hat{\mathbb{P}}(\rho_{j_0} > n \mid X_0 = i) = \frac{1}{h(i)} \mathbb{P}(n < \rho_{j_0} < \infty \mid X_0 = i)$$

for all $n \in \mathbb{N}$ and $i \neq j_0$.

(c) Starting from the result in (b), show that j_0 is recurrent for the chain determined by $\hat{\mathbf{P}}$.

⁷The "h" comes from the connection with harmonic functions.

⁸By Exercise 2.4.2, this is possible only if S in infinite.

Chapter 3 Stationary Probabilities

It is important to understand when stationary probabilities exist and how to compute them when they do. In this chapter I will develop some methods for addressing these questions.

3.1 Classification of States

In this section we deal with a topic which was hinted at but not explicitly discussed in Chap. 2. Namely, a transition probability **P** determines a relationship structure on the state space. To be precise, given a pair (i, j) of states, recall that we write $i \rightarrow j$ and say that j is accessible from i if, with positive probability, the chain can go from state i to state j. That is, $(\mathbf{P}^n)_{ij} > 0$ for some $n \in \mathbb{N}$. Notice that accessibility is transitive in the sense that

$$i \to j \text{ and } j \to \ell \implies i \to \ell.$$
 (3.1.1)

Indeed, if $(\mathbf{P}^m)_{ij} > 0$ and $(\mathbf{P}^n)_{j\ell} > 0$, then

$$\left(\mathbf{P}^{m+n}\right)_{i\ell} = \sum_{k} \left(\mathbf{P}^{m}\right)_{ik} \left(\mathbf{P}^{n}\right)_{k\ell} \ge \left(\mathbf{P}^{m}\right)_{ij} \left(\mathbf{P}^{n}\right)_{j\ell} > 0.$$

If *i* and *j* are accessible from one another in the sense that $i \rightarrow j$ and $j \rightarrow i$, then we write $i \leftrightarrow j$ and say that *i* communicates with *j*. It should be clear that " \leftrightarrow " is an equivalence relation. To wit, because $(\mathbf{P}^0)_{ii} = 1$, $i \leftrightarrow i$, and it is trivial that $j \leftrightarrow i$ if $i \leftrightarrow j$. Finally, if $i \leftrightarrow j$ and $j \leftrightarrow \ell$, then (3.1.1) makes it obvious that $i \leftrightarrow \ell$. Thus, " \leftrightarrow " leads to a partitioning of the state space into equivalence classes made up of communicating states. That is, for each state *i*, the *communicating equivalence class* [i] of *i* is the set of states *j* such that $i \leftrightarrow j$; and, for every pair (i, j), either [i] = [j]or $[i] \cap [j] = \emptyset$. In the case when every state communicates with every other state, we say that the chain is *irreducible*.

D.W. Stroock, *An Introduction to Markov Processes*, Graduate Texts in Mathematics 230, 49 DOI 10.1007/978-3-642-40523-5_3, © Springer-Verlag Berlin Heidelberg 2014

3.1.1 Classification, Recurrence, and Transience

In this subsection we will show that recurrence and transience are *communicating class properties*. That is, either all members of a communicating equivalence class are recurrent or all members are transient.

Recall (cf. Sect. 2.3.2) that ρ_j is the time of first return to j, and observe that $i \rightarrow j$ if and only if $\mathbb{P}(\rho_j < \infty | X_0 = 1) > 0$. Also, remember that we say j is recurrent or transient according to whether $\mathbb{P}(\rho_j < \infty | X_0 = j)$ is equal to or strictly less than 1.

Theorem 3.1.2 Assume that *i* is recurrent and that $j \neq i$. Then $i \rightarrow j$ if and only if $\mathbb{P}(\rho_j < \rho_i | X_0 = i) > 0$. Moreover, if $i \rightarrow j$, then $\mathbb{P}(\rho_k < \infty | X_0 = \ell) = 1$ for any $(k, \ell) \in \{i, j\}^2$. In particular, $i \rightarrow j$ implies that $i \leftrightarrow j$ and that *j* is recurrent.

Proof Given $j \neq i$ and $n \ge 1$, set (cf. (2.3.5))

$$G_n(k_0,\ldots,k_n) = \left(F_{n-1,i}(k_0,\ldots,k_{n-1}) - F_{n,i}(k_0,\ldots,k_n)\right)F_{n,j}(k_0,\ldots,k_n)$$

If $\{\rho_i^{(m)} : m \ge 0\}$ are defined as in Sect. 2.3.2, then, by (2.1.11),

$$\begin{split} & \mathbb{P}(\rho_{i}^{(m+1)} < \rho_{j} \mid X_{0} = i) \\ &= \sum_{\ell=1}^{\infty} \mathbb{P}(\rho_{i}^{(m)} = \ell \And \rho_{i}^{(m+1)} < \rho_{j} \mid X_{0} = i) \\ &= \sum_{\ell=1}^{\infty} \sum_{n=1}^{\infty} \mathbb{P}(\rho_{i}^{(m)} = \ell, \ \rho_{i}^{(m+1)} = \ell + n < \rho_{j} \mid X_{0} = i) \\ &= \sum_{\ell,n=1}^{\infty} \mathbb{E}[G_{n}(X_{\ell}, \dots, X_{\ell+n}), \rho_{i}^{(m)} = \ell < \rho_{j} \mid X_{0} = i] \\ &= \sum_{\ell,n=1}^{\infty} \mathbb{E}[G_{n}(X_{0}, \dots, X_{n}) \mid X_{0} = i] \mathbb{P}(\rho_{i}^{(m)} = \ell < \rho_{j} \mid X_{0} = i) \\ &= \sum_{\ell,n=1}^{\infty} \mathbb{P}(\rho_{i} = n < \rho_{j} \mid X_{0} = i) \mathbb{P}(\rho_{i}^{(m)} = \ell < \rho_{j} \mid X_{0} = i) \\ &= \mathbb{P}(\rho_{i} < \rho_{j} \mid X_{0} = i) \mathbb{P}(\rho_{i}^{(m)} < \rho_{j} \mid X_{0} = i), \end{split}$$

and so

$$j \neq i \implies \mathbb{P}(\rho_i^{(m)} < \rho_j \mid X_0 = i) = \mathbb{P}(\rho_i < \rho_j \mid X_0 = i)^m.$$
 (3.1.3)

Now suppose that $i \rightarrow j$ but $\mathbb{P}(\rho_j < \rho_i | X_0 = i) = 0$. Then, because

$$\mathbb{P}(\rho_j \neq \rho_i | X_0 = i) \ge \mathbb{P}(\rho_i < \infty | X_0 = i) = 1,$$

 $\mathbb{P}(\rho_i < \rho_j | X_0 = i) = 1$. By (3.1.3), this means that $\mathbb{P}(\rho_i^{(m)} < \rho_j | X_0 = i) = 1$ for all $m \ge 1$, which, since $\rho^{(m)} \ge m$, leads to $\mathbb{P}(\rho_j = \infty | X_0 = i) = 1$ and therefore rules out $i \rightarrow j$. Hence, we have already shown that

$$i \rightarrow j \implies \mathbb{P}(\rho_j < \rho_i | X_0 = i) > 0,$$

and the opposite implication needs no comment.

To prove that $i \rightarrow j \Longrightarrow \mathbb{P}(\rho_i < \infty | X_0 = j) = 1$, first observe that

$$\begin{aligned} &\mathbb{P}(\rho_{j} < \rho_{i} < \infty \mid X_{0} = i) \\ &= \lim_{n \to \infty} \sum_{m=1}^{\infty} \mathbb{P}(\rho_{j} = m < \rho_{i} \le m + n \mid X_{0} = i) \\ &= \lim_{n \to \infty} \sum_{m=1}^{\infty} \mathbb{E} \Big[1 - F_{n,i}(X_{m}, \dots, X_{m+n}), \ \rho_{j} = m < \rho_{i} \mid X_{0} = i \Big] \\ &= \lim_{n \to \infty} \sum_{m=1}^{\infty} \mathbb{P}(\rho_{j} = m < \rho_{i} \mid X_{0} = i) \mathbb{E} \Big[1 - F_{n,i}(X_{0}, \dots, X_{n}) \mid X_{0} = j \Big] \\ &= \mathbb{P}(\rho_{j} < \rho_{i} \mid X_{0} = i) \mathbb{P}(\rho_{i} < \infty \mid X_{0} = j). \end{aligned}$$

Thus, after combining this with $\mathbb{P}(\rho_i < \infty | X_0 = i) = 1$, we have

$$\mathbb{P}(\rho_j < \rho_i \mid X_0 = i) = \mathbb{P}(\rho_j < \rho_i \mid X_0 = i)\mathbb{P}(\rho_i < \infty \mid X_0 = j),$$

which, because $\mathbb{P}(\rho_j < \rho_i | X_0 = i) > 0$, is possible only if $\mathbb{P}(\rho_i < \infty | X_0 = j) = 1$. In particular, we have now proved that $j \rightarrow i$ and therefore that $i \leftrightarrow j$.

Similarly,

$$\mathbb{P}(\rho_j < \infty \mid X_0 = i)$$

= $\mathbb{P}(\rho_j < \rho_i \mid X_0 = i) + \mathbb{P}(\rho_i < \rho_j < \infty \mid X_0 = i)$
= $\mathbb{P}(\rho_j < \rho_i \mid X_0 = i) + \mathbb{P}(\rho_i < \rho_j \mid X_0 = i)\mathbb{P}(\rho_j < \infty \mid X_0 = i),$

and so, since $\mathbb{P}(\rho_i = \rho_j | X_0 = i) \le \mathbb{P}(\rho_i = \infty | X_0 = i) = 0$,

$$\mathbb{P}(\rho_j < \infty \mid X_0 = i)\mathbb{P}(\rho_j < \rho_i \mid X_0 = i) = \mathbb{P}(\rho_j < \rho_i \mid X_0 = i).$$

Thus, $i \to j \Longrightarrow \mathbb{P}(\rho_j < \infty | X_0 = i) = 1$. Finally,

$$\mathbb{P}(\rho_i < \rho_j < \infty \mid X_0 = j) = \mathbb{P}(\rho_j < \infty \mid X_0 = i)\mathbb{P}(\rho_i < \rho_j \mid X_0 = j).$$

Hence, because we now know that

$$i \to j \implies \mathbb{P}(\rho_i < \infty | X_0 = j) = 1 = \mathbb{P}(\rho_j < \infty | X_0 = i),$$

we see that $i \rightarrow j$ implies

$$\mathbb{P}(\rho_j < \infty \mid X_0 = j)$$

= $\mathbb{P}(\rho_j < \rho_i \mid X_0 = j) + \mathbb{P}(\rho_i < \rho_j < \infty \mid X_0 = j)$
= $\mathbb{P}(\rho_j < \rho_i \mid X_0 = j) + \mathbb{P}(\rho_j < \infty \mid X_0 = i)\mathbb{P}(\rho_i < \rho_j \mid X_0 = j) = 1,$

since $\mathbb{P}(\rho_i = \rho_j | X_0 = j) \le \mathbb{P}(\rho_i = \infty | X_0 = j) = 0.$

As an immediate consequence of Theorem 3.1.2 we have the following corollary.

Corollary 3.1.4 If $i \leftrightarrow j$, then j is recurrent (transient) if and only if i is. Moreover, if i is recurrent, then $\mathbb{P}(\rho_j < \infty | X_0 = i)$ is either 1 or 0 according whether or not i communicates with j. In particular, if i is recurrent, then $(\mathbf{P}^n)_{ij} = 0$ for all $n \ge 0$ and all j which do not communicate with i.

When a chain is irreducible, all or none of its states possess any particular communicating class property. Hence, when a chain is irreducible, we will say that it is *recurrent* or *transient* if any one, and therefore all, of its states is.

These considerations allow us to prove the following property of stationary probabilities.

Theorem 3.1.5 If π is a stationary probability for the transition probability \mathbf{P} on \mathbb{S} , then $(\pi)_i = 0$ for all transient $i \in \mathbb{S}$. Furthermore, if i is recurrent and $(\pi)_i > 0$, then $(\pi)_i > 0$ for all j that communicate with i.

Proof First observe that, for any $n \ge 1$,

$$n(\boldsymbol{\pi})_{i} = \sum_{m=0}^{n-1} (\boldsymbol{\pi} \mathbf{P}^{m})_{i} = \sum_{j \in \mathbb{S}} (\boldsymbol{\pi})_{j} \mathbb{E} \left[\sum_{m=0}^{n-1} \mathbf{1}_{\{i\}} (X_{m}) \, \middle| \, X_{0} = j \right]$$
$$\leq \sum_{j \in \mathbb{S}} (\boldsymbol{\pi})_{j} \mathbb{E} \left[\sum_{m=0}^{\infty} \mathbf{1}_{\{i\}} (X_{m}) \, \middle| \, X_{0} = j \right] \leq \mathbb{E} [T_{i} \, \middle| \, X_{0} = i],$$

where the final inequality comes from (2.3.7). Hence, if *i* is transient and therefore $\mathbb{E}[T_i | X_0 = i] < \infty$, then $(\pi)_i = 0$.

Next suppose that *i* is recurrent and that $(\pi)_i > 0$. Then, for any *j* in the communicating class [*i*] of *i*, $\mathbf{P}_{ij}^n > 0$ for some $n \ge 0$, and therefore $(\pi)_j = \sum_{k \in \mathbb{S}} (\pi)_k \mathbf{P}_{ki}^n \ge (\pi)_i \mathbf{P}_{ij}^n > 0$.

3.1.2 Criteria for Recurrence and Transience

There are many tests which can help determine whether a state is recurrent, but no one of them works in all circumstances. In this subsection, we will develop a few

of the most common of these tests. Throughout, we will use **u** to denote the column vector determined by a function $u : \mathbb{S} \longrightarrow \mathbb{R}$.

We begin with a criterion for transience.

Theorem 3.1.6 If *u* is a non-negative function on \mathbb{S} with the property that $(\mathbf{Pu})_i \leq (\mathbf{u})_i$ for all $i \in \mathbb{S}$, then $(\mathbf{Pu})_i < (\mathbf{u})_i$ for some $j \in \mathbb{S}$ implies that *j* is transient.

Proof Set $\mathbf{f} = \mathbf{u} - \mathbf{P}\mathbf{u}$, and note that, for all $n \ge 1$,

$$u(j) \ge (\mathbf{u})_j - (\mathbf{P}^n \mathbf{u})_j = \sum_{m=0}^{n-1} ((\mathbf{P}^m \mathbf{u})_j - (\mathbf{P}^{m+1} \mathbf{u})_j)$$
$$= \sum_{m=0}^{n-1} (\mathbf{P}^m \mathbf{f})_j \ge (\mathbf{f})_j \sum_{m=0}^{n-1} (\mathbf{P}^m)_{jj}$$

Thus $\mathbb{E}[T_j|X_0 = j] = \sum_{m=0}^{\infty} (\mathbf{P}^m)_{jj} \le \frac{u(j)}{(\mathbf{f})_j} < \infty$, which, by (2.3.7), means that j is transient.

In order to prove our next criterion, we will need the following special case of a general result known as *Doob's stopping time theorem*.

Lemma 3.1.7 Assume that $u : \mathbb{S} \longrightarrow \mathbb{R}$ is bounded below and that Γ is a nonempty subset of \mathbb{S} . If $(\mathbf{Pu})_i \le u(i)$ for all $i \notin \Gamma$ and $\rho_{\Gamma} \equiv \inf\{n \ge 1 : X_n \in \Gamma\}$, then

$$\mathbb{E}\left[u(X_{n\wedge\rho_{\Gamma}}) \mid X_0=i\right] \leq u(i) \quad for \ all \ n \geq 0 \ and \ i \in \mathbb{S} \setminus \Gamma.$$

Moreover, if the inequality in the hypothesis is replaced by equality, then the inequality in the conclusion can be replaced by equality.

Proof Set $A_n = \{\rho_{\Gamma} > n\}$. Then, A_n is measurable with respect to (X_0, \ldots, X_n) , and so, by (2.1.1), for any $i \notin \Gamma$,

$$\mathbb{E}\left[u(X_{(n+1)\wedge\rho_{\Gamma}}) \mid X_{0}=i\right]$$

$$=\mathbb{E}\left[u(X_{n\wedge\rho_{\Gamma}}), A_{n}\mathbb{C} \mid X_{0}=i\right] + \sum_{k\notin\Gamma} \mathbb{E}\left[u(X_{n+1}), A_{n}\cap\{X_{n}=k\} \mid X_{0}=i\right]$$

$$=\mathbb{E}\left[u(X_{n\wedge\rho_{\Gamma}}), A_{n}\mathbb{C} \mid X_{0}=i\right] + \sum_{k\notin\Gamma} \mathbb{E}\left[(\mathbf{Pu})_{k}, A_{n}\cap\{X_{n}=k\} \mid X_{0}=i\right]$$

$$\leq \mathbb{E}\left[u(X_{n\wedge\rho_{\Gamma}}), A_{n}\mathbb{C} \mid X_{0}=i\right] + \mathbb{E}\left[u(X_{n\wedge\rho_{\Gamma}}), A_{n} \mid X_{0}=i\right]$$

$$=\mathbb{E}\left[u(X_{n\wedge\rho_{\Gamma}}) \mid X_{0}=i\right].$$

Clearly, the same argument works just as well in the case of equality.

Theorem 3.1.8 Assume that *j* is recurrent, and set C = [j]. If $u : \mathbb{S} \longrightarrow [0, \infty)$ is a bounded function and either $u(i) = (\mathbf{Pu})_i$ or $u(j) \ge u(i) \ge (\mathbf{Pu})_i$ for all $i \in C \setminus \{j\}$, then *u* is constant on *C*. On the other hand, if *j* is transient, then the function *u* given by

$$u(i) = \begin{cases} 1 & \text{if } i = j \\ \mathbb{P}(\rho_j < \infty | X_0 = i) & \text{if } i \neq j \end{cases}$$

is a non-negative, bounded, non-constant solution to $u(i) = (\mathbf{Pu})_i$ for all $i \neq j$.

Proof In proving the first part, we will assume, without loss in generality, that C = S. Now suppose that j is recurrent and that $u(i) = (\mathbf{Pu})_i$ for $i \neq j$. By applying Lemma 3.1.7 with $\Gamma = \{j\}$, we see that, for $i \neq j$,

$$u(i) = u(j)\mathbb{P}(\rho_j \le n \,|\, X_0 = i) + \mathbb{E}[u(X_n), \rho_j > n \,|\, X_0 = i].$$

Hence, since, by Theorem 3.1.2, $\mathbb{P}(\rho_j < \infty | X_0 = i) = 1$ and *u* is bounded, we get u(i) = u(j) after letting $n \to \infty$. Next assume that $u(j) \ge u(i) \ge (\mathbf{Pu})_i$ for all $i \ne j$. Then, again by Lemma 3.1.7, we have

$$u(j) \ge u(i) \ge u(j)\mathbb{P}(\rho_j \le n \mid X_0 = i) + \mathbb{E}[u(X_n), \rho_j > n \mid X_0 = i],$$

which leads to the required conclusion when $n \to \infty$.

To prove the second part, let u be given by the prescription described, and begin by observing that, because j is transient,

$$1 > \mathbb{P}(\rho_j < \infty \mid X_0 = j) = \mathbf{P}_{jj} + \sum_{i \neq j} \mathbf{P}_{ji} u(i) \ge \mathbf{P}_{jj} + (1 - \mathbf{P}_{jj}) \inf_{i \neq j} u(i)$$

Because *j* is transient, and therefore $(\mathbf{P})_{jj} < 1$, this proves that $\inf_{i \neq j} u(i) < u(j)$ and therefore that *u* is non-constant. At the same time, when $i \neq j$, by conditioning on what happens at time 1, we know that

$$u(i) = \mathbb{P}(\rho_j < \infty \mid X_0 = i) = \mathbf{P}_{ij} + \sum_{k \neq j} \mathbf{P}_{ik} \mathbb{P}(\rho_j < \infty \mid X_0 = k) = (\mathbf{P}\mathbf{u})_i.$$

Lemma 3.1.9 If **P** is irreducible on \mathbb{S} , then, for any finite subset $F \neq \mathbb{S}$, $\mathbb{P}(\rho_{\mathbb{S}\setminus F} < \infty | X_0 = i) = 1$ for all $i \in F$.

Proof Set $\tau = \rho_{S\setminus F}$. By irreducibility, $\mathbb{P}(\tau < \infty | X_0 = i) > 0$ for each $i \in F$. Hence, because F is finite, there exists a $\theta \in (0, 1)$ and an $N \ge 1$ such that $\mathbb{P}(\tau > N | X_0 = i) \le \theta$ for all $i \in F$. But this means that, for each $i \in F$, $\mathbb{P}(\tau > (\ell + 1)N | X_0 = i)$ equals

$$\sum_{k \in F} \mathbb{P} (\tau > (\ell + 1)N \& X_{\ell N} = k | X_0 = i)$$

= $\sum_{k \in F} \mathbb{P} (X_n \in F \text{ for } \ell N + 1 \le n \le (\ell + 1)N, \ \tau > \ell N, \ \& X_{\ell N} = k | X_0 = i)$
= $\sum_{k \in F} \mathbb{P} (\tau > N | X_0 = k) \mathbb{P} (\tau > \ell N \& X_{\ell N} = k | X_0 = i)$
 $\le \theta \mathbb{P} (\tau > \ell N | X_0 = i).$

Thus, $\mathbb{P}(\tau > \ell N | X_0 = i) \le \theta^{\ell}$, and so $\mathbb{P}(\tau = \infty | X_0 = j) = 0$ for all $i \in F$. \Box

Theorem 3.1.10 Assume that **P** is irreducible on \mathbb{S} , and let $u : \mathbb{S} \longrightarrow [0, \infty)$ be a function with the property that $\{k : u(k) \le L\}$ is finite for each $L \in (0, \infty)$. If, for some $j \in \mathbb{S}$, $(\mathbf{Pu})_i \le u(i)$ for all $i \ne j$, the chain determined by **P** is recurrent on \mathbb{S} .

Proof If S is finite, then (cf., for example, Exercise 2.4.2) at least one state is recurrent, and therefore, by irreducibility, all are. Hence, we will assume that S is infinite.

Given $i \neq j$, set $F_L = \{k : u(k) \leq u(i) + u(j) + L\}$ for $L \in \mathbb{N}$, and denote by ρ_L the first return time $\rho_{(\mathbb{S} \setminus F_L) \cup \{j\}}$ to $(\mathbb{S} \setminus F_L) \cup \{j\}$. By Lemma 3.1.6,

$$u(i) \ge \mathbb{E}\left[u(X_{n \land \rho_L}) \mid X_0 = i\right] \ge \left(u(i) + u(j) + L\right) \mathbb{P}(\rho_{\mathbb{S} \backslash F_L} < n \land \rho_j \mid X_0 = i)$$

for all $n \ge 1$. Hence, after letting $n \to \infty$, we conclude that, for all $L \in \mathbb{N}$,

$$u(i) \ge (u(i) + u(j) + L) \mathbb{P}(\rho_{\mathbb{S}\setminus F_L} < \rho_j \mid X_0 = i)$$

$$\ge (u(i) + u(j) + L) \mathbb{P}(\rho_j = \infty \mid X_0 = i),$$

since, by Lemma 3.1.8, we know that $\mathbb{P}(\rho_{\mathbb{S}\setminus F_L} < \infty | X_0 = i) = 1$. Thus, after letting $L \to \infty$, we have now shown that $\mathbb{P}(\rho_i < \infty | X_0 = i) = 1$ for all $i \neq j$. Since

$$\mathbb{P}(\rho_j < \infty \mid X_0 = j) = (\mathbf{P})_{jj} + \sum_{i \neq j} \mathbb{P}(\rho_j < \infty \mid X_0 = i)(\mathbf{P})_{ji},$$

it follows that $\mathbb{P}(\rho_j < \infty | X_0 = j) = 1$, which means that *j* is recurrent.

Remark The preceding criteria are examples, of which there are many others, that relate recurrence of $j \in S$ to the existence or non-existence of certain types of functions that satisfy either $\mathbf{Pu} = u$ or $\mathbf{Pu} \le u$ on $S \setminus \{j\}$. All these criteria can be understood as mathematical implementations of the intuitive idea that

$$u(i) = (\mathbf{Pu})_i$$
 or $(\mathbf{Pu})_i \le u(i)$ for $i \ne j$,

implies that as long as $X_n \neq j$, $u(X_n)$ will be "nearly constant" or "nearly nonincreasing" as *n* increases. The sense in which these "nearly's" should be interpreted is the subject of *martingale theory*, and our proofs of these criteria would have been simplified had we been able to call on that theory.

3.1.3 Periodicity

Periodicity is another important communicating class property. In order to describe this property, we must recall Euclid's concept of the *greatest common divisor* gcd(S) of a non-empty subset $S \subseteq \mathbb{Z}$. Namely, we say that $d \in \mathbb{Z}^+$ is a common divisor of *S* and write d|S if $\frac{s}{d} \in \mathbb{Z}$ for every $s \in S$. Clearly, if $S = \{0\}$, then d|S for every $d \in \mathbb{Z}^+$, and so we take $gcd(\{0\}) = \infty$. On the other hand, if $S \neq \{0\}$, then no common divisor of *S* can be larger than $\min\{|s| : s \in S \setminus \{0\}\}$, and so we know that $gcd(S) < \infty$.

Our interest in this concept comes from the role it plays in the ergodic theory of Markov chains. Namely, as we will see in Chap. 4, it allows us to distinguish between the chains for which powers of the transition probability matrix converge and those for which it is necessary to take averages. More precisely, given a state i, set

$$S(i) = \{n \ge 0 : (\mathbf{P}^n)_{ii} > 0\} \text{ and } d(i) = \gcd(S(i)).$$
(3.1.11)

Then d(i) is called the *period* of the state *i*, and, *i* is said to be *aperiodic* if d(i) = 1. In Chap. 4, we will see that averaging is required unless *i* is aperiodic. However, here we will only take care of a few more mundane matters. In the first place, the period is a communicating class property:

$$i \leftrightarrow j \implies d(i) = d(j).$$
 (3.1.12)

To see this, assume that $(\mathbf{P}^m)_{ij} > 0$ and $(\mathbf{P}^n)_{ji} > 0$, and let *d* be a common divisor of *S*(*i*). Then for any $k \in S(j)$,

$$\left(\mathbf{P}^{m+k+n}\right)_{ii} \ge \left(\mathbf{P}^{m}\right)_{ij} \left(\mathbf{P}^{k}\right)_{jj} \left(\mathbf{P}^{n}\right)_{ji} > 0,$$

and so $m + k + n \in S(i)$. Hence $d|\{m + k + n : k \in S(j)\}$. But, because $m + n \in S(i)$, and therefore d|(m + n), this is possible only if d divides S(j), and so we now know that $d(i) \leq d(j)$. After reversing the roles of i and j, one sees that $d(j) \leq d(i)$, which means that d(i) must equal d(j).

We next need the following elementary fact from number theory.

Theorem 3.1.13 Given $\emptyset \neq S \subseteq \mathbb{Z}$ with $S \neq \{0\}$,

$$\gcd(S) \le \min\{|s| : s \in S \setminus \{0\}\}$$

and equality holds if and only if $\{\gcd(S), -\gcd(S)\} \cap S \neq \emptyset$. More generally, there always exists an $M \in \mathbb{Z}^+$, $\{a_m\}_1^M \subseteq \mathbb{Z}$, and $\{s_m\}_1^M \subseteq S$ such that $\gcd(S) = \sum_{i=1}^{M} a_m s_m$. Finally, if $S \subseteq \mathbb{N}$ and $(s_1, s_2) \in S^2 \Longrightarrow s_1 + s_2 \in S$, then there exists an $M \in \mathbb{Z}^+$ such that

$$\left\{m \operatorname{gcd}(S) : m \ge M\right\} = \left\{s \in S : s \ge M \operatorname{gcd}(S)\right\}.$$

Proof The first assertion needs no comment. To prove the second assertion, let \hat{S} be the smallest subset of \mathbb{Z} which contains S and has the property that $(s_1, s_2) \in \hat{S}^2 \Longrightarrow s_1 \pm s_2 \in \hat{S}$. As is easy to check, \hat{S} coincides with the subset of \mathbb{Z} whose elements can be expressed in the form $\sum_{1}^{M} a_m s_m$ for some $M \ge 1$, $\{a_m\}_{1}^{M} \subseteq \mathbb{Z}$, and $\{s_m\}_{1}^{M} \subseteq S$. In particular, this means that $gcd(S)|\hat{S}$, and so $gcd(S) \le gcd(\hat{S})$. On the other hand, because $S \subseteq \hat{S}$, $gcd(\hat{S})|S$. Hence, $gcd(S) = gcd(\hat{S})$, and so, by the first part, we will be done once we show that $gcd(\hat{S}) \in \hat{S}$. To this end, let $m = \min\{s \in \mathbb{Z}^{+} : s \in \hat{S}\}$. We already know that $gcd(\hat{S}) \le m$. Thus, to prove that $gcd(\hat{S}) \in \hat{S}$, we can write s = am + r for some $(a, r) \in \mathbb{Z}^2$ with $0 \le r < m$. In particular, $r = s - am \in \hat{S}$. Hence, if $r \ne 0$, then r would contradict the condition that m is the smallest positive element of \hat{S} .

To prove the final assertion, first note that it suffices to prove that there is an $M \in \mathbb{Z}^+$ such that $\{m \operatorname{gcd}(S) : m \ge M\} \subseteq S$. To this end, begin by checking that, under the stated hypothesis, $\hat{S} = \{s_2 - s_1 : s_1, s_2 \in S \cup \{0\}\}$. Thus $\operatorname{gcd}(S) = s_2 - s_1$ for some $s_1 \in S \cup \{0\}$ and $s_2 \in S \setminus \{0\}$. If $s_1 = 0$, then $m \operatorname{gcd}(S) = ms_2 \in S$ for all $m \in \mathbb{Z}^+$, and so we can take M = 1. If $s_1 \neq 0$, choose $a \in \mathbb{Z}^+$ so that $s_1 = a \operatorname{gcd}(S)$. Then, for any $(m, r) \in \mathbb{N}^2$ with $0 \le r < a$,

$$(a^2 + ma + r)$$
gcd $(S) = ms_1 + rs_2 + (a - r)s_1 = (m + a - r)s_1 + rs_2 \in S.$

Hence, after another application of the Euclidean algorithm, we see that we can take $M = a^2$.

As an immediate consequence of Theorem 3.1.13, we see that

$$d(i) < \infty \implies (\mathbf{P}^{nd(i)})_{ii} > 0 \text{ for all sufficiently large } n \in \mathbb{Z}^+.$$
 (3.1.14)

In particular,¹

i is aperiodic
$$\iff (\mathbf{P}^n)_{ii} > 0$$
 for all sufficiently large $n \in \mathbb{Z}^+$. (3.1.15)

We close this subsection with an application of these considerations to stationary probabilities of transition probabilities on a finite state space.

Corollary 3.1.16 Suppose that **P** is an transition probability matrix on a finite state space S. If there is an aperiodic state $j_0 \in S$ such that $i \rightarrow j_0$ for every $i \in S$, then there exists an $M \in \mathbb{Z}^+$ and an $\epsilon > 0$ such that $(\mathbf{P}^M)_{ij_0} \ge \epsilon$ for all $i \in S$. In particular, if π is the unique stationary probability guaranteed by Theorem 2.2.1, then

$$\| \boldsymbol{\mu} \mathbf{P}^n - \boldsymbol{\pi} \|_{\mathbf{v}} \leq 2(1-\epsilon)^{\lfloor \frac{n}{M} \rfloor}$$
 for all $n \in \mathbb{Z}^+$ and initial distributions $\boldsymbol{\mu}$.

¹The "if" part of the following statement depends on the existence of infinitely many prime numbers.

Proof Because j_0 is aperiodic, we know that there is an $M_0 \in \mathbb{N}$ such that $(\mathbf{P}^n)_{j_0j_0} > 0$ for all $n \ge M_0$. Further, because $i \to j_0$, there exists an $m(i) \in \mathbb{Z}^+$ such that $(\mathbf{P}^{m(i)})_{ij_0} > 0$. Hence, $(\mathbf{P}^n)_{ij_0} > 0$ for all $n \ge m(i) + M_0$. Finally, take $M = M_0 + \max_{i \in \mathbb{S}} m(i), \epsilon = \min_{i \in \mathbb{S}} (\mathbf{P}^M)_{ij_0}$, and apply (2.2.3).

3.2 Computation of Stationary Probabilities

It is important to know how to compute stationary probabilities, and in this section we will develop a procedure that works when the state space is finite.

3.2.1 Preliminary Results

Whether or not the state space is finite, a stationary probability for a transition probability **P** is a probability vector π that satisfies $\pi \mathbf{P} = \pi$. Equivalently, $\pi (\mathbf{I} - \mathbf{P}) = \mathbf{0}$ and $\pi \mathbf{1} = 1$. If one ignores the necessity of finding a solution with non-negative components, when the state space \mathbb{S} is finite general principles of linear algebra guarantee that one can always solve these equations. Indeed, $(\mathbf{I} - \mathbf{P})\mathbf{1} = \mathbf{0}$ and therefore²

$$\det((\mathbf{I} - \mathbf{P})^{\top}) = \det((\mathbf{I} - \mathbf{P})) = 0,$$

which means that there must exist a row vector **v** such that $\mathbf{v}(\mathbf{I} - \mathbf{P}) = \mathbf{0}$ with $\mathbf{v1} = 1$. For us, a better way to see that solutions exist is the following. Clearly, the set of probability vectors on \mathbb{S} is a compact subset of \mathbb{R}^N , where *N* is the cardinality of \mathbb{S} . Thus, for any probability vector $\boldsymbol{\mu}$, the sequence { $\boldsymbol{\mu}\mathbf{A}_M : M \ge 1$ } admits a subsequence { $\boldsymbol{\mu}\mathbf{A}_{M_k} : k \ge 1$ } that converges to a probability vector $\boldsymbol{\pi}$. Moreover,

$$\boldsymbol{\pi} \mathbf{P} = \lim_{k \to \infty} \frac{1}{M_k} \sum_{m=0}^{M_k - 1} \boldsymbol{\mu} \mathbf{P}^{m+1} = \lim_{k \to \infty} \left(\frac{\boldsymbol{\mu} \mathbf{P}^{M_k} - \boldsymbol{\mu}}{M_k} + \boldsymbol{\mu} \mathbf{A}_{M_k} \right) = \boldsymbol{\pi}.$$

Hence π is a stationary probability for **P**.

The following result, which contains a partial converse of Theorem 2.2.5, exploits the preceding construction.

Theorem 3.2.1 Suppose that **P** is a transition probability on a finite state space S. Then there is at least one recurrent $i \in S$. In addition, if j is transient, then $j \rightarrow i$ for some recurrent i. Finally, suppose that **P** admits precisely one stationary probability π . Then all recurrent states communicate with one another and there is

²If **A** is a matrix, then \mathbf{A}^{\top} is its transpose. That is, $\mathbf{A}_{ij}^{\top} = \mathbf{A}_{ji}$. Thus, if **v** is a row vector, then \mathbf{v}^{\top} is a column vector and $\mathbf{v}\mathbf{A} = \mathbf{A}^{\top}\mathbf{v}^{\top}$.

an $M \in \mathbb{Z}^+$ and an $\epsilon > 0$ such that $(\mathbf{A}_M)_{ji} \ge \epsilon$ for all $j \in \mathbb{S}$ and all recurrent $i \in \mathbb{S}$. Thus, when \mathbb{S} is finite, \mathbf{P} admits precisely one stationary measure if and only if \mathbf{A}_M satisfies Doeblin's condition for some $M \ge 1$.

Proof Begin by noting that, no matter where it starts, the chain can spend only a finite amount of time in the set \mathcal{T} of transient states. Indeed, if $i \in \mathbb{S}$ and $j \in \mathcal{T}$, then (cf. (2.3.7)) $\mathbb{E}[T_j | X_0 = i] < \infty$ and therefore, since \mathcal{T} is finite, $\mathbb{E}[\sum_{n=0}^{\infty} \mathbf{1}_{\mathcal{T}}(X_n) | X_0 = i] < \infty$. Hence, there must exist at least one recurrent state, and the chain must visit one of them no matter where it starts. In particular, if j is transient, then there is a recurrent i for which $j \to i$.

Now assume that there is only one stationary probability π for **P**. We will show first that all the recurrent states must communicate with one another. To see this, suppose that i_1 and i_2 were recurrent states that do not communicate. Then $[i_1] \cap [i_2] = \emptyset$ and the restrictions of **P** to $[i_1]$ and $[i_2]$ are transition probabilities. Hence, by the preceding, we would be able to construct stationary probabilities π_1 and π_2 for these restrictions. After extending π_k to \mathbb{S} by taking it to be 0 off of $[i_k]$, we would then have two unequal stationary probabilities for **P**. Thus, $i_1 \leftrightarrow i_2$.

We now know that if *i* is a recurrent state, then $j \rightarrow i$ for all *j*. In other words, for each $j \in \mathbb{S}$ and recurrent *i*, $(\mathbf{P}^n)_{ji} > 0$ for some *n*. Thus, because \mathbb{S} is finite, this completes the proof that *M* and ϵ exist. As for the concluding assertion, it follows immediately from the result just proved combined with Theorem 2.2.5.

3.2.2 Computations via Linear Algebra

Our next goal is to find a formula for the stationary probability when there is only one, and the procedure that we will give uses elementary linear algebra. Recall that if **A** is a square matrix indexed by an *N* element set \mathbb{S} , then its *eigenvalues* are the roots of its characteristic polynomial $\lambda \in \mathbb{C} \mapsto \det(\lambda I - \mathbf{A}) \in \mathbb{C}$. An eigenvalue that is different from all the other is said to be a *simple eigenvalue*.

Given an $N \times N$ -matrix **A**, let cof(**A**) denote its *cofactor matrix*. That is, cof(**A**)_{*ij*} is $(-1)^{i+j}$ times the determinant of the $(N-1) \times (N-1)$ -matrix obtained by removing the *i*th column and the *j*th row from **A**. By Cramer's formula, Acof(**A**) = cof(**A**)**A** equals det(**A**) times the identity matrix. From this it is easy to see that cof(**A**)_{*ji*} is equal to the derivative of det(**A**) with respect to its (i, j)th entry (**A**)_{*ij*}.

Lemma 3.2.2 If each row of **A** sums to 0, then $cof(\mathbf{A})_{ij} = cof(\mathbf{A})_{jj}$ for all $1 \le i, j \le N$ and

$$\sum_{i=1}^{N} \operatorname{cof}(\mathbf{A})_{ii} = \Pi_{\mathbf{A}},$$

where $\Pi_{\mathbf{A}}$ equals the product of the non-zero eigenvalues of \mathbf{A} if 0 is a simple eigenvalue of \mathbf{A} and is equal to 0 otherwise.

Proof Given $i_1 \neq i_2$ and j, let **B** be the matrix whose (j, i_1) st entry is 1, whose (j, i_2) th entry is -1, and whose other entries are all 0. Then, for each $t \in \mathbb{R}$, all the rows of $\mathbf{A} + t\mathbf{B}$ sum to 0 and therefore det $(\mathbf{A} + t\mathbf{B}) = 0$. Hence

$$\operatorname{cof}(\mathbf{A})_{i_2 j} - \operatorname{cof}(\mathbf{A})_{i_1 j} = \frac{d}{dt} \operatorname{det}(\mathbf{A} + t\mathbf{B}) \bigg|_{t=0} = 0.$$

To prove the second assertion, consider the polynomial $\lambda \rightarrow \det(\lambda \mathbf{I} + \mathbf{A})$. Because $\det(\mathbf{A}) = 0$, this polynomial vanishes at 0. Furthermore, using Cramer's formula, it is easy to check that $\sum_{i=1}^{N} \operatorname{cof}(\mathbf{A})_{ii}$ is the coefficient of λ . At the same time, $\det(\lambda \mathbf{I} + \mathbf{A}) = \prod_{k=1}^{N} (\lambda + \lambda_k)$, where $\lambda_1, \ldots, \lambda_N$ are the eigenvalues of \mathbf{A} . Hence, since 0 is one of these eigenvalues, we can take $\lambda_1 = 0$ and thereby see that

$$\prod_{k=2}^{N} \lambda_k = \lim_{\lambda \to 0} \frac{\det(\lambda \mathbf{I} + \mathbf{A})}{\lambda} = \sum_{i=1}^{N} \operatorname{cof}(\mathbf{A})_{ii}.$$

Theorem 3.2.3 Assume that each row of **A** sums to 0 and that 0 is a simple eigenvalue of **A**. Determine the row vector $\boldsymbol{\pi}$ by

$$(\boldsymbol{\pi})_i = \frac{\operatorname{cof}(\mathbf{A})_{ii}}{\Pi_{\mathbf{A}}}, \quad 1 \le i \le N,$$

where $\Pi_{\mathbf{A}}$ is the product of the non-zero eigenvalues of \mathbf{A} . Then π is the unique row vector $\mathbf{v} \in \mathbb{C}^N$ such that $\mathbf{v}\mathbf{A} = 0$ and $\sum_{i=1}^N (\mathbf{v})_i = 1$.

Proof Since the null space of \mathbf{A}^{\top} has the same dimension as that of \mathbf{A} and the latter has dimension 1, we know that there is at most one **v**. Furthermore, by Lemma 3.2.2, we know that $\sum_{i=1}^{N} (\boldsymbol{\pi})_i = 1$. Finally, by that same lemma and Cramer's formula,

$$\Pi_{\mathbf{A}} \sum_{i=1}^{N} (\boldsymbol{\pi})_{i} a_{ij} = \sum_{i=1}^{N} \operatorname{cof}(\mathbf{A})_{ji} a_{ij} = \det(\mathbf{A}) = 0 \quad \text{for } 1 \le j \le N.$$

In what follows, it will be useful to have introduced the following notation. Given a matrix **A** indexed by some finite set \mathbb{S} and $\Delta \subseteq \mathbb{S}$, when $\Delta \neq \mathbb{S}$ let \mathbf{A}^{Δ} be the matrix indexed by $\mathbb{S} \setminus \Delta$ obtained by eliminating the rows and columns with indices in Δ and when $\Delta = \mathbb{S}$ define $\mathbf{A}^{\Delta} = 1$. In this notation, for each $i \in \mathbb{S}$, det($\mathbf{A}^{\{i\}}$) equals cof(\mathbf{A})_{*ii*}, and therefore, if $i \notin \Delta$ and det(\mathbf{A}^{Δ}) $\neq 0$, then Cramer's formula says that

$$((\mathbf{A}^{\Delta})^{-1})_{ii} = \frac{\det(\mathbf{A}^{\Delta \cup \{i\}})}{\det(\mathbf{A}^{\Delta})}.$$

More generally, if $\Delta \subsetneq \Delta'$ and (i_1, \ldots, i_K) is an ordering of the elements of $\Delta' \setminus \Delta$, set $\Delta(1) = \Delta$, $\Delta(k) = \Delta \cup \{i_1, \ldots, i_{k-1}\}$ for $2 \le k \le K$, and assume that

 $det(\mathbf{A}^{\Delta(k)}) \neq 0$ for all $1 \leq k \leq K$. Then, by induction on K, one sees that

$$\frac{\det(\mathbf{A}^{\Delta'})}{\det(\mathbf{A}^{\Delta})} = \prod_{k=1}^{K} \left(\left(\mathbf{A}^{\Delta(k)} \right)^{-1} \right)_{i_k i_k}.$$
(3.2.4)

In conjunction with the next lemma, Theorem 3.2.3 will give us a formula for computing stationary probabilities when there is only one.

Lemma 3.2.5 Let **P** be a transition probability on a finite state space S. Then all the eigenvalues of **P** lie inside the closed unit disk \mathbb{D} in \mathbb{C} centered at 0. Also, det $((\mathbf{I} - \mathbf{P})^{[i]}) \ge 0$ for all $i \in \mathbb{S}$. Finally, there is only one stationary probability π for **P** if and only if 0 is a simple eigenvalue for $\mathbf{I} - \mathbf{P}$, in which case $(\pi)_i > 0$ if and only if i is recurrent.

Proof Let *N* be the cardinality of S.

Suppose that λ is an eigenvalue of **P**. Then there is a non-zero $\mathbf{v} \in \mathbb{C}^N$ satisfying $\|\mathbf{v}\|_{\mathbf{u}} = 1$ and $\mathbf{P}\mathbf{v} = \lambda \mathbf{v}$. Hence $1 \ge \|\mathbf{P}\mathbf{v}\|_{\mathbf{u}} \ge |\lambda|$.

Knowing that the eigenvalues of **P** lie in \mathbb{D} , we know that $\alpha \mathbf{I} - \mathbf{P}$ is invertible for all $\alpha > 1$. Furthermore, for any $\alpha > 1$ and **v**, the series $\sum_{m=0}^{\infty} \alpha^{-m-1} \mathbf{P}^m \mathbf{v}$ converges in uniform norm to a limit \mathbf{v}_{α} which satisfies $(\alpha \mathbf{I} - \mathbf{P})\mathbf{v}_{\alpha} = \mathbf{v}$. Hence, $\mathbf{v}_{\alpha} = (\alpha \mathbf{I} - \mathbf{P})^{-1}\mathbf{v}$. In particular, this means that

$$((\alpha \mathbf{I} - \mathbf{P})^{-1})_{ii} = \sum_{m=0}^{\infty} \alpha^{-m-1} (\mathbf{P}^m)_{ii} > 0 \text{ for all } \alpha > 1 \text{ and } i \in \mathbb{S}.$$

In addition, because all the real eigenvalues of $\mathbf{I} - \mathbf{P}$ are non-negative and the ones with non-zero imaginary parts must come in conjugate pairs, we have that $\det(\alpha \mathbf{I} - \mathbf{P}) > 0$, and therefore (3.2.4) now implies that

$$\det((\alpha \mathbf{I} - \mathbf{P})^{\{i\}}) = ((\alpha \mathbf{I} - \mathbf{P})^{-1})_{ii} \det(\alpha \mathbf{I} - \mathbf{P}) > 0 \text{ for all } \alpha > 1 \text{ and } i \in \mathbb{S}.$$

Hence, after letting $\alpha \searrow 1$, we conclude that $\det((\mathbf{I} - \mathbf{P})^{\{i\}}) \ge 0$.

If 0 is a simple eigenvalue of $\mathbf{I} - \mathbf{P}$, then the null space of $\mathbf{I} - \mathbf{P}$ and, therefore of $(\mathbf{I} - \mathbf{P})^{\top}$, is one dimensional. Hence, there can be only one stationary probability.

Now assume that there is only one stationary probability $\boldsymbol{\pi}$. If 0 were not a simple eigenvalue of $\mathbf{I} - \mathbf{P}$, then there would exist a non-zero column vector $\mathbf{v} \in \mathbb{R}^{\mathbb{S}}$ and an $1 \leq n \leq N$ such that $\boldsymbol{\pi} \mathbf{v} = 0$, $(\mathbf{I} - \mathbf{P})^{n-1}\mathbf{v} \neq \mathbf{0}$, and $(\mathbf{I} - \mathbf{P})^n\mathbf{v} = \mathbf{0}$. If n = 1, then $\mathbf{v} = \mathbf{A}_M \mathbf{v}$ for all $M \geq 1$, and so, by Theorems 3.2.1 and 2.2.5, we would get the contradiction $\mathbf{v} = \boldsymbol{\pi} \mathbf{v} \mathbf{1} = \mathbf{0}$. Next suppose that $n \geq 2$. For $\alpha > 1$, construct \mathbf{v}_{α} as above. Then $\|\mathbf{v}_{\alpha}\|_{u} \leq (\alpha - 1)^{-1} \|\mathbf{v}\|_{u}$. At the same time, if $\mathbf{w} = \sum_{m=0}^{n-1} (-1)^m (\alpha - 1)^m (\mathbf{I} - \mathbf{P})^m \mathbf{v}$, then, because $(\mathbf{I} - \mathbf{P})^n \mathbf{v} = \mathbf{0}$, $(\alpha \mathbf{I} - \mathbf{P})\mathbf{w} = \mathbf{v}$, and therefore $\mathbf{v}_{\alpha} = \mathbf{w}$. But this means that

$$\begin{aligned} (\alpha - 1)^{-1} \|\mathbf{v}\|_{\mathbf{u}} &\geq \left\| \sum_{m=0}^{n-1} (-1)^m (\alpha - 1)^{-m-1} (\mathbf{I} - \mathbf{P})^m \mathbf{v} \right\|_{\mathbf{u}} \\ &\geq (\alpha - 1)^{-n} \left\| (\mathbf{I} - \mathbf{P})^{n-1} \mathbf{v} \right\| - \sum_{m=0}^{n-2} (\alpha - 1)^{-m-1} \left\| (\mathbf{I} - \mathbf{P})^m \mathbf{v} \right\|_{\mathbf{v}} \end{aligned}$$

for all $\alpha > 1$. Hence, after multiplying both sides by $(\alpha - 1)^n$ and letting $\alpha \searrow 1$, we would get the contradiction that $(\mathbf{I} - \mathbf{P})^{n-1}\mathbf{v} = \mathbf{0}$.

Turning to the final assertion, Theorem 3.1.5 says that $(\pi)_i = 0$ if *i* is transient. Hence, all that remains is to check that $(\pi)_i > 0$ if *i* is recurrent. But we know that $(\pi)_i > 0$ for some recurrent *i*, and, by Theorem 3.2.1, we know that all recurrent states communicate with one another. Hence, the desired conclusion follows from Theorem 3.1.5.

Theorem 3.2.6 Let **P** be a transition probability on a finite state space S, and assume that π is the one and only stationary probability for **P**. Then 0 is a simple eigenvalue of $\mathbf{I} - \mathbf{P}$ and, if $\Pi_{\mathbf{I}-\mathbf{P}}$ is the product of the other eigenvalues of $\mathbf{I} - \mathbf{P}$, then $\Pi_{\mathbf{I}-\mathbf{P}} > 0$ and

$$(\boldsymbol{\pi})_i = \frac{\det((\mathbf{I} - \mathbf{P})^{\{i\}})}{\Pi_{\mathbf{I} - \mathbf{P}}} = \frac{\det((\mathbf{I} - \mathbf{P})^{\{i\}})}{\sum_{j \in \mathbb{S}} \det((\mathbf{I} - \mathbf{P})^{\{j\}})} \quad \text{for all } i \in \mathbb{S}.$$
(3.2.7)

In particular, *i* is transient if and only if $det((\mathbf{I} - \mathbf{P})^{\{i\}}) = 0$ and *i* is recurrent if and only if $det((\mathbf{I} - \mathbf{P})^{\{i\}}) > 0$.

Proof To prove (3.2.7), take $\mathbf{A} = \mathbf{I} - \mathbf{P}$. Clearly the rows of \mathbf{A} each sum to 0, and, by Lemma 3.2.5, 0 is a simple eigenvalue of \mathbf{A} . In addition, again by Lemma 3.2.5, the π in (3.2.7) is a probability vector and, by Theorem 3.2.3, it is the one and only stationary probability for \mathbf{P} .

The following theorem can sometimes simplify calculations.

Corollary 3.2.8 Let **P** and π be as in Theorem 3.2.6, and denote by \mathcal{T} the set of transient states. Then $(\pi)_i = 0$ for $i \in \mathcal{T}$ and

$$(\boldsymbol{\pi})_i = \frac{\det((\mathbf{I} - \mathbf{P})^{\mathcal{T} \cup \{i\}})}{\Pi_{(\mathbf{I} - \mathbf{P})^{\mathcal{T}}}} = \frac{\det((\mathbf{I} - \mathbf{P})^{\mathcal{T} \cup \{i\}})}{\sum_{j \in \mathbb{S}} \det((\mathbf{I} - \mathbf{P})^{\mathcal{T} \cup \{j\}})} \quad \text{for } i \in \mathbb{S} \setminus \mathcal{T}.$$

Proof First notice that $\mathbf{P}^{\mathcal{T}}$ is a transition probability on $\mathbb{S} \setminus \mathcal{T}$ and that $\mathbf{\pi}^{\mathcal{T}} = \mathbf{\pi}^{\mathcal{T}} \mathbf{P}^{\mathcal{T}}$ if $\mathbf{\pi}^{\mathcal{T}}$ is the restriction of $\mathbf{\pi}$ to $\mathbb{S} \setminus \mathcal{T}$. Conversely, observe that if $\boldsymbol{\mu}^{\mathcal{T}}$ is a stationary probability for $\mathbf{P}^{\mathcal{T}}$ and $\boldsymbol{\mu}$ is the extention of $\boldsymbol{\mu}^{\mathcal{T}}$ that is 0 on \mathcal{T} , then $\boldsymbol{\mu}$ is a stationary probability for \mathbf{P} . Hence, $\mathbf{\pi}^{\mathcal{T}}$ is the unique stationary probability for $\mathbf{P}^{\mathcal{T}}$ and therefore, by (3.2.7) applied to $\mathbf{P}^{\mathcal{T}}$, the assertion is proved.

Another interesting application of the ideas here is contained in the following theorem.

Theorem 3.2.9 Let **P** and \mathbb{S} be as in Theorem 3.2.6. If *i* is recurrent and $i \in \Delta \subsetneq \mathbb{S}$, then $(\mathbf{I} - \mathbf{P})^{\Delta}$ is invertible and

$$\left(\left((\mathbf{I}-\mathbf{P})^{\Delta}\right)^{-1}\right)_{jk} = \mathbb{E}\left[\sum_{m=0}^{\zeta^{\Delta}-1} \mathbf{1}_{\{k\}}(X_m) \,\middle|\, X_0 = j\right] \quad for \ j, k \in \mathbb{S} \setminus \Delta, \qquad (3.2.10)$$

where $\zeta^{\Delta} \equiv \inf\{n \ge 0 : X_n \in \Delta\}$ is the time of the first visit by $\{X_n : n \ge 0\}$ to Δ . In particular,

$$\mathbb{E}\left[\sum_{m=0}^{\zeta^{\Delta}-1} \mathbf{1}_{\{j\}}(X_m) \,\middle|\, X_0 = j\right] = \frac{\det((\mathbf{I} - \mathbf{P})^{\Delta \cup \{j\}})}{\det((\mathbf{I} - \mathbf{P})^{\Delta})} \quad \text{for } j \in \mathbb{S} \setminus \Delta.$$
(3.2.11)

(See Exercise 4.2.8 for an interesting application of (3.2.11).)

Proof First observe that

$$\mathbb{P}(X_m = j_m \text{ for } 0 \le m \le n \& \zeta^{\Delta} > n \mid X_0 = j_0)$$
$$= \mathbf{1}_{\Delta \mathsf{C}}(j_0) \prod_{m=1}^n \mathbf{1}_{\Delta \mathsf{C}}(j_m) (\mathbf{P})_{j_{m-1}j_m}$$

and therefore that

$$\mathbb{P}(X_n = k \& \zeta^{\Delta} > n \mid X_0 = j) = (\mathbf{P}^{\Delta})_{jk}^n \quad \text{for } j, k \in \mathbb{S} \setminus \Delta.$$

From this it follows that

$$\mathbb{E}\left[\sum_{m=0}^{n\wedge\zeta^{\Delta}-1}\mathbf{1}_{\{k\}}(X_m)\,\Big|\,X_0=j\right] = \sum_{m=0}^{n-1} \left(\mathbf{P}^{\Delta}\right)_{jk}^m \quad \text{for all } n \ge 1. \tag{*}$$

Hence, since $\mathbb{E}[\rho_i | X_0 = j] \ge \mathbb{E}[\zeta^{\Delta} | X_0 = j]$ and, by Theorems 2.3.8 and 3.2.1, we know that $\mathbb{E}[\rho_i | X_0 = j] < \infty$, we have that

$$\sum_{k\in\mathbb{S}}\sum_{m=0}^{\infty} (\mathbf{P}^{\Delta})_{jk} = \sum_{k\in\mathbb{S}} \mathbb{E}\left[\sum_{m=0}^{\zeta^{\Delta}-1} \mathbf{1}_{\{k\}}(X_m) \mid X_0 = j\right] = \mathbb{E}\left[\zeta^{\Delta} \mid X_0 = j\right] < \infty$$

for all $j \in \mathbb{S}$. In particular, if $\mathbf{v} \in \mathbb{R}^{\mathbb{S}\setminus \Delta}$, then $\sum_{m=0}^{\infty} (\mathbf{P}^{\Delta})^m \mathbf{v}$ converges to a $\mathbf{w} \in \mathbb{R}^{\mathbb{S}}$, and it is easy to see that $(\mathbf{I} - \mathbf{P})^{\Delta}\mathbf{w} = \mathbf{v}$. This proves that $(\mathbf{I} - \mathbf{P})^{\Delta}$ is invertible and, after letting $n \to \infty$ in (*), that (3.2.10) holds. Furthermore, (3.2.11) is just (3.2.10) followed by an application of (3.2.4).

3.3 Wilson's Algorithm and Kirchhoff's Formula

A graph Γ is an ordered pair (V, \mathcal{E}) , where V is a finite set and \mathcal{E} is a collection of unordered pairs of elements from V. An element of V is called a *vertex* and an element of \mathcal{E} is called an *edge*. Without further notice, we will be considering graphs in which no vertex has an edge connecting it to itself and no pair of vertices has more than one edge connecting them. When an edge exists between vertices v and w, we will write $v \sim w$.

A directed path is an ordered n-tuple $(v_1, \ldots, v_n) \in V^n$ such that $v_m \sim v_{m-1}$ for each $2 \leq m \leq n$, and a graph is said to be *connected* if every pair of distinct vertices has a directed path connecting them. A directed path (v_1, \ldots, v_n) is said to be a *loop* if $v_1 = v_n$, and a graph is called a *tree* if it is connected and no directed path in it is a loop. Equivalently, a graph is a tree if every pair of distinct vertices has precisely one directed path running from one to the other.

Given a graph $\Gamma = (V, \mathcal{E})$, a *subgraph* is a graph $\Gamma' = (V', \mathcal{E}')$, where $V' \subseteq V$ and \mathcal{E}' is a subset of $\{\{v, w\} \in \mathcal{E} : v, w \in V'\}$. Finally, given a connected graph, a *spanning tree* is a subgraph that is a tree in which all the vertices of the graph are present. A challenging combinatorial problem is that of determining how many spanning trees a graph has. When a graph with N elements is *complete*, in the sense that there is an edge between every pair of distinct vertices, G. Cayley showed that there are N^{N-2} spanning trees, although the problem had been solved in the general case somewhat earlier by G. Kirchhoff.

A related problem is that of finding an algorithm for constructing the uniform measure on the set of all spanning trees. Such an algorithm based on Markov chains was discovered by D. Wilson. In this section, I will present Wilson's algorithm and from it derive Cayley's and Kirchhoff's results.

3.3.1 Spanning Trees and Wilson Runs

In order to explain Wilson's algorithm, I must first give a method for labeling the spanning tree in terms of what I will call *Wilson runs*.

A Wilson run is an *L*-tuple (P_1, \ldots, P_L) where:

- (1) for each $1 \le \ell \le L$, $P_{\ell} = (w_{1,\ell}, \dots, w_{K_{\ell},\ell})$ is loop free (i.e., $w_{k,\ell} \ne w_{k',\ell}$ if $k \ne k'$) directed path of length $K_{\ell} \ge 2$,
- (2) $V = \{w_{k,\ell} : 1 \le \ell \le L \& 1 \le k \le K_\ell\},\$
- (3) if $2 \le \ell \le L$, then

$$\{w_{k,\ell}: 1 \le k \le K_{\ell}\} \cap \{w_{k,j}: 1 \le j < \ell \& 1 \le k \le K_{j}\} = \{w_{K_{\ell},\ell}\}.$$

Each Wilson run determines a unique spanning tree. Namely, the spanning tree whose edges are $\{\{w_{k,\ell}, w_{k+1,\ell}\}: 1 \le \ell \le L \& 1 \le k < K_\ell\}$.

Although a Wilson run determines a unique spanning tree, there are many Wilson runs that determine the same spanning tree. For example, if Γ is a complete graph with three vertices v_1, v_2, v_3 and T is the spanning tree with edges

 $\{\{v_1, v_2\}, \{v_2, v_3\}\}\$, then both (v_1, v_2, v_3) and $((v_1, v_2), (v_3, v_2))$ are Wilson runs that determine *T*. Nonetheless, one can make the correspondence one-to-one by specifying an ordering (v_1, \ldots, v_N) of the vertices in *V* and insisting that the initial segment in the Wilson run start at v_1 and end at v_2 and that later segments always start at the first v_m that is not in any of the earlier segments. To be precise, given a spanning tree *T*, take P_1 to be the directed path in *T* that begins at v_1 and ends at v_2 . Because *T* is a tree, there is only one such path. If all the vertices in *V* appear in P_1 , stop. Otherwise, take P_2 to be the directed path in *T* that starts at the first v_m not in P_1 and ends when it first hits P_1 . Again, because *T* is a tree, P_2 is uniquely determined. More generally, if one has constructed (P_1, \ldots, P_ℓ) and there is a vertex that is not in³ $\bigcup_{j=1}^{\ell} P_j$, take $P_{\ell+1}$ to be the unique directed path in *T* that connects the first v_m that is not in $\bigcup_{j=1}^{\ell} P_j$ to $\bigcup_{j=1}^{\ell} P_j$. After $L \leq N - 1$ steps, one will have produced a Wilson run (P_1, \ldots, P_L) that determines *T*, and this Wilson run will be the only one that respects the ordering (v_1, \ldots, v_N) .

Finally, I must describe the procedure that will be used to erase loops from a directed path. For this purpose, suppose that $J \ge 2$ and that (X_1, \ldots, X_J) is a directed path in the graph with the property that $X_J \notin \{X_0, \ldots, X_{J-1}\}$. We will say that (w_1, \ldots, w_K) is the *loop erasure* of (X_0, \ldots, X_J) if it is obtained by sequentially erasing the loops from (X_1, \ldots, X_J) . To be precise, $w_1 = X_0$, $w_K = X_J$, and, for $2 \le k \le K$, $w_k = X_{\ell_{k-1}+1}$, where $\ell_1 < \cdots < \ell_{K-1}$ are determined by $\ell_1 = \max\{j > 0 : X_j = w_1\}$ and, for $2 \le k \le K$, $\ell_k = \max\{j > \ell_{k-1} : X_j = w_k\}$. Clearly, (w_1, \ldots, w_K) is a loop free directed path.

3.3.2 Wilson's Algorithm

Let $\Gamma = (V, \mathcal{E})$ be a connected graph with $N \ge 2$ vertices. For each $v \in V$, the *degree* of v is the number d_v of $w \in V$ such that $\{v, w\} \in \mathcal{E}$. Because Γ is connected, $d_v \ge 1$ for all $v \in V$. Next, let \mathcal{A} be the *adjacency matrix*. That is, $\mathcal{A}_{vw} = 1$ if $\{v, w\} \in \mathcal{E}$ and $\mathcal{A}_{vw} = 0$ otherwise. Because of the running assumptions that I made about our graphs, the diagonal entries of \mathcal{A} are 0 and $\sum_{w \in V} \mathcal{A}_{vw} = d_v = \sum_w \mathcal{A}_{wv}$ for each $v \in V$. Further, because Γ is connected, $\sum_{n=0}^N (\mathcal{A}^n)_{vw} \ge 1$ for all $(v, w) \in V^2$.

Now let \mathcal{D} be the diagonal matrix $\mathcal{D}_{vw} = \delta_{v,w} d_v$, and set $\mathbf{P} = \mathcal{D}^{-1}\mathcal{A}$. Then \mathbf{P} is a transition probability, and the paths $\{X_n : n \ge 0\}$ followed by a Markov process with transition probability \mathbf{P} will have the property that $X_n \sim X_{n-1}$ for all $n \ge 1$. In addition, $\sum_{n=0}^{N} (\mathbf{P}^n)_{vw} \ge \epsilon$ for all $(v, w) \in V^2$ and some $\epsilon > 0$. Hence, by Theorem 2.2.5, there is precisely one stationary probability $\boldsymbol{\pi}$ for \mathbf{P} , $(\boldsymbol{\pi})_v \ge \epsilon$ for all v, and therefore every v is recurrent.

Wilson's algorithm is the following. Choose and fix an ordering $(v_1, ..., v_N)$ of the vertices, and start the Markov process $\{X_n : n \ge 0\}$ with transition probability **P**

³Although the P_{ℓ} 's are tuples and not sets, the following union is the set consisting of the entries of the tuples over which it is taken.

at v_1 . Set $\Delta_1(\omega) = \{v_2\}$ and, assuming that (cf. Theorem 3.2.9) $\zeta^{\Delta_1}(\omega) < \infty$, and take $P_1(\omega)$ to be the loop erasure of $(X_0(\omega), \ldots, X_{\zeta^{\Delta_1}}(\omega))$. If $P_1(\omega)$ contains all the vertices in V, the algorithm terminates at time $\zeta^{\Delta_1}(\omega)$. Otherwise, take $\Delta_2(\omega) =$ $P_1(\omega)$, start the Markov chain at the first v_m that is not in $\Delta_2(\omega)$, and, assuming that $\zeta^{\Delta_2}(\omega) < \infty$, take $P_2(\omega)$ to be the loop erasure of $(X_0(\omega), \ldots, X_{\zeta^{\Delta_2}}(\omega))$. More generally, assuming that $\zeta^{\Delta_{\ell-1}}(\omega) < \infty$ and that $\Delta_{\ell}(\omega) \equiv \bigcup_{j=1}^{\ell-1} P_j(\omega) \neq V$, one starts the chain at the first $v_m \notin \Delta_{\ell}(\omega)$ and takes $P_{\ell}(\omega)$ to be the loop erasure of $(X_0(\omega), \ldots, X_{\zeta^{\Delta_{\ell}}}(\omega))$. Assuming that $\zeta^{\Delta_{\ell}} < \infty$ for all ℓ , one continues in this way until one generates a Wilson run $(P_1(\omega), \ldots, P_L(\omega))$ that respects the ordering (v_1, \ldots, v_N) . Of course, one has to worry about the fact that one or more of the $\zeta^{\Delta_{\ell}}$'s may be infinite. However, by Theorem 2.3.8, with probability 1, none of them will be. Thus, with probability 1, Wilson's algorithm will generate a Wilson run that respects the ordering (v_1, \ldots, v_N) .

The problem now is to compute the probability $\mathcal{P}(T)$ that Wilson's algorithm generates the Wilson run corresponding to a specific spanning tree T. To solve this problem, for a given $\emptyset \neq \Delta \subsetneq V$ and a loop free directed path $P = (w_1, \ldots, w_K) \in$ $(V \setminus \Delta)^{K-1} \times \Delta$, let $\mathcal{P}^{\Delta}(P)$ be the probability that, when Markov chain is started at $w_1, \zeta^{\Delta} < \infty$ and the loop erasure of $(X_0, \ldots, X_{\zeta^{\Delta}})$ is P. If the Wilson run determining T is (P_1, \ldots, P_L) , then

$$\mathcal{P}(T) = \prod_{\ell=1}^{L} \mathcal{P}^{\Delta_{\ell}}(P_{\ell}), \qquad (3.3.1)$$

where $\Delta_1 = \{v_2\}$ and $\Delta_\ell = \bigcup_{j=1}^{\ell-1} P_j$ for $2 \le \ell \le L$. Thus, what we need to learn is how to compute $\mathcal{P}^{\Delta}(P)$, and for this purpose I will use a strategy that was introduced by Greg Lawler.

Lemma 3.3.2 Let $\emptyset \neq \Delta \subsetneq V$ and a loop free directed path

$$P = (w_1, \dots, w_K) \in (V \setminus \Delta)^{K-1} \times \Delta$$

be given. Set $\Delta(1) = \Delta$ and $\Delta(k) = \Delta \cup \{w_1, \dots, w_{k-1}\}$ for $2 \le k \le K$. Then

$$\mathcal{P}^{\Delta}(P) = \prod_{k=1}^{K-1} \frac{\left(\left((\mathbf{I} - \mathbf{P})^{\Delta(k)}\right)^{-1}\right)_{w_k, w_k}}{d_{w_k}} = \frac{\det((\mathcal{D} - \mathcal{A})^{\Delta \cup P})}{\det((\mathcal{D} - \mathcal{A})^{\Delta})}.$$
(3.3.3)

Proof Let $B = \{\omega : \zeta^{\Delta}(\omega) < \infty \& P(\omega) = P\}$, and, for

$$\mathbf{m} = (m_1, \ldots, m_{K-1}) \in \mathbb{N}^{K-1},$$

let $B^{\mathbf{m}}$ be the set of $\omega \in B$ such that, for each $1 \leq k < K$, precisely m_k loops are erased at step k in the passage from $(X_0(\omega), \ldots, X_{\zeta^{\Delta}(\omega)}(\omega))$ to P. Clearly

$$\mathbb{P}(B \mid X_0 = w_1) = \sum_{\mathbf{m} \in \mathbb{N}^{K-1}} \mathbb{P}(B^{\mathbf{m}} \mid X_0 = w_1).$$

To compute $\mathbb{P}(B^{\mathbf{m}} | X_0 = w_1)$, let $\sigma_1(\omega)$ be the time of the $(m_1 + 1)$ th visit to w_1 and $\zeta_1(\omega)$ the time of the first visit to $\Delta(1)$ by the path $\{X_n(\omega) : n \ge 0\}$. Proceeding by induction, for $2 \le k < K$, take $\sigma_k(\omega) = \zeta_k(\omega) = \infty$ if $\sigma_{k-1}(\omega) \ge \zeta_{k-1}(\omega)$ and otherwise let $\sigma_k(\omega)$ be the time of the $(m_k + 1)$ th visit to w_k and $\zeta_k(\omega)$ the time of the first visit to $\Delta(k)$ by $\{X_n(\omega) : n > \sigma_{k-1}(\omega)\}$. Then

$$B^{\mathbf{m}} = \left\{ \zeta^{\Delta} < \infty \right\} \cap \{ \sigma_k < \zeta_k < \infty \& X_{\sigma_k+1} = w_{k+1} \text{ for } 1 \le k < K \}.$$

We will now use an induction procedure to compute $\mathbb{P}(B^{\mathbf{m}} | X_0 = w_1)$. For $1 \le k < K$, set

$$B_k^{\mathbf{m}} = \left\{ \zeta^{\Delta} < \infty \right\} \cap \{ \sigma_j < \zeta_j < \infty \& X_{\sigma_j+1} = w_{j+1} \text{ for } 1 \le j \le k \}.$$

Obviously, $B^{\mathbf{m}} = B_{K-1}^{\mathbf{m}}$. By the Markov property,

$$\mathbb{P}(B_1^{\mathbf{m}} \mid X_0 = w_1)$$

$$= \sum_{n=0}^{\infty} \mathbb{P}(\sigma_1 = n < \zeta^{\Delta(1)} \& X_{n+1} = w_2 \mid X_0 = w_1)$$

$$= \sum_{n=0}^{\infty} \frac{\mathbb{P}(\rho_{w_1}^{(m_1)} = n < \zeta^{\Delta(1)} \mid X_0 = w_1)}{d_{w_1}} = \frac{\mathbb{P}(\rho_{w_1}^{(m_1)} < \zeta^{\Delta(1)} \mid X_0 = w_1)}{d_{w_1}},$$

where, for any $m \in \mathbb{N}$ and $v \in V$, $\rho_v^{(m)}(\omega)$ is the time of the *m*th return to *v* by $\{X_n(\omega) : n \ge 0\}$. Similarly, for $2 \le k < K$,

$$\mathbb{P}(B_k^{\mathbf{m}} \mid X_0 = w_1)$$

$$= \mathbb{P}(B_{k-1}^{\mathbf{m}} \cap \{\sigma_k < \zeta_k \& X_{\sigma_k+1} = w_{k+1}\} \mid X_0 = w_1)$$

$$= \mathbb{P}(B_{k-1}^{\mathbf{m}}) \frac{\mathbb{P}(\rho_{w_k}^{(m_k)} < \zeta^{\Delta(k)} \mid X_0 = w_k)}{d_{w_k}}.$$

Since $\rho_v^{(m)} < \zeta^{\Delta(k)} \iff \sum_{n=0}^{\zeta^{\Delta(k)-1}} \mathbf{1}_{\{v\}}(X_n) > m$,

$$\sum_{m_k=0}^{\infty} \mathbb{P}\left(\rho_{w_k}^{(m_k)} < \zeta^{\Delta(k)} \mid X_0 = w_k\right) = \mathbb{E}\left[\sum_{n=0}^{\zeta^{\Delta(k)}-1} \mathbf{1}_{\{w_k\}}(X_n) \mid X_0 = w_k\right],\$$

which, by Theorem 3.2.9, is equal to $(((\mathbf{I} - \mathbf{P})^{\Delta(k)})^{-1})_{w_k w_k}$, and so

$$\sum_{\mathbf{m}\in\mathbb{N}^{K-1}}\mathbb{P}(B^{\mathbf{m}})=\prod_{k=1}^{K-1}\frac{(((\mathbf{I}-\mathbf{P})^{\Delta(k)})^{-1})_{w_kw_k}}{d_{w_k}}$$

which is the first equality in (3.3.3). To prove the second one, note that, because $\mathcal{D} - \mathcal{A} = \mathcal{D}(\mathbf{I} - \mathbf{P})$,

$$\frac{(((\mathbf{I}-\mathbf{P})^{\Delta(k)})^{-1})_{w_k w_k}}{d_{w_k}} = \left(\left((\mathcal{D}-\mathcal{A})^{\Delta(k)}\right)^{-1}\right)_{w_k w_k},$$

and apply (3.2.4).

By combining this lemma with (3.3.1), we arrive at the following statement of Wilson's theorem.

Theorem 3.3.4 (Wilson) For any spanning tree T and any vertex $v \in V$,

$$\mathcal{P}(T) = \frac{1}{\det((\mathcal{D} - \mathcal{A})^{\{v\}})}$$

In particular, $\mathcal{P}(T)$ is the same for all spanning trees T.

Proof Choose an ordering $(v_1, ..., v_N)$ with $v_2 = v$, let $(P_1, ..., P_L)$ be the Wilson run determined by T that respects this ordering, define Δ_ℓ as in (3.3.1) for $1 \le \ell \le L$, and set $\Delta_{L+1} = V$. Then, by Lemma 3.3.2 and (3.3.1),

$$\mathcal{P}(T) = \prod_{\ell=1}^{L} \frac{\det((\mathcal{D} - \mathcal{A})^{\Delta_{\ell+1}})}{\det((\mathcal{D} - \mathcal{A})^{\Delta_{\ell}})} = \frac{1}{\det((\mathcal{D} - \mathcal{A})^{\{\nu\}})}.$$

3.3.3 Kirchhoff's Matrix Tree Theorem

Let everything be as in the preceding subsection.

Obviously each row and column of $\mathcal{D} - \mathcal{A}$ sums to 0. In addition, because 0 is a simple eigenvalue for $\mathbf{I} - \mathbf{P}$ and therefore for $\mathcal{D} - \mathcal{A}$, Theorem 3.2.3 applies. Further, because the columns of $\mathcal{D} - \mathcal{A}$ sum to 0, $N^{-1}\mathbf{1}$ is a left eigenvector, and therefore Theorem 3.2.3 says that

$$\det((\mathcal{D}-\mathcal{A})^{\{v\}}) = \frac{\Pi_{\mathcal{D}-\mathcal{A}}}{N}$$

Hence, by Theorem 3.3.4,

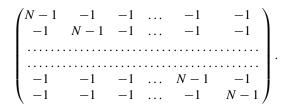
$$\mathcal{P}(T) = \frac{1}{\det((\mathcal{D} - \mathcal{A})^{\{v\}})} = \frac{N}{\Pi_{\mathcal{D} - \mathcal{A}}} \quad \text{for any spanning tree } T.$$
(3.3.5)

Since, with probability 1, Wilson's algorithm produces a spanning tree, it follows that, for all $v \in V$,

the number of spanning trees = det
$$((\mathcal{D} - \mathcal{A})^{\{v\}}) = \frac{\Pi_{\mathcal{D} - \mathcal{A}}}{N},$$
 (3.3.6)

which is one statement of Kirchhoff's renowned matrix tree theorem.

To derive Cayley's formula from (3.3.6), assume that the graph is complete, and let (v_1, \ldots, v_N) be an ordering of its vertices. Then $(\mathcal{D} - \mathcal{A})^{\{v_N\}}$ is equal to the $(N-1) \times (N-1)$ -matrix



By subtracting the last column from each of the columns 2 through N - 2, one sees that the determinant of $(\mathcal{D} - \mathcal{A})^{\{v_N\}}$ is the same as that of

$\binom{N-1}{-1}$	$0 \\ N$	0 0	 	0 0	$\begin{pmatrix} -1 \\ -1 \end{pmatrix}$	
1			 	 	· · · · · · · · · · · · · · · · · · ·	,
$\begin{pmatrix} -1\\ -1 \end{pmatrix}$	-N	-N	•••• •••	-N	$\begin{pmatrix} -1 \\ N-1 \end{pmatrix}$	

and by adding the sum of the columns 1 through N - 2 times $(N - 1)^{-1}$ to the last column, one sees that this has the same determinant as

$\binom{N-1}{-1}$	$0 \\ N$	0 0	 0 0	$\begin{pmatrix} 0 \\ 0 \end{pmatrix}$
-1	0	0	 N	0
−1	-N	-N	 -N	$\frac{N}{N-1}$

Hence there are N^{N-2} spanning trees in a complete graph with N vertices.

3.4 Exercises

Exercise 3.4.1 In Sect. 1.2.4 we worked quite hard to prove that the nearest neighbor, symmetric random walk on \mathbb{Z}^3 is transient, and one might hope that the criteria provided in Sect. 3.1.2 would allow us to avoid working so hard. However, even if one knows which function u ought to be plugged into Theorem 3.1.6 or 3.1.8, the computation to show that it works is rather delicate.

Let **P** be the transition probability for the simple, symmetric random walk in \mathbb{Z}^3 . That is,

$$(\mathbf{P})_{\mathbf{k}\boldsymbol{\ell}} = \begin{cases} \frac{1}{6} & \text{if } \sum_{i=1}^{3} |(\mathbf{k})_i - (\boldsymbol{\ell})_i| = 1\\ 0 & \text{otherwise.} \end{cases}$$

Show that if $\alpha \ge 1$ is sufficiently large and

$$u(\mathbf{k}) = \left(\alpha^2 + \sum_{i=1}^3 (\mathbf{k})_i^2\right)^{-\frac{1}{2}} \quad \text{for } \mathbf{k} \in \mathbb{Z}^3,$$

then $(\mathbf{Pu})_{\mathbf{k}} \le u(\mathbf{k}) \le u(\mathbf{0})$ when **u** is the column vector determined by the function *u*. By Theorem 3.1.8, this proves that **0** is transient.

What follows are some hints.

(a) Let $\mathbf{k} \in \mathbb{Z}^3$ be given, and set

$$M = 1 + \alpha^2 + \sum_{i=1}^{3} (\mathbf{k})_i^2$$
 and $x_i = \frac{(\mathbf{k})_i}{M}$ for $1 \le i \le 3$.

Show that $(\mathbf{Pu})_{\mathbf{k}} \leq u(\mathbf{k})$ if and only if

$$\left(1-\frac{1}{M}\right)^{-\frac{1}{2}} \ge \frac{1}{3} \sum_{i=1}^{3} \frac{(1+2x_i)^{\frac{1}{2}} + (1-2x_i)^{\frac{1}{2}}}{2(1-4x_i^2)^{\frac{1}{2}}}$$

(b) Show that $(1 - \frac{1}{M})^{-\frac{1}{2}} \ge 1 + \frac{1}{2M}$ and that

$$\frac{(1+\xi)^{\frac{1}{2}} + (1-\xi)^{\frac{1}{2}}}{2} \le 1 - \frac{\xi^2}{8} \quad \text{for } |\xi| < 1,$$

and conclude that $(\mathbf{Pu})_{\mathbf{k}} \leq u(\mathbf{k})$ if

$$1 + \frac{1}{2M} \ge \frac{1}{3} \sum_{i=1}^{3} \frac{1}{(1 - 4x_i^2)^{\frac{1}{2}}} - \frac{|\mathbf{x}|^2}{6}$$

where $|\mathbf{x}| = \sqrt{\sum_{i=1}^{3} x_i^2}$ is the Euclidean length of $\mathbf{x} = (x_1, x_2, x_3)$. (c) Show that there is a constant $C < \infty$ such that, as long as $\alpha \ge 1$,

$$\frac{1}{3}\sum_{i=1}^{3}\frac{1}{(1-4x_i^2)^{\frac{1}{2}}} \le 1 + \frac{2|\mathbf{x}|^2}{3} + C|\mathbf{x}|^4,$$

and put this together with the preceding to conclude that we can take any $\alpha \ge 1$ with $\alpha^2 + 1 \ge 2C$.

An analogous computation shows that, for each $d \ge 3$ and all sufficiently large $\alpha > 0$, the function $(\alpha^2 + \sum_{i=1}^{d} (\mathbf{k})_i^2)^{-\frac{d-2}{2}}$ can be used to prove that the nearest neighbor, symmetric random walk is transient in \mathbb{Z}^d . The reason why one does not get a contradiction when d = 2 and the walk is recurrent is that the non-constant, non-negative functions which satisfies $\mathbf{P}u \le u$ when d = 2 are of the form $\log(\alpha^2 + (\mathbf{k})_1^2 + (\mathbf{k})_2^2)$ and therefore do not achieve their maximum value at **0**.

Exercise 3.4.2 Let $\mathbf{P} = ((p_{i,j}))_{1 \le i,j \le 3}$ be a transition probability on $\{1, 2, 3\}$, and set

$$D(i, j, k) = p_{ji}(1 - p_{kk}) + p_{jk}p_{ki}.$$

Show that

$$det((\mathbf{I} - \mathbf{P})^{\{1\}}) = D(1, 2, 3), \qquad det((\mathbf{I} - \mathbf{P})^{\{2\}}) = D(2, 3, 1),$$

and
$$det((\mathbf{I} - \mathbf{P})^{\{3\}}) = D(3, 1, 2),$$

and conclude that **P** has a unique stationary distribution π if and only if

$$\Pi \equiv D(1, 2, 3) + D(2, 3, 1) + D(3, 1, 2) > 0,$$

in which case,

$$(\pi)_1 = \frac{D(1,2,3)}{\Pi}, \qquad (\pi)_2 = \frac{D(2,3,1)}{\Pi}, \quad \text{and} \quad (\pi)_3 = \frac{D(3,1,2)}{\Pi}.$$

Exercise 3.4.3 Given a transition probability **P** on a finite state space S, show that, for each $\Delta \subsetneq S$,

$$0 \le \det((\mathbf{I} - \mathbf{P})^{\Delta}) \le \det((\mathbf{I} - \mathbf{P})^{\Delta \cup \{i\}}) \quad \text{for all } i \in \mathbb{S} \setminus \Delta.$$

Conclude that if $\Delta \subseteq \Delta' \subseteq \mathbb{S}$, then

$$0 \le \det((\mathbf{I} - \mathbf{P})^{\Delta}) \le \det((\mathbf{I} - \mathbf{P})^{\Delta'}).$$

Hint: Show that for each $\epsilon > 0$ there is a \mathbf{P}^{ϵ} that satisfies Doeblin's condition and whose entries differ from those of \mathbf{P} by less that ϵ . In this way, reduce the problem to transition probabilities that satisfies Doeblin's condition.

Exercise 3.4.4 Let **P** the transition probability introduced in Sect. 3.2.2, and show that the unique stationary probability π for **P** is given by $(\pi)_v = \frac{d_v}{D}$, where $D = \sum_{w \in V} d_w$.

Chapter 4 More About the Ergodic Properties of Markov Chains

In Chap. 2 all of our considerations centered around one form or another of the Doeblin condition which says that there is a state which can be reached from any other state at a uniformly fast rate. Although there are lots of chains on an infinite state space which satisfy his condition, most do not. In fact, many chains on an infinite state space will not even admit a stationary probability distribution because the availability of infinitely many states means that there is enough room for the chain to "get lost and disappear." There are two ways in which this can happen. Namely, the chain can disappear because, like the a nearest neighbor, non-symmetric random walk in \mathbb{Z} (cf. (1.1.13)) or even the symmetric one in \mathbb{Z}^3 (cf. Sect. 1.2.4), it may have no recurrent states and, as a consequence, will spend a finite amount of time in any given state. A more subtle way for the chain to disappear is for it to be recurrent but not sufficiently recurrent for there to exist a stationary distribution. Such an example is the symmetric, nearest neighbor random walk in \mathbb{Z} which is recurrent, but just barely so. In particular, although this random walk returns infinitely often to the place where is begins, it does so at too sparse a set of times. More precisely, by (1.2.7) and (1.2.13), if **P** is the transition probability matrix for the symmetric, nearest neighbor random walk on \mathbb{Z} , then

$$\left(\mathbf{P}^{2n}\right)_{ij} \le \left(\mathbf{P}^{2n}\right)_{ii} = \mathbb{P}(X_{2n} = 0) \le A(1)n^{-\frac{1}{2}} \longrightarrow 0,$$

and so, if μ were a probability vector which was stationary for **P**, then, by Lebesgue's dominated convergence theorem, Theorem 7.1.11, we would have the contradiction

$$(\boldsymbol{\mu})_j = \lim_{n \to \infty} \sum_{i \in \mathbb{Z}} (\boldsymbol{\mu})_i (\mathbf{P}^{2n})_{ij} = 0 \text{ for all } j \in \mathbb{Z}.$$

In this chapter we will see that ergodic properties can exist in the absence of Doeblin's condition. However, as we will see, what survives does so in a weaker form. See part (c) of Exercise 4.2.5 for a comparison of the results here with the earlier ones.

D.W. Stroock, *An Introduction to Markov Processes*, Graduate Texts in Mathematics 230, 73 DOI 10.1007/978-3-642-40523-5_4, © Springer-Verlag Berlin Heidelberg 2014

4.1 Ergodic Theory Without Doeblin

In this section we will see to what extent the results obtained in Chap. 2 for Markov chains which satisfy Doeblin's condition can be reproduced without Doeblin's condition. The progression which we will adopt here runs in the opposite direction to that in Chap. 2. That is, here we will start with the most general but weakest form of the result and afterwards will see what can be done to refine it.

4.1.1 Convergence of Matrices

Because we will be looking at power series involving matrices that may have infinitely many entries, it will be important for us to be precise about what is the class of matrices with which we are dealing and in what sense our series are converging. For our purposes, the most natural class will be of matrices \mathbf{M} for which

$$\|\mathbf{M}\|_{\mathbf{u},\mathbf{v}} \equiv \sup_{i \in \mathbb{S}} \sum_{j \in \mathbb{S}} \left| (\mathbf{M})_{ij} \right|$$
(4.1.1)

is finite, and the set of all such matrices will be denoted by $M_{u,v}(\mathbb{S})$. An easy calculation shows that $M_{u,v}(\mathbb{S})$ is a vector space over \mathbb{R} and that $\|\cdot\|_{u,v}$ is a good norm on $M_{u,v}(\mathbb{S})$. That is,

$$\|\mathbf{M}\|_{\mathbf{u},\mathbf{v}} = 0 \quad \text{if and only if} \quad \mathbf{M} = \mathbf{0},$$
$$\|\alpha\mathbf{M}\|_{\mathbf{u},\mathbf{v}} = |\alpha|\|\mathbf{M}\|_{\mathbf{u},\mathbf{v}} \quad \text{for } \alpha \in \mathbb{R},$$

and

$$\|\mathbf{M} + \mathbf{M}'\|_{u,v} \le \|\mathbf{M}\|_{u,v} + \|\mathbf{M}'\|_{u,v}.$$

Slightly less obvious is the fact that

$$(\mathbf{M}, \mathbf{M}') \in M_{\mathbf{u}, \mathbf{v}}(\mathbb{S})^2 \implies \mathbf{M}\mathbf{M}' \text{ exists and} \\ \|\mathbf{M}\mathbf{M}'\|_{\mathbf{u}, \mathbf{v}} \leq \|\mathbf{M}\|_{\mathbf{u}, \mathbf{v}} \|\mathbf{M}'\|_{\mathbf{u}, \mathbf{v}}.$$
(4.1.2)

To see this, observe first that, since

$$\sum_{k} \left| (\mathbf{M})_{ik} \right| \left| \left(\mathbf{M}' \right)_{kj} \right| \leq \left(\sum_{k} \left| (\mathbf{M})_{ik} \right| \right) \sup_{k} \left| \left(\mathbf{M}' \right)_{kj} \right| < \infty,$$

the sum in

$$\left(\mathbf{M}\mathbf{M}'\right)_{ij} = \sum_{k} (\mathbf{M})_{ik} \left(\mathbf{M}'\right)_{kj}$$

is absolutely convergent. In addition, for each *i*,

$$\begin{split} \sum_{j} |(\mathbf{M}\mathbf{M}')_{ij}| &\leq \sum_{j} \sum_{k} |(\mathbf{M})_{ik}| |(\mathbf{M}')_{kj}| \\ &= \sum_{k} |(\mathbf{M})_{ik}| \left(\sum_{j} |(\mathbf{M}')_{kj}|\right) \leq \left(\sum_{k} |(\mathbf{M})_{ik}|\right) \|\mathbf{M}'\|_{\mathbf{u},\mathbf{v}}, \end{split}$$

and so the inequality in (4.1.2) follows.

We will show next that the metric that $\|\cdot\|_{u,v}$ determines on $M_{u,v}(\mathbb{S})$ is complete. In order to do this, it will be useful to know that if $\{\mathbf{M}_n\}_0^\infty \subseteq M_{u,v}(\mathbb{S})$, then

$$(\mathbf{M})_{ij} = \lim_{n \to \infty} (\mathbf{M}_n)_{ij} \quad \text{for each } (i, j)$$
$$\implies \|\mathbf{M}\|_{\mathbf{u}, \mathbf{v}} \le \lim_{n \to \infty} \|\mathbf{M}_n\|_{\mathbf{u}, \mathbf{v}}.$$
(4.1.3)

By Fatou's lemma, Theorem 7.1.10,

$$\sum_{j \in \mathbb{S}} \left| (\mathbf{M})_{ij} \right| \le \lim_{n \to \infty} \sum_{j \in \mathbb{S}} \left| (\mathbf{M}_n)_{ij} \right| \quad \text{for each } i \in \mathbb{S},$$

and so (4.1.3) is proved.

Knowing (4.1.3), the proof of completeness goes as follows. Assume that $\{\mathbf{M}_n\}_0^{\infty} \subseteq M_{u,v}(\mathbb{S})$ is $\|\cdot\|_{u,v}$ -Cauchy convergent:

$$\lim_{m\to\infty}\sup_{n>m}\|\mathbf{M}_n-\mathbf{M}_m\|_{\mathbf{u},\mathbf{v}}=0.$$

Obviously, for each $(i, j) \in \mathbb{S}^2$, $\{(\mathbf{M}_n)_{ij}\}_0^\infty$ is Cauchy convergent as a sequence in \mathbb{R} . Hence, there is a matrix \mathbf{M} such that, for each (i, j), $(\mathbf{M})_{ij} = \lim_{n \to \infty} (\mathbf{M}_n)_{ij}$. Furthermore, by (4.1.3),

$$\|\mathbf{M} - \mathbf{M}_m\|_{\mathbf{u},\mathbf{v}} \leq \lim_{n \to \infty} \|\mathbf{M}_n - \mathbf{M}_m\|_{\mathbf{u},\mathbf{v}} \leq \sup_{n > m} \|\mathbf{M}_n - \mathbf{M}_m\|_{\mathbf{u},\mathbf{v}}.$$

Hence, $\|\mathbf{M} - \mathbf{M}_m\|_{u,v} \longrightarrow 0.$

4.1.2 Abel Convergence

As I said in the introduction, we will begin with the weakest form of the convergence results at which we are aiming. That is, rather than attempting to prove the convergence of $(\mathbf{P}^n)_{ij}$ or even the Césaro means $\frac{1}{n} \sum_{m=0}^{n-1} (\mathbf{P}^m)_{ij}$ as $n \to \infty$, we will begin by studying the Abel sums $(1-s) \sum_{m=0}^{\infty} s^m (\mathbf{P}^m)_{ij}$ as $s \nearrow 1$.

I will say that a bounded sequence $\{x_n\}_0^\infty \subseteq \mathbb{R}$ is *Abel convergent* to *x* if

$$\lim_{s \nearrow 1} (1-s) \sum_{n=1}^{\infty} s^n x_n = x.$$

It should be clear that Abel convergence is weaker than (i.e., is implied by) ordinary convergence.¹ Indeed, if $x_n \rightarrow x$, then, since $(1 - s) \sum_{0}^{\infty} s^n = 1$, for any N:

$$\begin{vmatrix} x - (1-s) \sum_{n=0}^{\infty} s^n x_n \end{vmatrix} = (1-s) \begin{vmatrix} \sum_{n=0}^{\infty} s^n (x-x_n) \end{vmatrix}$$
$$\leq (1-s) \sum_{n=0}^{\infty} s^n |x-x_n|$$
$$\leq N(1-s) \sup_{n \in \mathbb{N}} |x-x_n| + \sup_{n > N} |x-x_n|$$

Hence,

$$\overline{\lim_{s \neq 1}} \left| x - (1 - s) \sum_{1}^{\infty} s^n x_n \right| \le \lim_{N \to \infty} \sup_{n \ge N} |x - x_n| = 0$$

if $x_n \longrightarrow x$. On the other hand, although $\{(-1)^n\}_1^\infty$ fails to converge to anything,

$$(1-s)\sum_{n=0}^{\infty}s^{n}(-1)^{n} = (1-s)\frac{1}{1+s} \longrightarrow 0 \text{ as } s \nearrow 1.$$

That is, Abel convergence does not in general imply ordinary convergence.

With the preceding in mind, we set

$$\mathbf{R}(s) = (1-s) \sum_{n=0}^{\infty} s^n \mathbf{P}^n \quad \text{for } s \in [0, 1).$$
 (4.1.4)

However, before accepting this definition, it is necessary to check that the above series converges. To this end, first note that, because \mathbf{P}^n is a transition probability matrix for each $n \ge 0$, $\|\mathbf{P}^n\|_{u,v} = 1$. Hence, for $0 \le m < n$,

$$\left\| (1-s) \sum_{\ell=0}^{n} s^{\ell} \mathbf{P}^{\ell} - (1-s) \sum_{\ell=0}^{m} s^{\ell} \mathbf{P}^{\ell} \right\|_{\mathbf{u},\mathbf{v}} \le (1-s) \sum_{\ell=m+1}^{n} s^{\ell} \left\| \mathbf{P}^{\ell} \right\|_{\mathbf{u},\mathbf{v}} \le s^{m},$$

and so, by the completeness proved above, the series in (4.1.4) converges with respect to the $\|\cdot\|_{u,v}$ -norm.

¹In Exercise 4.2.1 below, it is shown that Abel convergence is also weaker than Césaro convergence.

4.1 Ergodic Theory Without Doeblin

Our goal here is to prove that

$$\lim_{s \nearrow 1} (\mathbf{R}(s))_{ij} = \pi_{ij} \equiv \begin{cases} (\mathbb{E}[\rho_j | X_0 = j])^{-1} & \text{if } i = j \\ \mathbb{P}(\rho_j < \infty | X_0 = i)\pi_{jj} & \text{if } i \neq j, \end{cases}$$
(4.1.5)

and the key to our doing so lies in the renewal equation

$$\left(\mathbf{P}^{n}\right)_{ij} = \sum_{m=1}^{n} f(m)_{ij} \left(\mathbf{P}^{n-m}\right)_{jj} \quad \text{for } n \ge 1,$$
where $f(m)_{ij} \equiv \mathbb{P}(\rho_{j} = m \mid X_{0} = i),$

$$(4.1.6)$$

which is an elementary application of (2.1.11):

$$(\mathbf{P}^{n})_{ij} = \sum_{m=1}^{n} \mathbb{P}(X_{n} = j \& \rho_{j} = m \mid X_{0} = i)$$

=
$$\sum_{m=1}^{n} \mathbb{P}(X_{n-m} = j \mid X_{0} = j) \mathbb{P}(\rho_{j} = m \mid X_{0} = i).$$

Next, for $s \in [0, 1)$, set

$$\hat{f}(s)_{ij} \equiv \sum_{m=1}^{\infty} s^m f(m)_{ij} = \mathbb{E}\left[s^{\rho_j} \,\middle| \, X_0 = i\right],$$

and, starting from (4.1.6), conclude that

$$\begin{aligned} \left(\mathbf{R}(s) \right)_{ij} &= (1-s)\delta_{i,j} + (1-s)\sum_{n=1}^{\infty} s^n \left(\sum_{m=1}^n f(m)_{ij} \left(\mathbf{P}^{n-m} \right)_{jj} \right) \\ &= (1-s)\delta_{i,j} + (1-s)\sum_{m=1}^{\infty} s^m f(m)_{ij} \left(\sum_{n=m}^{\infty} s^{n-m} \left(\mathbf{P}^{n-m} \right)_{jj} \right) \\ &= (1-s)\delta_{i,j} + \hat{f}(s)_{ij} \left(\mathbf{R}(s) \right)_{jj}. \end{aligned}$$

That is,

$$\left(\mathbf{R}(s)\right)_{ij} = (1-s)\delta_{i,j} + \hat{f}(s)_{ij}\left(\mathbf{R}(s)\right)_{jj}$$
 for $s \in [0, 1).$ (4.1.7)

Given (4.1.7), (4.1.5) is easy. Namely,

$$\left(\mathbf{R}(s)\right)_{jj} = \frac{1-s}{1-\hat{f}(s)_{jj}}.$$

Hence, if j is transient, and therefore $\hat{f}(1)_{jj} < 1$,

$$\pi_{jj} = \lim_{s \neq 1} (\mathbf{R}(s))_{jj} = 0 = \frac{1}{\mathbb{E}[\rho_j | X_0 = j]}.$$

On the other hand, if j is recurrent, then, since

$$\frac{1-s^m}{1-s} = \sum_{\ell=0}^{m-1} s^\ell \nearrow m,$$

the monotone convergence theorem says that

$$\frac{1 - \hat{f}(s)_{jj}}{1 - s} = \sum_{m=1}^{\infty} \frac{1 - s^m}{1 - s} f(m)_{jj} \nearrow \sum_{m=1}^{\infty} mf(m)_{jj} = \mathbb{E}[\rho_j \mid X_0 = j]$$

as $s \nearrow 1$. At the same time, when $i \neq j$,

$$\left(\mathbf{R}(s)\right)_{ij} = \hat{f}(s)_{ij} \left(\mathbf{R}(s)\right)_{jj} \nearrow \mathbb{P}(\rho_j < \infty \mid X_0 = i)\pi_{jj}.$$

4.1.3 Structure of Stationary Distributions

We will say that a probability vector μ is **P**-*stationary* and will write $\mu \in \text{Stat}(\mathbf{P})$ if $\mu = \mu \mathbf{P}$. Obviously, if μ is stationary, then $\mu = \mu \mathbf{R}(s)$ for each $s \in [0, 1)$. Hence, by (4.1.5) and Lebesgue's dominated convergence theorem,

$$\boldsymbol{\mu} \in \operatorname{Stat}(\mathbf{P}) \implies (\boldsymbol{\mu})_j = \sum_i (\boldsymbol{\mu})_i \left(\mathbf{R}(s) \right)_{ij} \longrightarrow \sum_i (\boldsymbol{\mu})_i \pi_{ij} \text{ as } s \nearrow 1$$

If *j* is transient, then $\pi_{ij} = 0$ for all *i* and therefore $(\boldsymbol{\mu})_j = 0$. If *j* is recurrent, then either *i* is transient, and, by Theorem 3.1.5, $(\boldsymbol{\mu})_i = 0$, or *i* is recurrent, in which case, by Theorem 3.1.2, either $\pi_{ij} = 0$ or $i \leftrightarrow j$ and $\pi_{ij} = \pi_{jj}$. Hence, in either case, we have that

$$\boldsymbol{\mu} \in \operatorname{Stat}(\mathbf{P}) \implies (\boldsymbol{\mu})_j = \left(\sum_{i \leftrightarrow j} (\boldsymbol{\mu})_i\right) \pi_{jj}.$$
 (4.1.8)

We show next that

$$\pi_{jj} > 0 \quad \text{and} \quad C = [j] = \{i : i \leftrightarrow j\}$$

$$\implies \quad \pi^C \in \text{Stat}(\mathbf{P}) \quad \text{when} \left(\pi^C\right)_i \equiv \mathbf{1}_C(i)\pi_{ii}.$$
(4.1.9)

To do this, first note that $\pi_{jj} > 0$ only if j is recurrent. Thus, all $i \in C$ are recurrent and, for each $i \in C$ and $s \in (0, 1)$, $(\mathbf{R}(s))_{ik} > 0 \iff k \in C$. In particular, $(\pi^C)_i =$

 $\lim_{s \neq 1} (\mathbf{R}(s))_{ii}$ for all *i*, and therefore, by Fatou's lemma,

$$\sum_{i} \left(\boldsymbol{\pi}^{C} \right)_{i} \leq \underline{\lim}_{s \neq 1} \sum_{i} \left(\mathbf{R}(s) \right)_{ji} = 1.$$

Similarly, for any *i*,

$$\left(\boldsymbol{\pi}^{C}\mathbf{P}\right)_{i} = \sum_{k \in C} \pi_{kk}(\mathbf{P})_{ki} \leq \underline{\lim}_{s \neq 1} \sum_{k \in C} \left(\mathbf{R}(s)\right)_{jk}(\mathbf{P})_{ki} = \left(\boldsymbol{\pi}^{C}\right)_{i},$$

since

$$\sum_{k \in C} (\mathbf{R}(s))_{jk} (\mathbf{P})_{ki} = \frac{(\mathbf{R}(s) - (1 - s)\mathbf{I})_{ji}}{s} \longrightarrow \pi_{ji} = (\boldsymbol{\pi}^C)_i \quad \text{as } s \nearrow 1.$$

But if strict inequality were to hold for some i, then, by Fubini's theorem, Theorem 7.1.15, we would have the contradiction

$$\sum_{k} (\boldsymbol{\pi}^{C})_{k} = \sum_{k} (\boldsymbol{\pi}^{C})_{k} \left(\sum_{i} (\mathbf{P})_{ki} \right) = \sum_{i} \left(\sum_{k} (\boldsymbol{\pi}^{C})_{k} (\mathbf{P})_{ki} \right) < \sum_{i} (\boldsymbol{\pi}^{C})_{i}.$$

Hence, we now know that $\pi^C = \pi^C \mathbf{P}$. Finally, to prove that $\pi^C \in \text{Stat}(\mathbf{P})$, we still have to check that $\sum_i (\pi^C)_i = 1$. However, we have already checked that $\pi^C = \pi^C \mathbf{P}$, and so we know that $\pi^C = \pi^C \mathbf{R}(s)$. Therefore, since, as we already showed, $\sum_i (\pi^C)_i \leq 1$, Lebesgue's dominated convergence theorem justifies

$$0 < \pi_{jj} = \sum_{i} (\boldsymbol{\pi}^{C})_{i} (\mathbf{R}(s))_{ij} \longrightarrow \left(\sum_{i} (\boldsymbol{\pi}^{C})_{i}\right) \pi_{jj},$$

which is possible only if $\sum_{i} (\boldsymbol{\pi}^{C})_{i} = 1$.

Before summarizing the preceding as a theorem, we need to recall that a subset *A* of a linear space is said to be a *convex set* if $(1 - \theta)a + \theta a' \in A$ for all $a, a' \in A$ and $\theta \in [0, 1]$ and that $b \in A$ is an *extreme point* of the convex set *A* if $b = (1 - \theta)a + \theta a'$ for some $\theta \in (0, 1)$ and $a, a' \in A$ implies that b = a = a'. In addition, we need the notion of positive recurrence. Namely, we say that $j \in S$ is *positive recurrent* if $\mathbb{E}[\rho_j | X_0 = j] < \infty$. Obviously, only recurrent states can be positive recurrent. On the other hand, (1.1.13) together with (1.1.15) show that there can exist *null recurrent* states, those which are recurrent but not positive recurrent.

Theorem 4.1.10 Stat(**P**) is a convex subset of $\mathbb{R}^{\mathbb{S}}$. Moreover, Stat(**P**) $\neq \emptyset$ if and only if there is at least one positive recurrent state $j \in \mathbb{S}$. In fact, for any $\mu \in$ Stat(**P**), (4.1.8) holds, and μ is an extreme point in Stat(**P**) if and only if there is a communicating class *C* of positive recurrent states for which (cf. (4.1.9)) $\mu = \pi^{C}$. In particular, $(\mu)_{j} = 0$ for any transient state *j* and, for any recurrent state *j*, either $(\mu)_{i}$ is strictly positive or it is 0 simultaneously for all states *i*'s which communicate with *j*.

Proof The only statements not already covered are the characterization of the extreme points of $\text{Stat}(\mathbf{P})$ and the final assertion in the case when j is recurrent.

In view of (4.1.8), the final assertion when *j* is recurrent comes down to showing that if *j* is positive recurrent and $i \leftrightarrow j$, then *i* is positive recurrent. To this end, suppose that *j* is positive recurrent, set C = [j], and let $i \in C$ be given. Then $\pi^{C} \mathbf{P}^{n} = \pi^{C}$ for all $n \geq 0$, and therefore, by choosing *n* so that $(\mathbf{P}^{n})_{ji} > 0$, we see that $(\pi^{C})_{i} \geq (\pi^{C})_{j}(\mathbf{P}^{n})_{ji} > 0$.

To handle the characterization of extreme points, first suppose that $\mu \neq \pi^C$ for any communicating class *C* of positive recurrent states. Then, by (4.1.8), there must exist non-communicating, positive recurrent states *j* and *j'* for which $(\mu)_j > 0 < (\mu)_{j'}$. But, again by (4.1.8), this means that we can write $\mu = \theta \pi^C + (1 - \theta) \nu$, where C = [j], $\theta = \sum_{i \in C} (\mu)_i \in (0, 1)$, and $(\nu)_i$ equals 0 or $(1 - \theta)^{-1}(\mu)_i$ depending on whether *i* is or is not in *C*. Clearly $\nu \in \text{Stat}(\mathbf{P})$ and, because $\nu_{j'} > 0 = (\pi^C)_{j'}$, $\nu \neq \pi^C$. Hence, μ cannot be extreme. Equivalently, every extreme $\mu = \pi^C$ for some communicating class *C* of positive recurrent states.

Conversely, given a communicating class *C* of positive recurrent states, suppose that $\pi^{C} = (1 - \theta)\mu + \theta\nu$ for some $\theta \in (0, 1)$ and pair $(\mu, \nu) \in \text{Stat}(\mathbf{P})^{2}$. Then $(\mu)_{i} = 0$ for all $i \notin C$, and so, by (4.1.8), we see first that $\mu = \pi^{C}$ and then, as a consequence, that $\nu = \pi^{C}$ as well.

4.1.4 A Digression About Moments of Return Times²

The last part of Theorem 4.1.10 shows that *positive recurrence is a communicating class property*. In particular, if a chain is irreducible, we are justified in saying it is positive recurrent if any one of its states is. See Exercise 4.2.3 below for a criterion which guarantees positive recurrence.

I have often wondered whether the fact that positive recurrence is a communicating class property admits a direct proof, one that does not involve stationary probabilities. More generally, I wanted to know whether finiteness of other moments of return times is also a communicating class property. Affirmative answers to both these questions were given to me by Daniel Jerison, whose ideas provide the foundation on which everything in this subsection rests. To simplify the notation, I will use \mathbb{P}_i and \mathbb{E}_i here to denote conditional probability and expectation given $X_0 = i$.

Let *i* and *j* be a pair of distinct, recurrent, communicating states, and assume that $\mathbb{E}_i[\rho_i^p] < \infty$ for some $p \in (0, \infty)$. In order to show that $\mathbb{E}_j[\rho_j^p] < \infty$, set $A_m = \{\rho_i^{(m)} < \rho_j < \rho_i^{(m+1)}\}$ for $m \ge 1$. Then, by part (d) of Exercise 6.6.2 below,

²This section may be safely skipped since the material here is not used elsewhere.

4.1 Ergodic Theory Without Doeblin

$$\begin{split} \mathbb{E}_{j}[\rho_{j}^{p},A_{m}] &\leq 2^{(p-1)^{+}}\mathbb{E}_{j}[(\rho_{j}-\rho_{i}^{(m)})^{p},A_{m}] + 2^{(p-1)^{+}}\mathbb{E}_{j}[(\rho_{i}^{(m)})^{p},A_{m}] \\ &\leq 2^{(p-1)^{+}}\mathbb{E}_{j}[(\rho_{j}-\rho_{i}^{(m)})^{p},A_{m}] \\ &+ 2^{(p-1)^{+}}\mathbb{E}_{j}\left[\left(\sum_{\ell=1}^{m}(\rho_{i}^{(\ell)}-\rho_{i}^{(\ell-1)})\right)^{p},A_{m}\right] \\ &\leq 2^{(p-1)^{+}}\mathbb{E}_{j}[(\rho_{j}-\rho_{i}^{(m)})^{p},A_{m}] + 4^{(p-1)^{+}}\mathbb{E}_{j}[\rho_{i}^{p},A_{m}] \\ &+ (4(m-1))^{(p-1)^{+}}\sum_{2\leq\ell\leq m}\mathbb{E}_{j}[(\rho_{i}^{(\ell)}-\rho_{i}^{(\ell-1)})^{p},A_{m}]. \end{split}$$

By the Markov property,

$$\mathbb{E}_{j}\left[\left(\rho_{j}-\rho_{i}^{(m)}\right)^{p},A_{m}\right]=\mathbb{P}_{j}\left(\rho_{j}>\rho_{i}^{(m)}\right)\mathbb{E}_{i}\left[\rho_{j}^{p},\rho_{j}<\rho_{i}\right]\leq\mathbb{P}_{j}\left(\rho_{j}>\rho_{i}^{(m)}\right)\mathbb{E}_{i}\left[\rho_{i}^{p}\right]$$

and

$$\mathbb{E}_{j}[\rho_{i}^{p}, A_{m}] = \mathbb{E}_{j}[\rho_{i}^{p}, \rho_{j} > \rho_{i}]\mathbb{P}_{i}(A_{m-1}) \leq \mathbb{P}_{i}(A_{m-1})\mathbb{E}_{j}[\rho_{i}^{p}].$$

Similarly, for $2 \le \ell \le m$,

$$\begin{split} & \mathbb{E}_{j} \Big[\big(\rho_{i}^{(\ell)} - \rho_{i}^{(\ell-1)} \big)^{p}, A_{m} \Big] \\ &= \mathbb{E}_{j} \Big[\big(\rho_{i}^{(\ell)} - \rho_{i}^{(\ell-1)} \big)^{p}, \rho_{j} > \rho_{i}^{(\ell)} \Big] \mathbb{P}_{i}(A_{m-\ell}) \\ &= \mathbb{P}_{j} \big(\rho_{j} > \rho_{i}^{(\ell-1)} \big) \mathbb{E}_{i} \Big[\rho_{i}^{p}, \rho_{j} > \rho_{i} \Big] \mathbb{P}_{j}(A_{m-\ell}) = \mathbb{E}_{i} \Big[\rho_{i}^{p}, \rho_{j} > \rho_{i} \Big] \mathbb{P}_{j}(A_{m-1}) \\ &\leq \mathbb{P}_{j}(A_{m-1}) \mathbb{E}_{i} \Big[\rho_{i}^{p} \Big]. \end{split}$$

Hence,

$$\mathbb{E}_{j}[\rho_{j}^{p}, A_{m}] \leq 2^{(p-1)^{+}} \mathbb{P}_{j}(\rho_{j} > \rho_{i}^{(m)}) \mathbb{E}_{i}[\rho_{i}^{p}] + 4^{(p-1)^{+}} \mathbb{P}_{i}(A_{m-1}) \mathbb{E}_{j}[\rho_{i}^{p}] + 4^{(p-1)^{+}}(m-1)^{p \vee 1} \mathbb{P}_{j}(A_{m-1}) \mathbb{E}_{i}[\rho_{i}^{p}].$$

After summing over $m \in \mathbb{Z}^+$, we see that $\mathbb{E}_j[\rho_j^p, \rho_j > \rho_i]$ is dominated by

$$\sum_{m=1}^{\infty} \left(2^{(p-1)^{+}} \mathbb{P}_{j} \left(\rho_{j} > \rho_{i}^{(m)} \right) + 4^{(p-1)^{+}} (m-1)^{p \vee 1} \mathbb{P}_{j} (A_{m-1}) \right) \mathbb{E}_{i} \left[\rho_{i}^{p} \right]$$
$$+ 4^{(p-1)^{+}} \sum_{m=1}^{\infty} \mathbb{P}_{j} (A_{m-1}) \mathbb{E}_{j} \left[\rho_{i}^{p} \right].$$

Obviously $\sum_{m=1}^{\infty} \mathbb{P}_j(A_{m-1}) = 1$. In the proof of Theorem 3.1.2 we showed that

$$\alpha \equiv \mathbb{P}_i(\rho_j > \rho_i) \in (0, 1)$$
 and $\mathbb{P}_i(\rho_j > \rho_i^{(m)}) = \mathbb{P}_i(\rho_j > \rho_i)^m = \alpha^m$.

In addition, for $m \ge 1$.

$$\mathbb{P}_j(\rho_j > \rho_i^{(m)}) = \mathbb{P}_j(\rho_j > \rho_i)\mathbb{P}_i(\rho_j > \rho_i^{(m-1)}) \le \alpha^{m-1}.$$

Thus, since $\mathbb{P}_j(A_{m-1}) \leq \mathbb{P}_j(\rho_j > \rho_i^{(m-1)})$, we now know that there is a $C < \infty$ such that

$$\mathbb{E}_{j}\left[\rho_{j}^{p},\rho_{j}>\rho_{i}\right] \leq C\left(\mathbb{E}_{i}\left[\rho_{i}^{p}\right]+\mathbb{E}_{j}\left[\rho_{i}^{p}\right]\right)$$

and therefore that

$$\mathbb{E}_j[\rho_j^p] \le C \mathbb{E}_i[\rho_i^p] + (1+C) \mathbb{E}_j[\rho_i^p].$$

Finally, observe that

$$\mathbb{E}_i \left[\rho_i^p \right] \ge \mathbb{E}_i \left[(\rho_i - \rho_j)^p, \rho_j < \rho_i \right] = \mathbb{E}_j \left[\rho_i^p \right] \mathbb{P}_i (\rho_j < \rho_i),$$

and, as shown in the proof of Theorem 3.1.2, $\mathbb{P}_i(\rho_j < \rho_i) \in (0, 1)$. Hence, we have now shown the $\mathbb{E}_j[\rho_j^p] < \infty$.

4.1.5 A Small Improvement

Returning to our program, the next step will be the replacement of Abel convergence by Césaro convergence. That is, we will show that (4.1.5) can be replaced by

$$\lim_{n \to \infty} (\mathbf{A}_n)_{ij} = \pi_{ij}, \quad \text{where } \mathbf{A}_n = \frac{1}{n} \sum_{m=0}^{n-1} \mathbf{P}^m.$$
(4.1.11)

As is shown in Exercise 4.2.1, Césaro convergence does not, in general, follow from Abel convergence. In fact, general results which say when Abel convergence implies Césaro convergence can be quite delicate and, because the original one was proved by a man named Tauber, they are known as Tauberian theorems. Fortunately, the Tauberian theorem required here is quite straight-forward.

A key role in our proof will be played by the following easy estimate:

$$\{a_m\}_0^\infty \subseteq [0, 1] \& A_n = \frac{1}{n} \sum_{0}^{n-1} a_\ell$$

$$\implies |A_n - A_{n-m}| \le \frac{m}{n} \quad \text{for } 0 \le m < n.$$
(4.1.12)

The proof is:

$$A_n - A_{n-m} = \frac{1}{n} \sum_{\ell=n-m}^{n-1} a_\ell - \frac{m}{n(n-m)} \sum_{\ell=0}^{n-m-1} a_\ell \begin{cases} \geq -\frac{m}{n} \\ \leq \frac{m}{n} \end{cases}$$

Lemma 4.1.13 For all (i, j), $\overline{\lim}_{n\to\infty} (\mathbf{A}_n)_{ij} \leq e\pi_{jj}$. In addition, for any j and any subsequence $\{n_{\ell} : \ell \geq 0\} \subseteq \mathbb{N}$,

$$\lim_{\ell \to \infty} (\mathbf{A}_{n_{\ell}})_{jj} = \alpha \quad \Longrightarrow \quad \lim_{\ell \to \infty} (\mathbf{A}_{n_{\ell}})_{ij} = \mathbb{P}(\rho_j < \infty \mid X_0 = i)\alpha \quad \text{for all } i.$$

Proof To prove the first part, observe that

$$(\mathbf{A}_{n})_{ij} \leq \frac{1}{n} \left(1 - \frac{1}{n}\right)^{-n} \sum_{m=0}^{n-1} \left(1 - \frac{1}{n}\right)^{m} \left(\mathbf{P}^{m}\right)_{ij} \leq \left(1 - \frac{1}{n}\right)^{-n} \left(\mathbf{R}\left(1 - \frac{1}{n}\right)\right)_{ij},$$

which, together with (1.2.10) and (4.1.5), shows that $\overline{\lim}_{n\to\infty}(\mathbf{A}_n)_{ij} \le e\pi_{ij} \le e\pi_{jj}$. To handle the second part, use (4.1.6) to arrive at

$$(\mathbf{A}_n)_{ij} = \sum_{m=1}^{n-1} f(m)_{ij} \left(1 - \frac{m}{n}\right) (\mathbf{A}_{n-m})_{jj} \quad \text{for } i \neq j.$$

Hence,

$$\left| (\mathbf{A}_n)_{ij} - \mathbb{P}(\rho_j < n \mid X_0 = i) \alpha \right| \leq \sum_{m=1}^{n-1} f(m)_{ij} \left(\frac{m}{n} + \left| (\mathbf{A}_{n-m})_{jj} - \alpha \right| \right)$$
$$\leq 2 \sum_{m=1}^{n-1} \frac{m}{n} f(m)_{ij} + \left| (\mathbf{A}_n)_{jj} - \alpha \right|,$$

where, in the second inequality, I have used (4.1.12) plus $\sum_{m=1}^{\infty} f(m)_{ij} \leq 1$. Finally, by Lebesgue's dominated convergence theorem, $\sum_{0}^{n-1} \frac{m}{n} f(m)_{ij}$ tends to 0 as $n \to \infty$, and therefore, by applying the above with $n = n_{\ell}$ and letting $\ell \to \infty$, we get the desired conclusion.

We can now complete the proof of (4.1.11). If $\pi_{jj} = 0$, then the first part of Lemma 4.1.13 guarantees that $\lim_{n\to\infty} (\mathbf{A}_n)_{ij} = 0 = \pi_{ij}$ for all *i*. Thus, assume that $\pi_{jj} > 0$. In this case Theorem 4.1.10 says that *j* must be positive recurrent and $\boldsymbol{\pi}^C \in \text{Stat}(\mathbf{P})$ when C = [j]. In particular, $\pi_{jj} = \sum_{i \in C} (\boldsymbol{\pi}^C)_i (\mathbf{A}_n)_{ij}$. At the same time, if $\alpha^+ = \overline{\lim}_{n\to\infty} (\mathbf{A}_n)_{jj}$ and the subsequence $\{n_\ell : \ell \ge 0\}$ is chosen so that $(\mathbf{A}_{n_\ell})_{jj} \longrightarrow \alpha^+$, then, by the second part of Lemma 4.1.13 and Corollary 3.1.4,

$$i \in C \implies \lim_{\ell \to \infty} (\mathbf{A}_{n_\ell})_{ij} = \alpha^+.$$

Hence, after putting these two remarks together, we arrive at

$$\pi_{jj} = \lim_{\ell \to \infty} \sum_{i \in C} (\boldsymbol{\pi}^C)_i (\mathbf{A}_{n_\ell})_{ij} = \alpha^+ \sum_{i \in C} (\boldsymbol{\pi}^C)_i = \alpha^+.$$

Similarly, if $\alpha^- = \underline{\lim}_{n \to \infty} (\mathbf{A}_n)_{jj}$, we can show that $\alpha^- = \pi_{jj}$, and so we now know that $\pi_{jj} > 0 \implies \underline{\lim}_{n \to \infty} (\mathbf{A}_n)_{jj} = \pi_{jj}$, which, after another application of the second part of Lemma 4.1.13, means that we have proved (4.1.11).

4.1.6 The Mean Ergodic Theorem Again

Just as we were able to use Theorem 2.2.5 in Sect. 2.3.1 to prove Theorem 2.3.4, so here we can use (4.1.11) to prove the following version of the mean ergodic theorem.

Theorem 4.1.14 Let C a communicating class of positive recurrent states. If $\mathbb{P}(X_0 \in C) = 1$, then

$$\lim_{n\to\infty} \mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m)-\pi_{jj}\right)^2\right]=0.$$

(See Exercises 4.2.10 and 4.2.12 below for a more refined statement.)

Proof Since $\mathbb{P}(X_m \in C \text{ for all } m \in \mathbb{N}) = 1$, without loss in generality we may and will assume that *C* is the whole state space. In keeping with this assumption, we will set $\pi = \pi^C$.

Next note that if $\mu_i = \mathbb{P}(X_0 = i)$, then

$$\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m) - \pi_{jj}\right)^2\right]$$
$$= \sum_{i \in \mathbb{S}} \mu_i \mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m) - \pi_{jj}\right)^2 \middle| X_0 = i\right],$$

and so it suffices to prove that

$$\lim_{n \to \infty} \mathbb{E}\left[\left(\frac{1}{n} \sum_{m=0}^{n-1} \mathbf{1}_{\{j\}}(X_m) - \pi_{jj}\right)^2 \, \middle| \, X_0 = i\right] = 0 \quad \text{for each } i \in \mathbb{S}.$$

Further, because $\pi_{ii} > 0$ for all $i \in \mathbb{S}$ and

$$\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m)-\pi_{jj}\right)^2 \middle| X_0=i\right]$$

$$\leq \frac{1}{\pi_{ii}}\sum_{k\in\mathbb{S}}\pi_{kk}\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m)-\pi_{jj}\right)^2 \middle| X_0=k\right],$$

it is enough to prove the result when π is the initial distribution of the Markov chain. Hence, from now on, we will be making this assumption along with C = S.

Now let **f** the column vector whose *i*th component is $\mathbf{1}_{\{j\}}(i) - \pi_{jj}$. Then, just as in the proof of Theorem 2.3.4,

$$\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m)-\pi_{jj}\right)^2\right] \leq \frac{2}{n^2}\sum_{k=0}^{n-1}(n-k)\mathbb{E}\left[(\mathbf{f})_{X_k}(\mathbf{A}_{n-k}\mathbf{f})_{X_k}\right].$$

Since $\pi \in \text{Stat}(\mathbf{P})$,

$$\mathbb{E}\big[(\mathbf{f})_{X_k}(\mathbf{A}_{n-k}\mathbf{f})_{X_k}\big] = \boldsymbol{\pi}(f\mathbf{A}_{n-k}\mathbf{f}),$$

and therefore the preceding becomes

$$\mathbb{E}\left[\left(\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m)-\pi_{jj}\right)^2\right] \leq \frac{2}{n^2}\sum_{m=1}^n m\boldsymbol{\pi}(f\mathbf{A}_m\mathbf{f}),$$

where $(f \mathbf{A}_m \mathbf{f})_i \equiv (\mathbf{f})_i (\mathbf{A}_m \mathbf{f})_i$.

Finally, because $\pi_{ij} = \pi_{jj}$ for all $i \in S$, (4.1.11) and Lebesgue's dominated convergence theorem say that, for each $\epsilon > 0$, there exists an $N_{\epsilon} \in \mathbb{Z}^+$ such that

$$\left|\boldsymbol{\pi}(f\mathbf{A}_{n}\mathbf{f})\right| \leq \sum_{i} (\boldsymbol{\pi})_{i} \left| (\mathbf{A}_{n})_{ij} - \boldsymbol{\pi}_{jj} \right| < \epsilon \quad \text{for all } n \geq N_{\epsilon}.$$

Hence, we find that

$$\overline{\lim_{n \to \infty}} \mathbb{E}\left[\left(\frac{1}{n} \sum_{m=0}^{n-1} \mathbf{1}_{\{j\}}(X_m) - \pi_{jj}\right)^2\right]$$

$$\leq \overline{\lim_{n \to \infty}} \frac{2}{n^2} \sum_{m=1}^{N_{\epsilon}} m \left| \boldsymbol{\pi}(f \mathbf{A}_m \mathbf{f}) \right| + \overline{\lim_{n \to \infty}} \frac{2\epsilon}{n^2} \sum_{m=N_{\epsilon}+1}^n m \leq \epsilon.$$

4.1.7 A Refinement in the Aperiodic Case

Our goal in this subsection is to prove that (cf. (4.1.5))

If *j* is transient or aperiodic, then $\lim_{n \to \infty} (\mathbf{P}^n)_{ij} = \pi_{ij}$ for all $i \in \mathbb{S}$. (4.1.15)

Of course, the case when j is transient requires very little effort. Indeed, by (2.3.7), we have that

j transient
$$\implies \sum_{n=0}^{\infty} (\mathbf{P}^n)_{ij} \le \mathbb{E}[T_j | X_0 = j] < \infty,$$

and therefore that $\lim_{n\to\infty} (\mathbf{P}^n)_{ij} = 0$. At the same time, because

$$\mathbb{P}(\rho_j = \infty | X_0 = j) > 0$$
 when j is transient,

 $\pi_{ij} \le \pi_{jj} = 0$. Thus, from now on, we will concentrate on the case when j is recurrent and aperiodic.

Our first step is the observation that (cf. Sect. 2.3.2)

if *j* is aperiodic, then there exists an $N \in \mathbb{Z}^+$ such that

$$\max_{1 \le m \le n} \mathbb{P}(\rho_j^{(m)} = n \mid X_0 = j) > 0 \quad \text{for all } n \ge N.$$
(4.1.16)

To check this, use (3.1.15) to produce an $N \in \mathbb{Z}^+$ such that $(\mathbf{P}^n)_{jj} > 0$ for all $n \ge N$. Then, since

$$\left(\mathbf{P}^{n}\right)_{jj} = \sum_{m=1}^{n} \mathbb{P}\left(\rho_{j}^{(m)} = n \mid X_{0} = j\right) \quad \text{for } n \ge 1,$$

(4.1.16) is clear.

The second, and key, step is contained in the following lemma.

Lemma 4.1.17 Assume that *j* is aperiodic and recurrent, and set $\alpha_j^- = \underline{\lim}_{n \to \infty} (\mathbf{P}^n)_{jj}$ and $\alpha_j^+ = \overline{\lim}_{n \to \infty} (\mathbf{P}^n)_{jj}$. Then there exist subsequences $\{n_\ell^- : \ell \ge 1\}$ and $\{n_\ell^+ : \ell \ge 1\}$ such that

$$\alpha_j^{\pm} = \lim_{\ell \to \infty} \left(\mathbf{P}^{n_{\ell}^{\pm} - r} \right)_{jj} \quad \text{for all } r \ge 0.$$

Proof Choose a subsequence $\{n_{\ell} : \ell \geq 1\}$ so that $(\mathbf{P}^{n_{\ell}})_{jj} \longrightarrow \alpha_{j}^{+}$, and, using (4.1.16), choose $N \geq 1$ so that $\max_{1 \leq m \leq n} \mathbb{P}(\rho_{j}^{(m)} = n | X_{0} = j) > 0$ for all $n \geq N$. Given $r \geq N$, choose $1 \leq m \leq r$ so that $\delta \equiv \mathbb{P}(\rho_{j}^{(m)} = r | X_{0} = j) > 0$. Now for any $M \in \mathbb{Z}^{+}$, observe that, when $n_{\ell} \geq M + r$, $(\mathbf{P}^{n_{\ell}})_{jj}$ is equal to

$$\begin{aligned} &\mathbb{P}(X_{n_{\ell}} = j \& \rho_{j}^{(m)} = r \mid X_{0} = j) + \mathbb{P}(X_{n_{\ell}} = j \& \rho_{j}^{(m)} \neq r \mid X_{0} = j) \\ &= \delta(\mathbf{P}^{n_{\ell} - r})_{jj} + \mathbb{P}(X_{n_{\ell}} = j \& n_{\ell} - M \ge \rho_{j}^{(m)} \neq r \mid X_{0} = j) \\ &+ \mathbb{P}(X_{n_{\ell}} = j \& \rho_{j}^{(m)} > n_{\ell} - M \mid X_{0} = j). \end{aligned}$$

Furthermore,

$$\mathbb{P}\left(X_{n_{\ell}}=j \& n_{\ell}-M \ge \rho_{j}^{(m)} \neq r \mid X_{0}=j\right)$$
$$= \sum_{\substack{k=1\\k \neq r}}^{n_{\ell}-M} \mathbb{P}\left(\rho_{j}^{(m)}=k \mid X_{0}=j\right) \left(\mathbf{P}^{n_{\ell}-k}\right)_{jj} \le (1-\delta) \sup_{n \ge M} \left(\mathbf{P}^{n}\right)_{jj},$$

while

$$\mathbb{P}(X_{n_{\ell}} = j \& \rho_j^{(m)} > n_{\ell} - M \mid X_0 = j) \le \mathbb{P}(\rho_j^{(m)} > n_{\ell} - M \mid X_0 = j).$$

Hence, since j is recurrent and therefore $\mathbb{P}(\rho_j^{(m)} < \infty | X_0 = j) = 1$, we get

$$\alpha_j^+ \le \delta \lim_{\ell \to \infty} \left(\mathbf{P}^{n_\ell - r} \right)_{jj} + (1 - \delta) \sup_{n \ge M} \left(\mathbf{P}^n \right)_{jj}$$

after letting $\ell \to \infty$. Since this is true for all $M \ge 1$, it leads to

$$\alpha_j^+ \leq \delta \lim_{\ell \to \infty} \left(\mathbf{P}^{n_\ell - r} \right)_{jj} + (1 - \delta) \alpha_j^+,$$

which implies $\underline{\lim}_{\ell\to\infty}(\mathbf{P}^{n_{\ell}-r})_{jj} \ge \alpha_{jj}^+$. But obviously $\overline{\lim}_{\ell\to\infty}(\mathbf{P}^{n_{\ell}-r})_{jj} \le \alpha_j^+$, and so we have shown that $\lim_{\ell\to\infty}(\mathbf{P}^{n_{\ell}-r})_{jj} = \alpha_j^+$ for all $r \ge N$. Now choose L so that $n_L \ge N$, take $n_{\ell}^+ = n_{\ell+L} - N$, and conclude that $\lim_{\ell\to\infty}(\mathbf{P}^{n_{\ell}^+-r})_{jj} = \alpha_j^+$ for all $r \ge 0$.

The construction of $\{n_{\ell}^{-} : \ell \ge 1\}$ is essentially the same and is left as an exercise.

Lemma 4.1.18 If *j* is aperiodic and recurrent, then $\overline{\lim}_{n\to\infty}(\mathbf{P}^n)_{jj} \leq \pi_{jj}$. Furthermore, if the subsequences $\{n_{\ell}^{\pm} : \ell \geq 1\}$ are the ones described in Lemma 4.1.17, then $\lim_{\ell\to\infty}(\mathbf{P}^{n_{\ell}^{\pm}})_{ij} = \alpha_j^{\pm}$ for any *i* with $i \leftrightarrow j$.

Proof To prove the second assertion, simply note that, by Lemma 4.1.17 and Lebesgue's dominated convergence theorem,

$$\left(\mathbf{P}^{n_{\ell}^{\pm}}\right)_{ij} = \sum_{r=1}^{n_{\ell}^{\pm}} \mathbb{P}(\rho_{j} = r | X_{0} = i) \left(\mathbf{P}^{n_{\ell}^{\pm} - r}\right)_{jj} \longrightarrow \mathbb{P}(\rho_{j} < \infty | X_{0} = i) \alpha_{j}^{\pm}.$$

Turning to the first assertion, we again use the result in Lemma 4.1.17 to obtain

$$\alpha_{j}^{+} \sum_{r=1}^{N} \mathbb{P}(\rho_{j} \ge r | X_{0} = j) = \lim_{\ell \to \infty} \sum_{r=1}^{N} \mathbb{P}(\rho_{j} \ge r | X_{0} = j) \left(\mathbf{P}^{n_{\ell}^{+} - r} \right)_{jj}$$

for all $N \ge 1$. Thus, if we show that

$$\sum_{r=1}^{N} \mathbb{P}(\rho_j \ge r \mid X_0 = j) \left(\mathbf{P}^{n-r} \right)_{jj} \le 1 \quad \text{for } n \ge N \ge 1,$$
 (*)

then we will know that

$$\alpha_{j}^{+}\mathbb{E}[\rho_{j} | X_{0} = j] = \alpha_{j}^{+} \sum_{r=1}^{\infty} \mathbb{P}(\rho_{j} \ge r | X_{0} = j) \le 1,$$

which is equivalent to the first assertion.

To prove (*), note that, for any $n \ge 1$,

$$\left(\mathbf{P}^{n}\right)_{jj} = \sum_{r=1}^{n} \mathbb{P}(\rho_{j} = r | X_{0} = j) \left(\mathbf{P}^{n-r}\right)_{jj}$$
$$= \sum_{r=1}^{n} \mathbb{P}(\rho_{j} \ge r | X_{0} = j) \left(\mathbf{P}^{n-r}\right)_{jj}$$
$$- \sum_{r=1}^{n} \mathbb{P}(\rho_{j} \ge r + 1 | X_{0} = j) \left(\mathbf{P}^{n-r}\right)_{jj}$$
$$= \sum_{r=1}^{n} \mathbb{P}(\rho_{j} \ge r | X_{0} = j) \left(\mathbf{P}^{n-r}\right)_{jj}$$
$$- \sum_{r=2}^{n+1} \mathbb{P}(\rho_{j} \ge r | X_{0} = j) \left(\mathbf{P}^{n+1-r}\right)_{jj},$$

and so, since $\mathbb{P}(\rho_j \ge 1 | X_0 = j) = 1$,

$$\sum_{r=1}^{n+1} \mathbb{P}(\rho_j \ge r \mid X_0 = j) \left(\mathbf{P}^{n+1-r} \right)_{jj} = \sum_{r=1}^n \mathbb{P}(\rho_j \ge r \mid X_0 = j) \left(\mathbf{P}^{n-r} \right)_{jj}$$

for all $n \ge 1$. But $\sum_{r=1}^{n} \mathbb{P}(\rho_j \ge r | X_0 = j) (\mathbf{P}^{n-r})_{jj} = 1$ when n = 1, and so we have now proved that

$$\sum_{r=1}^{N} \mathbb{P}(\rho_j \ge r \mid X_0 = j) (\mathbf{P}^{n-r})_{jj} \le \sum_{r=1}^{n} \mathbb{P}(\rho_j \ge r \mid X_0 = j) (\mathbf{P}^{n-r})_{jj} = 1$$

for all $n \ge N \ge 1$.

Remark Here is a more conceptual way to prove that

$$\sum_{r=1}^{n} \mathbb{P}(\rho_{j} \ge r \mid X_{0} = j) (\mathbf{P}^{n-r})_{jj} = 1.$$

Take $\rho_j^{(0)} = 0$ and, for $m \ge 1$, $\rho_j^{(m)}$ to be the time of the *m*th return to *j*. In addition, set $T_j^{(n-1)} = \sum_{\ell=0}^{n-1} \mathbf{1}_{\{j\}}(X_\ell)$. Then

$$X_0 = j \implies \{T_j^{(n-1)} = m+1\} = \{\rho_j^{(m)} < n \le \rho_j^{(m+1)}\}.$$

Hence,

$$\sum_{r=1}^{n} \mathbb{P}(\rho_{j} \ge r \mid X_{0} = j) (\mathbf{P}^{n-r})_{jj}$$

$$= \sum_{r=0}^{n-1} \mathbb{P}(\rho_{j} \ge n - r \mid X_{0} = j) (\mathbf{P}^{r})_{jj}$$

$$= \sum_{r=0}^{n-1} \sum_{m=0}^{r} \mathbb{P}(\rho_{j} \ge n - r \mid X_{0} = j) \mathbb{P}(\rho_{j}^{(m)} = r \mid X_{0} = j)$$

$$= \sum_{r=0}^{n-1} \sum_{m=0}^{r} \mathbb{P}(\rho_{j}^{(m+1)} \ge n \And \rho_{j}^{(m)} = r \mid X_{0} = j)$$

$$= \sum_{m=0}^{n-1} \sum_{r=m}^{n-1} \mathbb{P}(\rho_{j}^{(m+1)} \ge n \And \rho_{j}^{(m)} = r \mid X_{0} = j)$$

$$= \sum_{m=0}^{n-1} \mathbb{P}(\rho_{j}^{(m)} < n \le \rho_{j}^{(m+1)} \mid X_{0} = j)$$

$$= \sum_{m=0}^{n-1} \mathbb{P}(T_{j}^{(n-1)} = m + 1) = \mathbb{P}(T_{n}^{(n-1)} \le n \mid X_{0} = j) = 1.$$

We can now complete the proof of (4.1.15) when *j* is recurrent and aperiodic. By the first part of Lemma 4.1.18, we know that $\lim_{n\to\infty} (\mathbf{P}^n)_{jj} = 0$ if $\pi_{jj} = 0$. Thus, if $\pi_{jj} = 0$, then, by Lebesgue's dominated convergence theorem, for any *i*,

$$\lim_{n \to \infty} \left(\mathbf{P}^n \right)_{ij} = \lim_{n \to \infty} \sum_{r=1}^n \mathbb{P}(\rho_j = r | X_0 = i) \left(\mathbf{P}^{n-r} \right)_{jj} = 0 = \pi_{ij}.$$

In order to handle the case when *j* is positive recurrent, set C = [j], and take π^{C} accordingly. Then, $\pi^{C} \in \text{Stat}(\mathbf{P})$. In particular, by the last part of Lemma 4.1.18 and Lebesgue's dominated convergence theorem,

$$\pi_{jj} = \sum_{i \in C} (\boldsymbol{\pi}^C)_i (\mathbf{P}^{n_\ell^{\pm}})_{ij} \longrightarrow \alpha_j^{\pm} \sum_{i \in C} (\boldsymbol{\pi}^C)_i = \alpha_j^{\pm},$$

and so $(\mathbf{P}^n)_{jj} \longrightarrow \pi_{jj}$. Finally, if $i \neq j$, then, again by the Lebesgue's theorem,

$$\left(\mathbf{P}^{n}\right)_{ij} = \sum_{r=1}^{n} \mathbb{P}(\rho_{j} = r | X_{0} = i) \left(\mathbf{P}^{n-r}\right)_{jj} \longrightarrow \mathbb{P}(\rho_{j} < \infty | X_{0} = i) \pi_{jj} = \pi_{ij}.$$

4.1.8 Periodic Structure

The preceding result allows us to give a finer analysis even in the case when the period is not 1. Consider a Markov chain with transition probability matrix **P** which is irreducible and recurrent on \mathbb{S} , and assume that its period is $d \ge 2$. The basic result in this subsection is that there exists a partition of \mathbb{S} into subsets \mathbb{S}_r , $0 \le r < d$, with the properties that

- (1) $(\mathbf{P}^{md+r})_{ik} > 0 \implies r(k) r(j) = r \mod d$,
- (2) $r(k) r(j) = r \mod d \implies (\mathbf{P}^{md+r})_{jk} > 0 < (\mathbf{P}^{md-r})_{kj}$ for all sufficiently large $m \ge 1$,
- (3) for each $0 \le r < d$, the restriction of \mathbf{P}^d to \mathbb{S}_r is an aperiodic, recurrent, irreducible transition probability matrix,

where I have used r(j) to denote the $0 \le r < d$ for which $j \in S_r$.

To prove that this decomposition exists, we begin by noting that, for any $0 \le r < d$,

$$\exists m \ge 0 \quad \left(\mathbf{P}^{md+r}\right)_{ij} > 0 \quad \Longrightarrow \quad \exists n \ge 1 \quad \left(\mathbf{P}^{nd-r}\right)_{ji} > 0. \tag{(*)}$$

Indeed, by irreducibility and the Euclidean algorithm, we know that there exists an $m' \ge 0$ and $0 \le r' < d$ such that $(\mathbf{P}^{m'd+r'})_{ji} > 0$. Furthermore, by (3.1.14), $(\mathbf{P}^{(m'+m'')d+r'})_{ji} \ge (\mathbf{P}^{m'd+r'})_{ji} (\mathbf{P}^{m''d})_{ii} > 0$ for all sufficiently large m'''s, and so we may and will assume that $m' \ge 1$. But then $(\mathbf{P}^{(m+m')d+(r+r')})_{ii} > 0$, and so d|(r+r'), which, because $0 \le r, r' < d$, means that r' = 0 if r = 0 and that r' = d - r if $r \ge 1$.

Starting from (*) it is easy to see that, for each pair $(i, j) \in \mathbb{S}^2$, there is a unique $0 \le r < d$ such that $(\mathbf{P}^{md+r})_{ij} > 0$ for some $m \ge 0$. Namely, suppose that $(\mathbf{P}^{md+r})_{ij} > 0 < (\mathbf{P}^{m'd+r'})_{ij}$ for some $m, m' \in \mathbb{N}$ and $0 \le r, r' < d$. Then, by (*), there exists an $n \ge 1$ such that $(\mathbf{P}^{nd-r})_{ji} > 0$ and therefore $(\mathbf{P}^{(m'+n)d+(r'-r)})_{ii} > 0$. Since this means that d|(r'-r), we have proved that r = r'.

Now, let i_0 be a fixed reference point in \mathbb{S} , and, for each $0 \le r < d$, define \mathbb{S}_r to be the set of j such that there exists an $m \ge 0$ for which $(\mathbf{P}^{md+r})_{i_0j} > 0$. By the preceding, we know that the \mathbb{S}_r 's are mutually disjoint. In addition, by irreducibility and the Euclidean algorithm, $\mathbb{S} = \bigcup_{r=0}^{d-1} \mathbb{S}_r$. Turning to the proof of property (1), use (*) to choose $n \ge 0$ and $n' \ge 1$ so that $(\mathbf{P}^{nd+r(j)})_{i_0j} > 0 < (\mathbf{P}^{n'd-r(k)})_{k_{i_0}}$. Then $(\mathbf{P}^{(n+m+n')d+(r(j)+r-r(k))})_{i_0i_0} > 0$, and so d|(r(j) + r - r(k)). Equivalently, $r(k) - r(j) = r \mod d$. Conversely, if $r(k) - r(j) = r \mod d$, choose $n \ge 0$ and $n' \ge 1$ so that $(\mathbf{P}^{nd+r(k)})_{i_0k} > 0$ and $(\mathbf{P}^{n'd-r(j)})_{j_{i_0}} > 0$. Then $(\mathbf{P}^{(n+m+n')d+r})_{j_k} > 0$ for any $m \ge 1$ satisfying $(\mathbf{P}^{md})_{i_0i_0} > 0$. Since, by (3.1.14), $(\mathbf{P}^{md})_{i_0i_0} > 0$ for all sufficiently large m's, this completes the left hand inequality in (2), and the right hand inequality in (2) is proved in the same way. Finally, to check (3), note that, from (1), the restriction of \mathbf{P}^d to \mathbb{S}_r is a transition probability matrix, and by (2), it is both irreducible and aperiodic.

The existence of such a partition has several interesting consequences. In the first place, it says that the chain proceeds through the state space in a cyclic fashion: if

it starts from *i*, then after *n* steps it is in $S_{r(i)+n}$, where the addition in the subscript should be interpreted modulo *d*. In fact, with this convention for addition, we have that

$$\mathbf{P}^n \mathbf{1}_{\mathbb{S}_r} = \mathbf{1}_{\mathbb{S}_{r+n}}.\tag{4.1.19}$$

To see this, simply observe that, on the one hand, $\mathbf{P}^n \mathbf{1}_{\mathbb{S}_r}(i) = 0$ unless $i \in \mathbb{S}_{r+n}$, while, on the other hand, $\sum_{r'=0}^{d-1} \mathbf{P}^n \mathbf{1}_{\mathbb{S}_{r'}} = \mathbf{1}$. Hence, $i \notin \mathbb{S}_{r+n} \implies (\mathbf{P}^n \mathbf{1}_{\mathbb{S}_r})_i = 0$, whereas $i \in \mathbb{S}_{r+n} \implies 1 = \sum_{r'=0}^{d-1} (\mathbf{P}^n \mathbf{1}_{\mathbb{S}_{r'}})_i = (\mathbf{P}^n \mathbf{1}_{\mathbb{S}_{r+n}})_i$.

Secondly, because, for each $0 \le r < d$, the restriction of \mathbf{P}^d to \mathbb{S}_r is an irreducible, recurrent, and aperiodic transition probability matrix, we know that, for each $0 \le r < d$ and $j \in \mathbb{S}_r$, there exists a $\pi_{jj}^{(r)} \in [0, 1]$ with the property that $(\mathbf{P}^{md})_{ij} \longrightarrow \pi_{jj}^{(r)}$ for all $(i, j) \in \mathbb{S}_r^2$. More generally,

$$\lim_{m \to \infty} \left(\mathbf{P}^{md+s} \right)_{ij} = \begin{cases} \pi_{jj}^{(r(j))} & \text{if } r(j) - r(i) = s \mod d \\ 0 & \text{otherwise.} \end{cases}$$
(4.1.20)

In particular, if $(i, j) \in (\mathbb{S}_r)^2$, then

$$(\mathbf{A}_{nd})_{ij} = \frac{1}{nd} \sum_{m=0}^{n-1} \sum_{s=0}^{d-1} (\mathbf{P}^{md+s})_{ij} = \frac{1}{nd} \sum_{m=0}^{n-1} (\mathbf{P}^{md})_{ij} \longrightarrow \frac{\pi_{jj}^{(r)}}{d}.$$

Hence, since we already know that $(\mathbf{A}_n)_{ij} \longrightarrow \pi_{jj}$, it follows that

$$\pi_{jj}^{(r)} = d\pi_{jj} \quad \text{for } 0 \le r < d \text{ and } j \in \mathbb{S}_r.$$

$$(4.1.21)$$

In the case when **P** is positive recurrent on \mathbb{S} , so is the restriction of \mathbf{P}^d to each \mathbb{S}_r , and therefore $\sum_{j \in \mathbb{S}_r} \pi_{jj}^{(r)} = 1$. Thus, in the positive recurrent case, (4.1.21) leads to the interesting conclusion that $\pi^{\mathbb{S}}$ assigns probability $\frac{1}{d}$ to each \mathbb{S}_r . See Exercises 4.2.6 and 6.6.5 below for other applications of these considerations.

4.2 Exercises

Exercise 4.2.1 Just as Césaro convergence is strictly weaker (i.e., it is implied by but does not imply) than ordinary convergence, so, in this exercise, we will show that Abel convergence is strictly weaker than Césaro convergence.

(a) Assume that the radius of convergence of $\{a_n\}_0^\infty \subseteq \mathbb{R}$ is less than or equal to 1. That is, $\overline{\lim}_{n\to\infty} |a_n|^{\frac{1}{n}} \leq 1$. Set $R(s) = (1-s) \sum_0^\infty s^n a_n$ for $s \in [0, 1)$ and $A_n = \frac{1}{n} \sum_0^{n-1} a_m$ for $n \geq 1$. Show that $\overline{\lim}_{n\to\infty} |A_n| \leq 1$ and that $R(s) = (1-s)^2 \sum_1^\infty n s^{n-1} A_n$ for $s \in [0, 1)$. Use this to conclude that

$$\lim_{n \to \infty} A_n = a \in \mathbb{R} \implies \lim_{s \neq 1} R(s) = a$$

(b) Take a_n = (−1)ⁿ⁺¹n for n ≥ 0, check that the radius of convergence of {a_n}₀[∞] is 1, and show that

$$\frac{1}{n}\sum_{m=0}^{n-1} a_m = \begin{cases} \frac{1}{2} & \text{if } n \text{ is even} \\ -\frac{1}{2} + \frac{1}{2n} & \text{if } n \text{ is odd} \end{cases} \quad \text{and} \quad (1-s)\sum_{m=0}^{\infty} s^m a_m = \frac{s(1-s)}{(1+s)^2}$$

Hence, $\{a_n\}_0^\infty$ is Abel convergent to 0 but is Césaro divergent.

Exercise 4.2.2 Recall the queuing model in Exercise 1.3.12. Show that $\{Q_n : n \ge 0\}$ is an \mathbb{N} -valued Markov chain conditioned to start from 0, and write down the transition probability matrix for this chain. Further, for this chain: show that 0 is transient if $\mathbb{E}[B_1] > 0$, 0 is null recurrent if $\mathbb{E}[B_1] = 0$, and that 0 is positive recurrent if $\mathbb{E}[B_1] < 0$. In order to handle the case when $\mathbb{E}[B_1] = 0$, you might want to refer to Exercise 1.3.11.

Exercise 4.2.3 Here is a test for positive recurrence. Given a transition probability matrix **P** and an element *j* of the state space \mathbb{S} , set C = [j]. Assume that *u* is a non-negative function on *C* with the properties that

$$(\mathbf{P}u)_j < \infty$$
 and $u(i) \ge (\mathbf{P}u)_i + \epsilon$ for all $i \in C \setminus \{j\}$ and some $\epsilon > 0$.

(a) Begin by showing that

$$\mathbb{E}\left[u(X_{(n+1)\wedge\rho_j})\,\middle|\,X_0=j\right] \le \mathbb{E}\left[u(X_{n\wedge\rho_j})\,\middle|\,X_0=j\right] - \epsilon \mathbb{P}(\rho_j > n\,|\,X_0=j),$$

and use this to conclude that j is positive recurrent.

(b) Suppose that $\mathbb{S} = \mathbb{Z}$ and that $|i| \ge \sum_{j} |j|(\mathbf{P})_{ij} + \epsilon$ for all $i \in \mathbb{Z} \setminus \{0\}$. Show that 0 is positive recurrent for the chain determined by **P**.

Exercise 4.2.4 Consider the nearest neighbor random walk on \mathbb{Z} which moves forward with probability $p \in (\frac{1}{2}, 1)$ and backward with probability q = 1 - p. In other words, we are looking at the Markov chain on \mathbb{Z} whose transition probability matrix **P** is given by $(\mathbf{P})_{ij} = p$ if j = i + 1, $(\mathbf{P})_{ij} = q$ if j = i - 1, and $(\mathbf{P})_{ij} = 0$ if $|j - i| \neq 1$. Obviously, this chain is irreducible, and the results in Sects. 1.2.2–1.2.1 show that 0 is transient. Thus, the considerations in Exercise 2.4.10 apply.

(a) Construct $\hat{\mathbf{P}}$ from \mathbf{P} by the prescription in Exercise 2.4.10 when $j_0 = 0$, and, using (1.1.12), show that

$$(\hat{\mathbf{P}})_{ij} = \begin{cases} p & \text{if } i \le 0 \& j = i+1 \text{ or } i \ge 1 \& j = i-1 \\ q & \text{if } i \le 0 \& j = i-1 \text{ or } i \ge 1 \& j = i+1 \\ 0 & \text{otherwise.} \end{cases}$$

(b) On the basis of Exercise 2.4.10, we know that 0 is recurrent for the chain determined by $\hat{\mathbf{P}}$. Moreover, by part (b) of Exercise 4.2.3, one can check that it

4.2 Exercises

is positive recurrent. In fact, by combining part (b) of Exercise 2.4.10 with the computations in Sect. 1.1.4, show that

$$\hat{\mathbb{E}}[\rho_0|X_0=0] = \frac{2p}{p-q},$$

where $\hat{\mathbb{E}}$ is used to indicate that the expectation value is taken relative to the chain determined by $\hat{\mathbf{P}}$.

(c) Since **P** is irreducible, so is $\hat{\mathbf{P}}$. Hence, since 0 is positive recurrent for the chain determined by $\hat{\mathbf{P}}$, there is a unique stationary probability vector for $\hat{\mathbf{P}}$. Find this vector and use it to show that

$$\mathbb{E}^{\hat{\mathbb{P}}}[\rho_j|X_0=j] = \begin{cases} \frac{2p}{p-q} & \text{if } j \in \{0,1\}\\ \frac{2pq^j}{p^j(p-q)} & \text{if } j \le -1\\ \frac{2p^j}{q^{j-1}(p-q)} & \text{if } j \ge 2. \end{cases}$$

Exercise 4.2.5 As I said at the beginning of this chapter, we would be content here with statements about the convergence of either $\{(\mathbf{P}^n)_{ij} : n \ge 0\}$ or $\{(\mathbf{A}_n)_{ij} : n \ge 0\}$ for each $(i, j) \in \mathbb{S}^2$. However, as we are about to see, there are circumstances in which the pointwise results we have just obtained self-improve.

- (a) Assume that *j* is positive recurrent, and set C = [j]. Given a probability vector $\boldsymbol{\mu}$ with the property that $\sum_{i \notin C} (\boldsymbol{\mu})_i = 0$, show that, in general, $(\boldsymbol{\mu} \mathbf{A}_n)_i \longrightarrow \pi_{ii}$ and, when *j* is aperiodic, $(\boldsymbol{\mu} \mathbf{P}^n)_i \longrightarrow \pi_{ii}$ for each $i \in C$.
- (b) Here is an interesting fact about convergence of series. Namely, for each m ∈ N, let {a_{m,n} : n ≥ 0} be a sequence of real numbers which converges to a real number b_m as n → ∞. Further, assume that, for each n ∈ N, the sequence {a_{m,n} : m ≥ 0} is absolutely summable. Finally, assume that

$$\sum_{m=0}^{\infty} |a_{m,n}| \longrightarrow \sum_{m=0}^{\infty} |b_m| < \infty \quad \text{as } n \to \infty.$$

Show that

$$\lim_{n\to\infty}\sum_{m=0}^{\infty}|a_{m,n}-b_m|=0.$$

Hint: Using the triangle inequality, show that

$$||a_{m,n}| - |b_m| - |a_{m,n} - b_m|| \le 2|b_m|,$$

and apply Lebesgue's dominated convergence theorem to conclude that

$$\sum_{m=0}^{\infty} |a_{m,n} - b_m| \le \left| \sum_{m=0}^{\infty} (|a_{m,n}| - |b_m|) \right| + \sum_{m=0}^{\infty} ||a_{m,n}| - |b_m| - |a_{m,n} - b_m| | \longrightarrow 0$$
as $n \to \infty$.

(c) Return to the setting in (a), and use (b) together with the result in (a) to show that, in general, $\|\mu \mathbf{A}_n - \boldsymbol{\pi}^C\|_v \longrightarrow 0$, and $\|\mu \mathbf{P}^n - \boldsymbol{\pi}^C\|_v \longrightarrow 0$ when *j* is aperiodic. In particular, for each probability vector $\boldsymbol{\mu}$ with $\sum_{i \in C} (\boldsymbol{\mu})_i = 1$,

$$\lim_{n\to\infty}\sup\{\left|\boldsymbol{\mu}\mathbf{P}^{n}\mathbf{f}-\boldsymbol{\pi}^{C}\mathbf{f}\right|:\|f\|_{u}\leq 1\}=0,$$

where **f** is the column vector determined by a function f. Of course, this is still far less than what we had under Doeblin's condition since his condition provided us with a rate of convergence which was independent of μ . In general, no such uniform rate of convergence will exist.

Exercise 4.2.6 When *j* is aperiodic for **P**, we know that $\lim_{n\to\infty} (\mathbf{P}^n)_{ij}$ exists for all $i \in \mathbb{S}$ and is 0 unless *j* is positive recurrent. When d(j) > 1 and *j* is positive recurrent, show that $\lim_{n\to\infty} (\mathbf{P}^n)_{jj}$ will fail to exist. On the other hand, even if d(j) > 1, show that $\lim_{n\to\infty} (\mathbf{P}^n)_{ij} = 0$ for any $j \in \mathbb{S}$ which is not positive recurrent.

Exercise 4.2.7 Here is an important interpretation of π^C when *C* is a positive recurrent communicating class. Let *i* be a recurrent state and, for $k \in S$, let μ_k be the expected number of times the chain visits *k* before returning to *i* given that the chain started from *i*:

$$\mu_k = \mathbb{E}\left[\sum_{m=0}^{\rho_i - 1} \mathbf{1}_{\{k\}}(X_m) \, \middle| \, X_0 = i\right] \in [0, \infty].$$

Determine the row vector $\boldsymbol{\mu} \in [0, \infty]^{\mathbb{S}}$ by $(\boldsymbol{\mu})_k = \mu_k$.

(a) Show that, for all $j \in \mathbb{S}$,

$$(\boldsymbol{\mu}\mathbf{P})_j = \mathbb{E}\left[\sum_{m=1}^{\rho_i} \mathbf{1}_{\{j\}}(X_m) \,\Big|\, X_0 = i\right] = \mu_j.$$

Thus, without any further assumptions about *i*, μ is **P**-stationary in the sense that $\mu = \mu \mathbf{P}$.

(b) Clearly $\mu_i = 1$ and $\sum_j \mu_j = \infty$ unless *i* is positive recurrent. Nonetheless, show that $\mu_j = 0$ unless $i \leftrightarrow j$ and that $\mu_j \in (0, \infty)$ if $i \leftrightarrow j$.

Hint: Show that

$$\mathbb{P}(\rho_j^{(m)} < \rho_i \mid X_0 = i) = \mathbb{P}(\rho_j < \rho_i \mid X_0 = j)^{m-1} \mathbb{P}(\rho_j < \rho_i \mid X_0 = i).$$

(c) If i is positive recurrent, show that

$$\bar{\boldsymbol{\mu}} \equiv \frac{\boldsymbol{\mu}}{\sum_k \mu_k} = \boldsymbol{\pi}^C.$$

Equivalently, when *i* is positive recurrent,

$$(\boldsymbol{\pi}^{C})_{j} = \frac{\mathbb{E}[\sum_{m=0}^{\rho_{i}-1} \mathbf{1}_{\{j\}}(X_{m}) | X_{0} = i]}{\mathbb{E}[\rho_{i} | X_{0} = i]}.$$

In words, $(\pi^{C})_{j}$ is the relative expected amount of time that the chains spends at *j* before returning to *i*.

Exercise 4.2.8 Let **P** be a transition probability on a finite state space \mathbb{S} , and assume that π is the one and only stationary probability for **P**. Given a recurrent state $i \in \mathbb{S}$ and a $j \neq i$, show that (cf. Sect. 3.2.1)

$$\mathbb{P}(\rho_j < \rho_i \mid X_0 = i) = \frac{\det((\mathbf{I} - \mathbf{P})^{\{j\}})}{\det((\mathbf{I} - \mathbf{P})^{\{i,j\}})}.$$

You might want to take the following steps.

(a) Since S is finite, we know that *i* is positive recurrent. Now use the uniqueness of π , Exercises 4.2.7, and 3.2.7 to show that

$$\frac{\det((\mathbf{I}-\mathbf{P})^{\{j\}})}{\det((\mathbf{I}-\mathbf{P})^{\{i\}})} = \mathbb{E}\left[\sum_{m=0}^{\rho_i-1} \mathbf{1}_{\{j\}}(X_m) \mid X_0 = i\right].$$

(b) Use the Markov property to show that

$$\mathbb{E}\left[\sum_{m=0}^{\rho_{i}-1}\mathbf{1}_{\{j\}}(X_{m}) \mid X_{0}=i\right] = \mathbb{P}(\rho_{j} < \rho_{i} \mid X_{0}=i)\mathbb{E}\left[\sum_{m=0}^{\rho_{i}-1}\mathbf{1}_{\{j\}}(X_{m}) \mid X_{0}=j\right].$$

(c) Use (3.2.11) together with (b) to arrive at the desired conclusion.

Exercise 4.2.9 We continue with the program initiated in Exercise 4.2.7 but assume now that the reference point *i* is null recurrent. In this case, $\sum_{j} (\boldsymbol{\mu})_{j} = \infty$ when $\boldsymbol{\mu} \in [0, \infty)^{\mathbb{S}}$ is the **P**-stationary measure introduced in Exercise 4.2.7. In this exercise we will show that, up to a multiplicative constant, $\boldsymbol{\mu}$ is the only **P**-stationary $\boldsymbol{\nu} \in [0, \infty)^{\mathbb{S}}$ with the property that $(\boldsymbol{\nu})_{j} = 0$ unless $i \rightarrow j$ (and therefore $i \leftrightarrow j$). Equivalently, given such a $\boldsymbol{\nu}, \boldsymbol{\nu} = (\boldsymbol{\nu})_{i} \boldsymbol{\mu}$.

(a) Assume that $\mathbf{v} \in [0, \infty)^{\mathbb{S}}$ satisfies $\mathbf{v} = \mathbf{v}\mathbf{P}$. If $(\mathbf{v})_i = 1$, show that, for all $j \in \mathbb{S}$ and $n \ge 0$,

$$(\mathbf{v})_{j} = \sum_{k \neq i} (\mathbf{v})_{k} \mathbb{P}(X_{n} = j \& \rho_{i} > n \mid X_{0} = k) \\ + \mathbb{E} \left[\sum_{m=0}^{n \land (\rho_{i} - 1)} \mathbf{1}_{\{j\}}(X_{m}) \mid X_{0} = i \right].$$

Hint: Work by induction on $n \ge 0$. When n = 0 there is nothing to do. To carry out the inductive step, use (2.1.1) and Fubini's theorem to show that

$$\sum_{k \neq i} (\mathbf{v})_k \mathbb{P}(X_n = j \& \rho_i > n \mid X_0 = k)$$

= $\sum_{k \neq i} (\mathbf{v} \mathbf{P})_k \mathbb{P}(X_n = j \& \rho_i > n \mid X_0 = k)$
= $\sum_{\ell} (\mathbf{v})_{\ell} \mathbb{P}(X_{n+1} = j \& \rho_i > n+1 \mid X_0 = \ell).$

(b) Assuming that $\mathbf{v} = \mathbf{v}\mathbf{P}$ and $(\mathbf{v})_i = 0$ unless $i \to j$, show that $\mathbf{v} = (\mathbf{v})_i \boldsymbol{\mu}$.

Hint: First show that $\mathbf{v} = \mathbf{0}$ if $(\mathbf{v})_i = 0$, and thereby reduce to the case when $(\mathbf{v})_i = 1$. Starting from the result in (a), apply the monotone convergence theorem to see that $(\mathbf{v})_j \ge (\boldsymbol{\mu})_j$ for all $j \in \mathbb{S}$. Now consider $\boldsymbol{\omega} \equiv \mathbf{v} - \boldsymbol{\mu}$, and conclude that $\boldsymbol{\omega} = \mathbf{0}$.

Exercise 4.2.10 Let *C* be a communicating class of positive recurrent states. The reason why Theorem 4.1.14 is called a "mean ergodic theorem" is that the asserted convergence is taking place in the sense of mean square convergence. Of course, mean square convergence implies convergence in probability, but, in general, it cannot be used to get convergence with probability 1. Nonetheless, as you are to show here, when $\mathbb{P}(X_0 \in C) = 1$ and $j \in C$,

$$\lim_{n \to \infty} \frac{1}{n} \sum_{m=0}^{n-1} \mathbf{1}_{\{j\}}(X_m) = \pi_{jj} \quad \text{with probability 1.}$$
(4.2.11)

Observe that in Sect. 2.3.3 we proved the individual ergodic theorem (2.3.10) under Doeblin's condition, and (4.2.11) says that the same sort of individual ergodic theorem holds even when Doeblin's condition is not present. In fact, there is a very general result, of which (4.2.11) is a very special case, which was proved originally by G.D. Birkhoff. However, we will not follow Birkhoff and instead, as we did in Sect. 2.3.3, we will base our proof on the strong saw of large numbers, although this time we need the full statement which holds (cf. Theorem 3.3.10 in [8])) for averages of mutually independent, identically distributed, integrable random variables.

(a) Show that it suffices to prove the result when P(X₀ = i) = 1 for some i ∈ C.
(b) Set ρ_i⁽⁰⁾ = 0, and use ρ_i^(m) to denote the time of the *m*th return to *i*. If

$$au_m =
ho_i^{(m)} -
ho_i^{(m-1)}$$
 and $Y_m = \sum_{\ell=
ho_i^{(m-1)}}^{
ho_i^{(m)}-1} \mathbf{1}_{\{j\}}(X_\ell),$

show that, conditional on $X_0 = i$, both $\{\tau_m : m \ge 1\}$ and $\{Y_m : m \ge 1\}$ are sequences of mutually independent, identically distributed, non-negative, integrable random variables. In particular, as an application of the strong law of

4.2 Exercises

large numbers and (4.1.5) plus the result in Exercise 4.2.7, conclude that, conditional on $X_0 = i$,

$$\lim_{m \to \infty} \frac{\rho_i^{(m)}}{m} = \frac{1}{\pi_{ii}} \quad \text{and} \quad \lim_{m \to \infty} \frac{1}{m} \sum_{\ell=0}^{\rho_i^{(m)} - 1} \mathbf{1}_{\{j\}}(X_m) = \frac{\pi_{jj}}{\pi_{ii}} \tag{(*)}$$

with probability 1. Hence, $\lim_{m\to\infty} \frac{1}{\rho_i^{(m)}} \sum_{\ell=0}^{\rho_i^{(m)}-1} \mathbf{1}_{\{j\}}(X_m) = \pi_{jj}$ with probability 1.

(c) In view of the results in (a) and (b), we will be done once we check that, conditional on X₀ = i,

$$\lim_{n \to \infty} \left| \frac{1}{n} \sum_{0}^{n-1} \mathbf{1}_{\{j\}}(X_{\ell}) - \frac{1}{\rho_i^{(m_n)}} \sum_{0}^{\rho_i^{(m_n)} - 1} \mathbf{1}_{\{j\}}(X_{\ell}) \right| = 0$$

with probability 1, where m_n is the \mathbb{Z}^+ -valued random variable determined so that $\rho_i^{(m_n-1)} \leq n < \rho_i^{(m_n)}$. To this end, first show that

$$\left|\frac{1}{n}\sum_{0}^{n-1}\mathbf{1}_{\{j\}}(X_{\ell}) - \frac{1}{\rho_{i}^{(m_{n})}}\sum_{0}^{\rho_{i}^{(m_{n})}-1}\mathbf{1}_{\{j\}}(X_{\ell})\right| \leq \frac{\tau_{m_{n}}}{m_{n}}$$

Next, from the first part of (*), show that $\mathbb{P}(\lim_{n\to\infty} m_n = \infty | X_0 = i) = 1$. Finally, check that, for any $0 < \epsilon < 1$,

$$\mathbb{P}\left(\sup_{m \ge M} \frac{\tau_m}{m} \ge \epsilon \mid X_0 = i\right)$$

$$\leq \sum_{m=M}^{\infty} \mathbb{P}(\rho_i \ge m\epsilon \mid X_0 = i) \le \frac{1}{\epsilon} \mathbb{E}[\rho_i, \rho_i \ge M\epsilon \mid X_0 = i],$$

and use this to complete the proof of (4.2.11).

(d) Introduce the *empirical measure* \mathbf{L}_n , which is the random probability vector measuring the average time spent at points. That is, $(\mathbf{L}_n)_i = \frac{1}{n} \sum_{0}^{n-1} \mathbf{1}_{\{i\}}(X_m)$. By combining the result proved here with the one in (b) of Exercise 4.2.5, conclude that $\lim_{n\to\infty} \|\mathbf{L}_n - \boldsymbol{\pi}^C\|_v = 0$ with probability 1 when $\mathbb{P}(X_0 \in C) = 1$.

Exercise 4.2.12 Although the statement in Exercise 4.2.10 applies only to positive recurrent states, it turns out that there is a corresponding limit theorem for states which are not positive recurrent. Namely, show that if j is not positive recurrent, then, no matter what the initial distribution,

$$\mathbb{P}\left(\lim_{n\to\infty}\frac{1}{n}\sum_{m=0}^{n-1}\mathbf{1}_{\{j\}}(X_m)=0\right)=1.$$

When *j* is transient, $\mathbb{E}[\sum_{0}^{\infty} \mathbf{1}_{\{j\}}(X_m)] < \infty$ and therefore the result is trivial. To handle *j* that are null recurrent, begin by noting that it suffices to handle the case when $\mathbb{P}(X_0 = j) = 1$. Next, note that, conditional on $X_0 = j$, $\{\rho_j^{(m+1)} - \rho_j^{(m)} : m \ge 0\}$ is a sequence of independent, identically distributed, \mathbb{Z}^+ -valued random variables, and apply the strong law of large numbers to see that, for any $R \ge 1$,

$$\mathbb{P}\left(\lim_{m \to \infty} \frac{\rho_j^{(m)}}{m} \ge H(R) \mid X_0 = j\right)$$
$$\ge \mathbb{P}\left(\lim_{m \to \infty} \frac{\rho_j^{(m)} \land R}{m} \ge H(R) \mid X_0 = j\right) = 1,$$

where $H(R) \equiv \frac{1}{2}\mathbb{E}[\rho_j \wedge R | X_0 = j] \nearrow \infty$ as $R \nearrow \infty$. Hence, given $X_0 = j$, $\frac{\rho_j^{(m)}}{m} \longrightarrow \infty$ with probability 1. Finally, check that, for any $\epsilon > 0$,

$$\mathbb{P}\left(\sup_{n\geq N}\frac{1}{n}\sum_{0}^{n-1}\mathbf{1}_{\{j\}}(X_m)\geq\epsilon\mid X_0=j\right)$$
$$\leq \mathbb{P}\left(\sup_{n\geq N}\frac{\rho_j^{(\lfloor n\epsilon \rfloor)}}{n}\leq\frac{1}{\epsilon}\mid X_0=j\right)$$
$$\leq \mathbb{P}\left(\sup_{m\geq n\epsilon}\frac{\rho_j^{(m)}}{m}\leq\frac{1}{\epsilon}\mid X_0=j\right),$$

and combine this with the preceding to reach the desired conclusion.

Chapter 5 Markov Processes in Continuous Time

Up until now we have been dealing with Markov processes for which time has a discrete parameter $n \in \mathbb{N}$. In this chapter we will introduce Markov processes for which time has a continuous parameter $t \in [0, \infty)$, even though our processes will continue to take their values in a countable state space S.

5.1 Poisson Processes

Just as Markov chains with independent, identically distributed increments (a.k.a. random walks) on \mathbb{Z}^d are the simplest discrete parameter Markov processes, so the simplest continuous time Markov processes are those whose increments are mutually independent and *homogeneous* in the sense that the distribution of an increment depends only on the length of the time interval over which the increment is taken. More precisely, we will be dealing in this section with \mathbb{Z}^d -valued stochastic processes $\{X(t) : t > 0\}$ with the property that $\mathbb{P}(X(0) = \mathbf{0}) = 1$ and

$$\mathbb{P}(X(t_1) - X(t_0) = j_1, \dots, X(t_n) - X(t_{n-1}) = j_n)$$
$$= \prod_{m=1}^n \mathbb{P}(X(t_m - t_{m-1}) = j_m)$$

for all $n \ge 1, 0 \le t_0 < \cdots < t_n$, and $(j_1, \dots, j_n) \in (\mathbb{Z}^d)^n$.

5.1.1 The Simple Poisson Process

The *simple Poisson process* is the N-valued stochastic process $\{N(t) : t \ge 0\}$ that starts from 0, sits there for a unit exponential holding time E_1 (i.e., N(t) = 0 for $t \in [0, E_1)$ and, for t > 0, $\mathbb{P}(E_1 > t) = e^{-t}$), at time E_1 moves to 1 where it remains

for an independent unit exponential holding time E_2 , moves to 2 at time $E_1 + E_2$, etc. More precisely, if $\{E_n : n \ge 1\}$ is a sequence of mutually independent, unit exponential random variables and

$$J_n = \begin{cases} 0 & \text{when } n = 0\\ \sum_{m=1}^n E_m & \text{when } n \ge 1, \end{cases}$$

then the stochastic process $\{N(t) : t \ge 0\}$ given by

$$N(t) = \max\{n \ge 0 : J_n \le t\}$$
(5.1.1)

is a simple Poisson process. When thinking about the meaning of (5.1.1), keep in mind that, because,

with probability 1,
$$E_n > 0$$
 for all $n \ge 1$ and $\sum_{m=1}^{\infty} E_m = \infty$,

with probability 1 the path $t \in [0, \infty) \mapsto N(t)$ is piecewise constant, right continuous, and, when it jumps, it jumps by +1: $N(t) - N(t-) \in \{0, 1\}$ for all t > 0, where $N(t-) \equiv \lim_{s \nearrow t} N(s)$ is the left limit of $N(\cdot)$ at t.

We now show that $\{N(t) : t \ge 0\}$ moves along in independent, homogeneous increments.¹ That is, for each $s, t \in [0, \infty)$, we will prove that the increment N(s+t) - N(s) is independent of (cf. Sect. 7.1.3) $\sigma(\{N(\tau) : \tau \in [0, s]\})$ and has the same distribution as N(t):

$$\mathbb{P}\big(N(s+t) - N(s) = n \mid N(\tau), \tau \in [0, s]\big) = \mathbb{P}\big(N(t) = n\big), \quad n \in \mathbb{N}.$$
(5.1.2)

Equivalently, what we have to check is that when $(s, t) \in [0, \infty)^2$ and $A \in \sigma(\{N(\tau) : \tau \in [0, s]\})$, $\mathbb{P}(\{N(s+t) - N(s) \ge n\} \cap A) = \mathbb{P}(N(t) \ge n)\mathbb{P}(A)$ for all $n \in \mathbb{N}$. Since this is trivial when n = 0, we will assume that $n \in \mathbb{Z}^+$. In addition, since we can always write A as the disjoint union of the sets $A \cap \{N(s) = m\}$, $m \in \mathbb{N}$, and each of these is again in $\sigma(\{N(\tau) : \tau \in [0, s]\})$, we may and will assume that, for some m, N(s) = m on A. But in that case we can write $A = \{J_{m+1} > s\} \cap B$ where $B \in \sigma(\{E_1, \ldots, E_m\})$ and $J_m \le s$ on B. Hence, since $N(s + t) \ge m + n \iff J_{m+n} \le s + t$, and $\sigma(\{J_\ell - J_m : \ell > m\})$ is independent of $\sigma(\{E_k : k \le m\})$, an application of (7.4.2) shows that

$$\mathbb{P}(\{N(s+t) - N(s) \ge n\} \cap A)$$

= $\mathbb{P}(\{J_{m+n} \le s+t\} \cap \{J_{m+1} > s\} \cap B)$
= $\mathbb{P}(\{J_{m+n} - J_m \le s+t - J_m\} \cap \{J_{m+1} - J_m > s - J_m\} \cap B)$
= $\mathbb{E}[v(J_m), B],$

¹I have chosen the following proof in order to develop a line of reasoning which will serve us well later on. A more straight-forward, but less revealing, proof that the simple Poisson process has independent, homogeneous increments is given in Exercise 5.5.1 below.

where, for $\xi \in [0, s]$,

$$v(\xi) \equiv \mathbb{P}(\{J_{m+n} - J_m \le s + t - \xi\} \cap \{J_{m+1} - J_m > s - \xi\})$$

= $\mathbb{P}(\{J_n \le s + t - \xi\} \cap \{E_1 > s - \xi\}) = \mathbb{P}(\{J_n \le s + t - \xi\} \cap \{E_n > s - \xi\})$
= $\mathbb{P}(\{J_{n-1} + E_n \le s + t - \xi\} \cap \{E_n > s - \xi\}) = \mathbb{E}[w(\xi, E_n), E_n > s - \xi]$

when $w(\xi, \eta) \equiv \mathbb{P}(J_{n-1} \le s + t - \xi - \eta)$ for $\xi \in [0, s]$ and $\eta \in [s - \xi, s + t - \xi]$.

Up to this point we have not used any property of exponential random variables other than that they are positive. However, in our next step we will use their characteristic property, namely, the fact that an exponential random variable E "has no memory." That is, $\mathbb{P}(E > a + b | E > a) = \mathbb{P}(E > b)$, from which it is an easy step to

$$\mathbb{E}[f(E), E > a] = e^{-a} \mathbb{E}[f(a+E)]$$
(5.1.3)

for any non-negative, $\mathcal{B}_{[0,\infty)}$ -measurable function f. In particular, this means that

$$v(\xi) = \mathbb{E} \Big[w(\xi, E_n), E_n > s - \xi \Big] = e^{-(s-\xi)} \mathbb{E} \Big[w(\xi, E_n + s - \xi) \Big]$$

= $e^{-(s-\xi)} \mathbb{P} (J_{n-1} \le t - E_n) = e^{-(s-\xi)} \mathbb{P} (J_n \le t) = e^{-(s-\xi)} \mathbb{P} \Big(N(t) \ge n \Big).$

That is, we have shown that

$$\mathbb{P}(\{N(s+t)-N(s)\geq n\}\cap A)=\mathbb{E}[e^{-(s-J_m)},B]\mathbb{P}(N(t)\geq n).$$

Finally, since

$$\mathbb{P}(A) = \mathbb{P}(\{J_{m+1} > s\} \cap B) = \mathbb{P}(\{E_{m+1} > s - J_m\} \cap B) = \mathbb{E}[e^{-(s-J_m)}, B],$$

the proof of (5.1.2) is complete. Hence, we now know that the simple Poisson process $\{N(t) : t \ge 0\}$ has homogeneous and mutually independent increments.

Before moving on, we must still find out what is the distribution of N(t). But, the sum of *n* mutually independent, unit exponential random variables has a $\Gamma(n)$ -distribution, and so

$$\mathbb{P}(N(t) = n) = \mathbb{P}(J_n \le t < J_{n+1}) = \mathbb{P}(J_n \le t) - \mathbb{P}(J_{n+1} \le t)$$
$$= \frac{1}{(n-1)!} \int_0^t \tau^{n-1} e^{-\tau} d\tau - \frac{1}{n!} \int_0^t \tau^n e^{-\tau} d\tau = \frac{t^n}{n!}.$$

In other words, N(t) is a Poisson random variable with mean t. More generally, when we combine this with (5.1.2), we get that, for all $s, t \in [0, \infty)$ and $n \in \mathbb{N}$,

$$\mathbb{P}\big(N(s+t) - N(s) = n \,|\, N(\tau), \, \tau \in [0, s]\big) = e^{-t} \frac{t^n}{n!}.$$
(5.1.4)

Alternatively, again starting from (5.1.2), we can now give the following Markovian description of $\{N(t) : t \ge 0\}$:

$$\mathbb{P}\big(N(s+t) = n \mid N(\tau), \tau \in [0,s]\big) = e^{-t} \frac{t^{n-N(s)}}{(n-N(s))!} \mathbf{1}_{[0,n]}\big(N(s)\big).$$
(5.1.5)

5.1.2 Compound Poisson Processes on \mathbb{Z}^d

Having constructed the simplest Poisson process, we can easily construct a rich class of processes which are the continuous time analogs of random walks on \mathbb{Z}^d . Suppose that μ is a probability vector on \mathbb{Z}^d which gives 0 mass to the origin **0**. Then the *compound Poisson process with jump distribution* μ *and rate* $R \in (0, \infty)$ is the stochastic process {**X**(t) : $t \ge 0$ } which starts at **0**, sits there for an exponential holding time having mean value R^{-1} , at which time it jumps by the amount $\mathbf{k} \in \mathbb{Z}^d$ with probability (μ)_{**k**}, sits where it lands for another, independent holding time with mean R^{-1} , jumps again a random amount with distribution μ , etc. Thus, the simple Poisson process is the case when d = 1, (μ)₁ = 1, and R = 1. Of course, the amount by which the simple Poisson process jumps is deterministic whereas the amount by which a compound Poisson process jumps will, in general, be random.

To construct a compound Poisson process, let $\{\mathbf{B}_n : n \ge 1\}$ be a sequence of mutually independent \mathbb{Z}^d -valued random variables with distribution μ , introduce the random walk $\mathbf{X}_0 = \mathbf{0}$ and $\mathbf{X}_n = \sum_{m=1}^n \mathbf{B}_m$ for $n \ge 1$, and define $\{\mathbf{X}(t) : t \ge 0\}$ so that $\mathbf{X}(t) = \mathbf{X}_{N(Rt)}$, where $\{N(t) : t \ge 0\}$ is a simple Poisson process that is independent of the \mathbf{B}_m 's. The existence of all these random variables is guaranteed (cf. the footnote 8 in Sect. 1.2.1) by Theorem 7.3.2. Obviously, $\mathbf{X}(0) = \mathbf{0}$ and $t \in [0, \infty) \longrightarrow \mathbf{X}(t) \in \mathbb{Z}^d$ is a piecewise constant, right continuous \mathbb{Z}^d -valued path. In addition, because $(\mu)_0 = 0$, it is clear² that the number of jumps that $t \rightsquigarrow \mathbf{X}(t)$ makes during a time interval (s, t] is precisely N(Rt) - N(Rs) and that $\mathbf{X}_n - \mathbf{X}_{n-1}$ is the amount by which $t \rightsquigarrow \mathbf{X}(t)$ jumps on its *n*th jump. Thus, if $J_0 \equiv 0$ and, for $n \geq 1$, J_n is the time of the *n*th jump of $t \rightsquigarrow \mathbf{X}(t)$, then $N(Rt) = n \iff J_n \leq Rt < J_{n+1}$, and $\mathbf{X}(J_n) - \mathbf{X}(J_{n-1}) = \mathbf{X}_n - \mathbf{X}_{n-1}$. Equivalently, if $\{E_n : n \ge 1\}$ denotes the sequence of unit exponential random variables out of which $\{N(t): t \ge 0\}$ is built, then $J_n - J_{n-1} = \frac{E_n}{R}$, $\mathbf{X}(t) - \mathbf{X}(t-) = \mathbf{0}$ for $t \in (J_{n-1}, J_n)$, and $\mathbf{X}(J_n) - \mathbf{X}(J_{n-1}) = \mathbf{B}_n$. Therefore { $\mathbf{X}(t) : t \ge 0$ } is indeed a compound Poisson process with jump distribution μ and rate R.

²This is the reason for my having assumed that $(\mu)_0 = 0$. However, one should realize that this assumption causes no loss in generality. Namely, if $(\mu)_0 = 1$, then the resulting compound process would be trivial: it would never move. On the other hand, if $(\mu)_0 \in (0, 1)$, then we could replace μ by $\bar{\mu}$, where $(\bar{\mu})_0 = 0$ and $(\bar{\mu})_k = (1 - (\mu)_0)^{-1}(\mu)_k$ when $k \neq 0$, and *R* by $\bar{R} = (1 - (\mu)_0)R$. The compound Poisson process corresponding to $\bar{\mu}$ and \bar{R} would have exactly the same distribution of the one corresponding to μ and *R*.

5.1 Poisson Processes

We show next that a compound process moves along in homogeneous, mutually independent increments:

$$\mathbb{P}\big(\mathbf{X}(s+t) - \mathbf{X}(s) = \mathbf{k} \,|\, \mathbf{X}(\tau), \tau \in [0, s]\big) = \mathbb{P}\big(\mathbf{X}(t) = \mathbf{k}\big), \quad \mathbf{k} \in \mathbb{Z}^d.$$
(5.1.6)

For this purpose, we use the representation $\mathbf{X}(t) = \mathbf{X}_{N(Rt)}$ introduced above. Given $A \in \sigma({\mathbf{X}(\tau) : \tau \in [0, s]})$, we need to show that

$$\mathbb{P}(\{\mathbf{X}(s+t) - \mathbf{X}(s) = \mathbf{k}\} \cap A) = \mathbb{P}(\{\mathbf{X}(s+t) - \mathbf{X}(s) = \mathbf{k}\})\mathbb{P}(A),$$

and, just as in the derivation of (5.1.2), we will, without loss in generality, assume that, for some $m \in \mathbb{N}$, N(Rs) = m on A. But then A is independent of $\sigma(\{\mathbf{X}_{m+n} - \mathbf{X}_m : n \ge 0\} \cup \{N(R(s+t)) - N(Rs)\})$, and so

$$\mathbb{P}\left(\left\{\mathbf{X}(s+t) - \mathbf{X}(s) = \mathbf{k}\right\} \cap A\right)$$

$$= \sum_{n=0}^{\infty} \mathbb{P}\left(\left\{\mathbf{X}(s+t) - \mathbf{X}(s) = \mathbf{k} \& N(R(s+t)) - N(Rs) = n\right\} \cap A\right)$$

$$= \sum_{n=0}^{\infty} \mathbb{P}\left(\left\{\mathbf{X}_{m+n} - \mathbf{X}_m = \mathbf{k} \& N(R(s+t)) - N(Rs) = n\right\} \cap A\right)$$

$$= \sum_{n=0}^{\infty} \mathbb{P}(\mathbf{X}_n = \mathbf{k}) \mathbb{P}(N(Rt) = n) \mathbb{P}(A)$$

$$= \sum_{n=0}^{\infty} \mathbb{P}\left(\mathbf{X}_n = \mathbf{k} \& N(Rt) = n\right) \mathbb{P}(A) = \mathbb{P}\left(\mathbf{X}(t) = \mathbf{k}\right) \mathbb{P}(A).$$

Hence, (5.1.6) is proved.

Finally, to compute the distribution of $\mathbf{X}(t)$, begin by recalling that the distribution of the sum of *n* independent, identically distributed random variables is the *n*-fold convolution of their distribution. Thus, $\mathbb{P}(\mathbf{X}_n = \mathbf{k}) = (\boldsymbol{\mu}^{\star n})_{\mathbf{k}}$, where $(\boldsymbol{\mu}^{\star 0})_{\mathbf{k}} = \delta_{0,\mathbf{k}}$ is the point mass at **0** and

$$(\boldsymbol{\mu}^{\star n})_{\mathbf{k}} = \sum_{\mathbf{j} \in \mathbb{Z}^d} (\boldsymbol{\mu}^{\star (n-1)})_{\mathbf{k}-\mathbf{j}}(\boldsymbol{\mu})_{\mathbf{j}} \text{ for } n \ge 1.$$

Hence,

$$\mathbb{P}(\mathbf{X}(s+t) = \mathbf{k}) = \sum_{n=0}^{\infty} \mathbb{P}(\mathbf{X}_n = \mathbf{k} \& N(Rt) = n) = e^{-Rt} \sum_{n=0}^{\infty} \frac{(Rt)^n}{n!} (\boldsymbol{\mu}^{\star n})_{\mathbf{k}}.$$

Putting this together with (5.1.6), we have that for $A \in \sigma(\{N(\tau) : \tau \in [0, s]\})$,

$$\mathbb{P}\left(\left\{\mathbf{X}(s+t) = \mathbf{k}\right\} \cap A\right)$$

= $\sum_{\mathbf{j} \in \mathbb{Z}^d} \mathbb{P}\left(\left\{\mathbf{X}(s+t) = \mathbf{k}\right\} \cap A \cap \left\{\mathbf{X}(s) = \mathbf{j}\right\}\right)$
= $\sum_{\mathbf{j} \in \mathbb{Z}^d} \mathbb{P}\left(\left\{\mathbf{X}(s+t) - \mathbf{X}(s) = \mathbf{k} - \mathbf{j}\right\} \cap A \cap \left\{\mathbf{X}(s) = \mathbf{j}\right\}\right)$
= $\sum_{\mathbf{j} \in \mathbb{Z}^d} \left(\mathbf{P}(t)\right)_{\mathbf{jk}} \mathbb{P}\left(A \cap \left\{\mathbf{X}(s) = \mathbf{j}\right\}\right) = \mathbb{E}\left[\left(\mathbf{P}(t)\right)_{\mathbf{X}(s)\mathbf{k}}, A\right],$

where

$$\left(\mathbf{P}(t)\right)_{\mathbf{k}\boldsymbol{\ell}} \equiv e^{-Rt} \sum_{m=0}^{\infty} \frac{(Rt)^m}{m!} \left(\boldsymbol{\mu}^{\star m}\right)_{\boldsymbol{\ell}-\mathbf{k}}.$$
(5.1.7)

Equivalently, we have proved that $\{\mathbf{X}(t) : t \ge 0\}$ is a continuous time Markov process with transition probability function $t \rightsquigarrow \mathbf{P}(t)$ in the sense that

$$\mathbb{P}(\mathbf{X}(s+t) = \mathbf{k} \,|\, \mathbf{X}(\tau), \, \tau \in [0, s]) = (\mathbf{P}(t))_{\mathbf{X}(s)\mathbf{k}}.$$
(5.1.8)

Observe that, as a consequence of (5.1.8), we find that $\{\mathbf{P}(t) : t \ge 0\}$ is a *semigroup*. That is, $\mathbf{P}(0) = \mathbf{I}, \mathbf{P}(t)$ is a transition probability for each $t \ge 0$, and $t \rightsquigarrow \mathbf{P}(t)$ satisfies the *Chapman–Kolmogorov* equation

$$\mathbf{P}(s+t) = \mathbf{P}(s)\mathbf{P}(t), \quad s, t \in [0, \infty).$$
(5.1.9)

Indeed, it is obvious that $\mathbf{P}(0) = \mathbf{I}$ and that $\mathbf{P}(t)$ is a transition probability for each $t \ge 0$. In addition,

$$\begin{aligned} \left(\mathbf{P}(s+t) \right)_{\mathbf{0k}} &= \sum_{\mathbf{j} \in \mathbb{Z}^d} \mathbb{P} \left(\mathbf{X}(s+t) = \mathbf{k} \& \mathbf{X}(s) = \mathbf{j} \right) \\ &= \sum_{\mathbf{j} \in \mathbb{Z}^d} \left(\mathbf{P}(t) \right)_{\mathbf{j\ell}} \left(\mathbf{P}(s) \right)_{\mathbf{0j}} = \left(\mathbf{P}(s) \mathbf{P}(t) \right)_{\mathbf{0k}}, \end{aligned}$$

from which the asserted matrix equality follows immediately when one remembers that $(\mathbf{P}(\tau))_{k\ell} = (\mathbf{P}(\tau))_{\mathbf{0}(\ell-\mathbf{k})}$.

5.2 Markov Processes with Bounded Rates

There are two directions in which one can generalize the preceding without destroying the Markov property. For one thing, one can make the distribution of jumps depend on where the process is at the time it makes the jump. This change comes down to replacing the random walk in the compound Poisson process by a more general Markov chain. The second way is to increase the randomness by making the rate of jumping depend on the place where the process is waiting before it jumps. That is, instead of the holding times all having the same mean value, the holding time at a particular state may depend on that state.

5.2.1 Basic Construction

Let S be a finite or countable state space and **P** a transition probability matrix on S. In addition, let $\mathfrak{R} = \{R_i : i \in S\}$ be a bounded set of non-negative numbers. We will show in this subsection that for each probability vector μ there exists a stochastic process $\{X(t) : t \ge 0\}$ with the properties that

- (a) $t \rightsquigarrow X(t)$ is right continuous and piecewise constant,
- (b) $\mathbb{P}(X(0) = i) = (\boldsymbol{\mu})_i$ for all $i \in \mathbb{S}$ and

$$\lim_{h \searrow 0} \sup_{\substack{j \in \mathbb{S} \\ t \in [0,\infty)}} \frac{1}{h} \Big| \mathbb{P} \Big(X(t+h) = j \,|\, X(\tau), \, \tau \in [0,t] \Big) \\ - (1 - hR_{X(t)}) \delta_{X(t),j} + hR_{X(t)} \mathbf{P}_{X(t)j} \Big| = 0.$$
(5.2.1)

Furthermore, we will show that $\{X(t) : t \ge 0\}$ is a Markov process and that its distribution is uniquely determined by (a) and (b).

Set $\mathbf{Q} = \mathbf{R}(\mathbf{P} - \mathbf{I})$, where **R** is the diagonal matrix with $(\mathbf{R})_{ii} = R_i$, and define $t \rightsquigarrow \mathbf{P}(t)$ by

$$\mathbf{P}(t) = \sum_{n=0}^{\infty} \frac{t^n}{n!} \mathbf{Q}^n.$$
(5.2.2)

Since $\|\mathbf{Q}\|_{u,v} \leq 2R$, where $R = \sup_{i \in \mathbb{S}} R_i$, it is easy to check that the series in its definition is converging in $\|\cdot\|_{u,v}$ -norm to $\mathbf{P}(t)$ uniformly for *t*'s in bounded intervals. Also, it is easy to see that $t \rightsquigarrow \mathbf{P}(t)$ satisfies the Chapman–Kolmogorov (5.1.9) equation. Less obvious is the fact that $\mathbf{P}(t)$ is a transition probability for each $t \geq 0.3$ To verify this, first note that there is nothing to do unless R > 0, since $\mathbf{P}(t) \equiv \mathbf{I}$ when R = 0. Thus, assume that R > 0 and define

$$(\tilde{\mathbf{P}})_{ij} = \begin{cases} 1 - \frac{R_i}{R} (1 - (\mathbf{P})_{ii}) & \text{if } j = i \\ \frac{R_i}{R} (\mathbf{P})_{ij} & \text{if } j \neq i. \end{cases}$$
(5.2.3)

³It is important to observe that $\|\mathbf{P}(t)\|_{u,v}$ can be as large as e^{2tR} when t < 0 and therefore will not be a transition probability matrix when R > 0 and t < 0.

Clearly $\tilde{\mathbf{P}}$ is a transition probability matrix and $\mathbf{Q} = R(\tilde{\mathbf{P}} - \mathbf{I})$. Hence, another expression for $\mathbf{P}(t)$ is

$$\sum_{n=0}^{\infty} \frac{(tR)^n}{n!} \sum_{m=0}^n \binom{n}{m} (-1)^{n-m} \tilde{\mathbf{P}}^m = \sum_{m=0}^{\infty} \frac{(tR)^m}{m!} \tilde{\mathbf{P}}^m \sum_{n=m}^{\infty} \frac{(-tR)^{n-m}}{(n-m)!}$$
$$= e^{-tR} \sum_{m=0}^{\infty} \frac{(tR)^m}{m!} \tilde{\mathbf{P}}^m,$$

from which it is clear that $\mathbf{P}(t)$ is a transition probability matrix and therefore that $t \rightsquigarrow \mathbf{P}(t)$ is a transition probability function.

We begin by showing that if $\{X(t) : t \ge 0\}$ satisfies (5.2.1) and has initial distribution μ , then, for any $K \ge 0$, $0 = t_0 < t_1 < \cdots < t_K$, and $(j_0, \ldots, j_K) \in \mathbb{S}^{K+1}$,

$$\mathbb{P}(X(t_k) = j_k \text{ for } 0 \le k \le K) = \mu_{j_0} \prod_{1 \le k \le K} (\mathbf{P}(t_k - t_{k-1}))_{j_{k-1}j_k},$$
(5.2.4)

where the product over the empty set is taken to be 1. To prove this, we will work by induction on *K*. When K = 0 there is nothing to do. Thus, suppose that $K \ge 1$ and that (5.2.4) holds for K - 1. Set $A = \{X(t_k) = j_k, 0 \le k < K\}$ and define $u(t, j) = \mathbb{P}(\{X(t_{k-1} + t) = j\} \cap A)$ for $t \ge 0$ and $j \in \mathbb{S}$. By (5.2.1),

$$\lim_{h \searrow 0} \sup_{\substack{j \in \mathbb{S} \\ t \in [0,\infty)}} \frac{1}{h} \Big| u(t+h,j) - u(t,j) - h \big(\mathbf{u}(t) \mathbf{Q} \big)_j \Big| = 0,$$

where $\mathbf{u}(t)$ is the row vector determined by $u(t, \cdot)$. Therefore $t \rightsquigarrow u(t, j)$ is continuously differentiable and its derivative $\dot{u}(t, j)$ is equal to $(\mathbf{Qu}(t))_j$. In addition, $u(0, j) = \delta_{j,j_{K-1}} \mathbb{P}(A)$. In other words, we now know that $(t, j) \rightsquigarrow u(t, j)$ is a bounded solution to

$$\dot{u}(t,j) = \left(\mathbf{Qu}(t)\right)_{j} \quad \text{with } u(0,j) = \delta_{j,j_{K-1}} \mathbb{P}(A). \tag{*}$$

This information together with the following lemma will allow us to complete the inductive step.

Lemma 5.2.5 A bounded function $w : [0, \infty) \times \mathbb{S} \longrightarrow \mathbb{R}$ that is continuously differentiable with respect to $t \ge 0$ satisfies

$$\dot{w}(t,j) = \left(\mathbf{Q}\mathbf{w}(t)\right)_{i},$$

where $\mathbf{w}(t)$ is the row vector determined by $w(t, \cdot)$, if and only if $w(t, j) = (\mathbf{w}(0)\mathbf{P}(t))_j$ for $(t, j) \in [0, \infty) \times \mathbb{S}$.

Proof Let **w** be a row vector with $\|\mathbf{w}\|_{u} < \infty$, and define $w(t, j) = (\mathbf{wP}(t))_{j}$ for $(t, j) \in [0, \infty) \times \mathbb{S}$. Then $|w(t, j)| \le \|\mathbf{w}\|_{u}$ and

$$w(t+h, j) - w(t, j) = \left(\left(\mathbf{P}(h) - \mathbf{I} \right) \mathbf{w}(t) \right)_j,$$

where $\mathbf{w}(t)$ is the row vector determined by $w(t, \cdot)$. Next note that

$$\lim_{h\searrow 0} \left\| \frac{\mathbf{P}(h) - \mathbf{I}}{h} - \mathbf{Q} \right\|_{\mathbf{u},\mathbf{v}} = 0,$$

and from this conclude that $t \rightsquigarrow w(t, j)$ is continuously differentiable and that $\dot{w}(t, j) = (\mathbf{wQ})_j$. Thus the "if" assertion is proved.

To prove the "only if" assertion, let t > 0 and consider the function $f(s, j) = (\mathbf{w}(t-s)\mathbf{P}(s))_j$ for $s \in [0, t]$. Clearly $f(\cdot, j)$ is continuously differentiable and

$$\dot{f}(s, j) = -(\mathbf{w}(t-s)\mathbf{QP}(s))_j + (\mathbf{w}(t-s)\mathbf{QP}(s))_j = 0$$

for $s \in (0, t)$. Hence, $w(t, j) = f(0, j) = f(t, j) = (\mathbf{w}(0)\mathbf{P}(t))_j$.

Combining (*) with Lemma 5.2.5, we see that

$$u(t_K - t_{K-1}, j_K) = \left(\mathbf{P}(t_K - t_{K-1}) \right)_{j_{K-1}j_K} \mathbb{P}(A),$$

and therefore, by the induction hypothesis, that (5.2.4) holds for K.

We next show that if (5.2.4) holds, then $\{X(t) : t \ge 0\}$ is a Markov process with transition probability function $t \rightsquigarrow \mathbf{P}(t)$. That is, we will show that if $s \ge 0$, then for t > s and $j \in \mathbb{S}$,

$$\mathbb{P}\left(\left\{X(t)=j\right\} \mid X(\tau), \tau \in [0,s]\right) = \left(\mathbf{P}(t-s)\right)_{X(s)j}.$$
(5.2.6)

Equivalently, we have to show that if $A \in \sigma(\{X(\tau) : \tau \in [0, s]\})$, then

$$\mathbb{P}(\{X(t)=j\}\cap A) = \mathbb{E}[(\mathbf{P}(t-s))_{X(s)j}, A],$$

and we have to do so only when $A = \{X(t_k) = j_k, 0 \le k < K\}$ for some $K \ge 1$, $0 = t_0 < \cdots < t_{K-1} = s$, and $(j_0, \ldots, j_{K-1}) \in \mathbb{S}^K$. But in that case, the required equation is just (5.2.4) with $t_K = t$ and $j_K = j$.

In order to prove that there exists a stochastic process satisfying conditions (a) and (b), it will be useful to know that (5.2.4) implies (5.2.1). To see this, note that, (5.2.4) implies (5.2.6) and (5.2.6) implies

$$\mathbb{P}\left(X(t+h) = j \mid X(\tau), \tau \in [0, t]\right)$$
$$= \left(\mathbf{P}(h)\right)_{X(t)j}$$
$$= \delta_{X(t),j} + h(\mathbf{Q})_{X(t)j} + h^2 \sum_{n=2}^{\infty} \frac{h^{n-2}(\mathbf{Q}^n)_{X(t)j}}{n!}$$

Therefore, since

$$\delta_{X(t)j} + h(\mathbf{Q})_{X(t)j} = (1 - hR_{X(t)})\delta_{X(t),j} + hR_{X(t)}(\mathbf{P})_{X(t)j},$$

 \square

and

$$\left\|\sum_{n=2}^{\infty} \frac{h^{n-2} \mathbf{Q}^n}{n!}\right\|_{\mathbf{u},\mathbf{v}} \le e^{2R} \quad \text{for } h \in [0,1],$$

(5.2.1) follows from (5.2.4).

We can now prove the existence of a stochastic process that satisfies (a) and (b). To this end, recall the transition function $\tilde{\mathbf{P}}$ in (5.2.3), and let $\{X_n : n \ge 0\}$ be a Markov chain with transition probability $\tilde{\mathbf{P}}$ and initial distribution $\boldsymbol{\mu}$. Next, let $\{N(t) : t \ge 0\}$ be a simple Poisson process that is independent of $\{X_n : n \ge 0\}$, and set $X(t) = X_{N(Rt)}$. Clearly $t \rightsquigarrow X(t)$ is right continuous and piecewise constant. In addition, if $K \ge 0$, $0 = t_0 < \cdots < t_K$, and $(j_0, \ldots, j_K) \in \mathbb{S}^{K+1}$, then

$$\begin{split} \mathbb{P} \Big(X(t_k) &= j_k, 0 \le k \le K \Big) \\ &= \sum_{0=m_0 \le \dots \le m_K} \mathbb{P} \big(X_{m_k} = j_k \& N(t_k) = m_k \text{ for } 0 \le k \le K \big) \\ &= \sum_{0=m_0 \le \dots \le m_K} \mathbb{P} (X_{m_k} = j_k, 0 \le k \le K) \mathbb{P} \big(N(t_k) = m_k, 0 \le k \le K \big) \\ &= (\mu)_{j_0} \sum_{0=m_0 \le \dots \le m_K} \prod_{k=1}^K \big(\tilde{\mathbf{P}}^{m_k - m_{k-1}} \big)_{j_{k-1}k_k} \\ &\qquad \times \frac{e^{-(t_k - t_{k-1})R} (R(t_k - t_{k-1}))^{m_k - m_{k-1}}}{(m_k - m_{k-1})!} \\ &= (\mu)_{j_0} \prod_{k=1}^K \big(\mathbf{P}(t_k - t_{k-1}) \big)_{j_{k-1}j_k}. \end{split}$$

Hence, $\{X(t) : t \ge 0\}$ has initial distribution μ and satisfies (5.2.4) and therefore (5.2.1).

The following theorem summarizes our progress thus far.

Theorem 5.2.7 Let a bounded set $\Re = \{R_i : i \in \mathbb{S}\} \subseteq [0, \infty)$ and a transition probability matrix \mathbf{P} be given. Then, for each probability vector $\boldsymbol{\mu}$, there is a stochastic process $\{X(t) : t \ge 0\}$ satisfying the conditions (a) and (b). Furthermore, if $\mathbf{Q} = \mathbf{R}(\mathbf{P} - \mathbf{I})$, where $(\mathbf{R})_{ij} = \delta_{i,j}R_i$, and $\mathbf{P}(t)$ is given by (5.2.2), then $\{X(t) : t \ge 0\}$ is a Markov process with transition probability function $t \rightsquigarrow \mathbf{P}(t)$ in the sense that it satisfies (5.2.6). In particular, the distribution of $\{X(t) : t \ge 0\}$ is uniquely determined by (5.2.4).

5.2.2 An Alternative Construction

In this subsection we will give an alternative construction of processes that satisfy conditions (a) and (b), one that reveals other properties.

5.2 Markov Processes with Bounded Rates

Given a bounded set $\Re = \{R_i : i \in \mathbb{S}\} \subseteq [0, \infty)$, define

$$\Phi^{\mathfrak{R}}:[0,\infty)\times(0,\infty)^{\mathbb{Z}^{+}}\times\mathbb{S}^{\mathbb{N}}\longrightarrow\mathbb{S}$$

by the prescription

$$\Phi^{\Re}(t; (e_1, \dots, e_n, \dots), (j_0, \dots, j_n, \dots)) = j_n \quad \text{for } \xi_n \le t < \xi_{n+1}$$

where $\xi_0 = 0$ and $\xi_n = \sum_{m=1}^n R_{j_{m-1}}^{-1} e_m$ when $n \ge 1$, (5.2.8)

with $\xi_n = \infty$ if $R_{j_{m-1}} = 0$ for some $1 \le m \le n$. Next, let $\{\tilde{E}_n : n \ge 1\}$ be a sequence of mutually independent random, unit exponential random variables, and without loss in generality, assume that, for each ω , $\tilde{E}_n(\omega) > 0$ for all *n*. Because,

$$\mathbb{E}\left[\exp\left(-\sum_{n=1}^{\infty}\tilde{E}_n\right)\right] = \lim_{N \to \infty} \left(\int_0^{\infty} e^{-2t} dt\right)^N = 0,$$

we may and will also assume that $\sum_{n=1}^{\infty} \tilde{E}_n(\omega) = \infty$ for all ω . Now let $\{\tilde{X}_n : n \ge 0\}$ be a Markov chain with transition probability **P** and initial distribution μ , assume that $\sigma(\{\tilde{X}_n : n \ge 0\})$ and $\sigma(\{\tilde{E}_n : n \ge 1\})$ are independent, and set

$$\tilde{J}_0 = 0$$
 and $\tilde{J}_n = \sum_{m=1}^n R_{\tilde{X}_{m-1}}^{-1} \tilde{E}_m$ for $n \ge 1$,

again with the understanding that $\tilde{J}_n = \infty$ if $R_{\tilde{X}_{m-1}} = 0$ for some $1 \le m \le n$. Because the R_i 's are bounded and $\sum_{m=1}^{\infty} \tilde{E}_m \equiv \infty$, $\tilde{J}_n \nearrow \infty$ as $n \to \infty$. Hence, we can define a right continuous, piecewise constant process $\{\tilde{X}(t): t \ge 0\}$ by

$$\tilde{X}(t) = \boldsymbol{\Phi}^{\mathfrak{R}}(t; (\tilde{E}_1, \dots, \tilde{E}_n, \dots), (\tilde{X}_0, \dots, \tilde{X}_n, \dots)).$$
(5.2.9)

Theorem 5.2.10 *Refer to the preceding. The process* $\{\tilde{X}(t) : t \ge 0\}$ *given by* (5.2.9) *satisfies conditions* (a) *and* (b) *at the beginning of Sect.* 3.2.1.

Proof Since it is obvious that (a) is satisfied, what remains to be shown is that $\{\tilde{X}(t) : t \ge 0\}$ satisfies (5.2.1). We begin by observing that

$$\Phi^{\mathfrak{R}}(t+h; (e_1, \dots, e_n, \dots), (j_0, \dots, j_n, \dots))$$

= $\Phi^{\mathfrak{R}}(h; (e_{m+1} - R_{j_m}(t-\xi_m), e_{m+2} \dots, e_{m+n}, \dots),$
 $\times (j_m, \dots, j_{m+n}, \dots)) \quad \text{for } \xi_m \le t < \xi_{m+1}.$ (5.2.11)

Now, let $\tilde{A} \in \sigma({\tilde{X}(\tau) : \tau \in [0, t]})$ be given, and assume that $\tilde{X}(t) = i$ on A. Set $\tilde{A}_m \equiv \tilde{A} \cap {\tilde{J}_m \le s < \tilde{J}_{m+1}}$ for $m \ge 0$. Then

$$\tilde{A}_m = \left\{ \tilde{E}_{m+1} > R_i(t - \tilde{J}_m) \right\} \cap \tilde{B}_m,$$

where $\{\tilde{J}_m \leq s\} \supseteq \tilde{B}_m \in \sigma(\{\tilde{E}_1, \ldots, \tilde{E}_m\} \cup \{\tilde{X}_0, \ldots, \tilde{X}_m\})$. In addition,

$$\mathbb{P}\left(\left\{\tilde{X}(t+j)=j\right\}\cap\tilde{A}\right) = \sum_{m=0}^{\infty} \mathbb{P}\left(\left\{\tilde{X}(t+h)=j\right\}\cap\tilde{A}_{m}\right)$$
$$= \sum_{m=0}^{\infty} \mathbb{P}\left(\left\{\tilde{X}(t+h)=j \& \tilde{E}_{m+1} > R_{i}(t-\tilde{J}_{m})\right\}\cap\tilde{B}_{m}\right),$$

and, by (5.2.9) and (5.2.11),

$$\mathbb{P}\left(\left\{\tilde{X}(t+h)=j \& \tilde{E}_{m+1} > R_i(t-\tilde{J}_m)\right\} \cap \tilde{B}_m\right)$$

= $\mathbb{P}\left(\left\{\Phi^{\mathfrak{R}}\left(h; \left(\tilde{E}_{m+1}-R_i(t-\tilde{J}_m), \tilde{E}_{m+2}, \dots, \tilde{E}_{m+n}, \dots\right), (i, \tilde{X}_{m+1}, \dots, \tilde{X}_{m+n}, \dots)\right)=j\right\} \cap \left\{\tilde{E}_{m+1} > R_i(t-\tilde{J}_m)\right\} \cap \tilde{B}_m\right).$

Since \tilde{E}_{m+1} is independent of all the other quantities in the last expression, we can apply (5.1.3) to see that it is equal to

$$\mathbb{P}\big(\tilde{X}(h)=j \mid \tilde{X}_0=i\big)\mathbb{E}\big[e^{-R_i(t-\tilde{J}_m)}, \tilde{B}_m\big]=\mathbb{P}\big(\tilde{X}(h)=j \mid \tilde{X}_0=i\big)\mathbb{P}(\tilde{A}_m)$$

Hence, we will be done once we show that

$$\lim_{h \searrow 0} \sup_{i, j \in \mathbb{S}} \frac{1}{h} \left| \mathbb{P} \left(\tilde{X}(h) = j \, | \, \tilde{X}_0 = i \right) - (1 - hR_i) \delta_{i, j} - hR_i(\mathbf{P})_{ij} \right| = 0. \tag{*}$$

Note that $\mathbb{P}(\tilde{X}(h) = j | \tilde{X}_0 = i)$ equals

$$\mathbb{P}(\tilde{X}_0 = j \& \tilde{J}_1 \le h \,|\, \tilde{X}_0 = i) + \mathbb{P}(\tilde{X}_1 = j \& \tilde{J}_1 \le h < \tilde{J}_2 \,|\, \tilde{X}_0 = i) \\ + \mathbb{P}(\tilde{X}(h) = j \& \tilde{J}_2 \le h \,|\, \tilde{X}_0 = i).$$

Obviously

$$\mathbb{P}(\tilde{X}_0 = j \& \tilde{J}_1 \le h \mid \tilde{X}_0 = i) = e^{-hR_i} \delta_{i,j}$$

and

$$\mathbb{P}(\tilde{X}_{1} = j \& \tilde{J}_{1} \le h < \tilde{J}_{2} | \tilde{X}_{0} = i) = \mathbb{P}(R_{i}^{-1}\tilde{E}_{1} \le h < R_{i}^{-1}\tilde{E}_{1} + R_{j}^{-1}\tilde{E}_{2})(\mathbf{P})_{ij}$$
$$= (1 - e^{-hR_{i}})(\mathbf{P})_{ij} - \mathbb{P}(R_{i}^{-1}\tilde{E}_{1} + R_{j}^{-1}\tilde{E}_{2} \le h)(\mathbf{P})_{ij}.$$

Hence

$$\left|\mathbb{P}\left(\tilde{X}(h)=j\,|\,\tilde{X}_{0}=i\right)-e^{-hR_{i}}\delta_{i,j}-\left(1-e^{-hR_{i}}\right)(\mathbf{P})_{ij}\right|\leq 2\mathbb{P}(\tilde{E}_{1}+\tilde{E}_{2}\leq Rh),$$

where $R = \sup_{i \in \mathbb{S}} R_i$. Finally,

$$\mathbb{P}(\tilde{E}_1 + \tilde{E}_2 \le Rh) = \int_0^{Rh} \tau e^{-\tau} d\tau \le (Rh)^2,$$

and so (*) follows and the proof of theorem is complete.

5.2.3 Distribution of Jumps and Jump Times

As a consequence of Theorem 5.2.7, we see that the distribution of $\{X(t) : t \ge 0\}$ is determined by its initial distribution μ and the matrix $\mathbf{Q} = \mathbf{R}(\mathbf{P} - \mathbf{I})$. For this reason, we say that a stochastic process satisfying conditions (a) and (b) is a Markov process generated \mathbf{Q} with initial distribution μ .

The distinguishing characteristics of \mathbf{Q} are that its off-diagonal entries are nonnegative and each of its rows sums to 0. For no better reason than that Q follows P in the alphabet, probabilists call such a matrix a Q-matrix. Given a Q-matrix \mathbf{Q} , there are many diagonal matrices \mathbf{R} with non-negative entries and transition matrices \mathbf{P} for which $\mathbf{Q} = \mathbf{R}(\mathbf{P} - \mathbf{I})$. Indeed, we already took advantage of this fact when we wrote $\mathbf{R}(\mathbf{P} - \mathbf{I})$ as (cf. (5.2.3)) $R(\tilde{\mathbf{P}} - \mathbf{I})$ with $R = \sup_{i \in \mathbb{S}} R_i$. Nonetheless, it is possible to remove this ambiguity and make a canonical choice of \mathbf{R} and \mathbf{P} . For reasons that will become clear shortly, we will make this canonical choice by taking $(\mathbf{R})_{ij} = \delta_{i,j} R_i$ and \mathbf{P} where

$$R_{i} = -(\mathbf{Q})_{ii}, \qquad (\mathbf{P})_{ii} = \mathbf{1}_{\{0\}}(R_{i}), \quad \text{and},$$

for $j \neq i$, $(\mathbf{P})_{ij} = \begin{cases} 0 & \text{if } R_{i} = 0 \\ R_{i}^{-1}(\mathbf{Q})_{ij} & \text{if } R_{i} > 0. \end{cases}$ (5.2.12)

When $\mathfrak{R} = \{R_i : i \in \mathbb{S}\}$ and **P** are those given by (5.2.12), we will call \mathfrak{R} the *rates determined by* **Q** and **P** the *transition probability matrix determined by* **Q**.

Our goal in this subsection is to prove the following theorem that describes how $\{X(t) : t \ge 0\}$ jumps when it jumps.

Theorem 5.2.13 Let $\{X(t) : t \ge 0\}$ be a Markov process generated by **Q** and let \Re and **P** be given by (5.2.12). Set $J_0 = 0$,

$$J_{n} = \inf\{t > J_{n-1} : X(t) \neq X(J_{n-1})\} \text{ for } n \ge 1$$
$$X_{n} = \begin{cases} X(J_{n}) & \text{if } J_{n} < \infty \\ X_{n-1} & \text{if } n \ge 1 \text{ and } J_{n} = \infty, \end{cases}$$
$$E_{n} = \begin{cases} R_{X_{n-1}}(J_{n} - J_{n-1}) & \text{if } J_{n} < \infty \\ 0 & \text{if } J_{n} = \infty, \end{cases}$$

and $\zeta = \inf\{n \ge 0 : R_{X_n} = 0\}$. Then $\{X_n : n \ge 0\}$ is a Markov chain with transition matrix **P** and, for any $K \ge 1$ and $\{t_k : 1 \le k \le K\} \subseteq (0, \infty)$,

$$\mathbb{P}(\{E_k > t_k, 1 \le k \le K\} \cap A) = e^{-\sum_{k=1}^{K} t_k} \mathbb{P}(A)$$

for any $A \in \sigma(\{X_n : n \ge 0\})$ contained in $\{\zeta > K\}$.

Proof Refer to Sect. 5.2.2 and let $\{\tilde{X}(t) : t \ge 0\}$ be the process in Theorem 5.2.10 when μ is the distribution of X(0). Because it has the same distribution as $\{X(t) : t \ge 0\}$, it suffices to prove that $\{\tilde{X}(t) : t \ge 0\}$ has the asserted properties.

Obviously, $\{J_n : n \ge 0\}$ and $\{X_n : n \ge 0\}$ are is related to $\{\overline{X}(t) : t \ge 0\}$ in the same way as $\{J_n : n \ge 0\}$ and $\{X_n : n \ge 0\}$ are related to $\{\overline{X}(t) : t \ge 0\}$. Furthermore, if $\tilde{\zeta} = \inf\{n \ge 0 : R_{\tilde{X}_{n-1}} = 0\}$, then $\tilde{\zeta}$ is related to $\{\overline{X}(t) : t \ge 0\}$ the same way that ζ is related to $\{X(t) : t \ge 0\}$. Finally, $\tilde{E}_k = R_{\tilde{X}_{k-1}}(\tilde{J}_k - \tilde{J}_{k-1})$ if $\tilde{\zeta} > k$. Now suppose that $\tilde{A} \in \sigma(\{\tilde{X}_n : n \ge 0\})$ is contained in $\{\tilde{\zeta} > K\}$. Then, because the \tilde{E}_k 's are independent of \tilde{A} ,

$$\mathbb{P}\left(\{\tilde{E}_k > t_k, 1 \le k \le K\} \cap \tilde{A}\right) = \mathbb{P}(\tilde{E}_k > t_k, 1 \le k \le K)\mathbb{P}(\tilde{A})$$
$$= e^{-\sum_{k=1}^{K} t_k}\mathbb{P}(\tilde{A}).$$

Remark It is important to appreciate what Theorem 5.2.13 is saying. Namely, if \mathfrak{R} and \mathbf{P} are the canonical rates and transition probability matrix determined by the Q-matrix \mathbf{Q} and if $\{X(t) : t \ge 0\}$ is a Markov process generated by \mathbf{Q} , then the distribution of $\{X(J_n \land \zeta) : n \ge 0\}$, where the J_n 's are the successive jump times of the process, has the same distribution as that of a Markov chain with transition matrix \mathbf{P} stopped at time ζ . Furthermore, conditional on $\{X(J_n \land \zeta) : n \ge 0\}$, the holding times $\{J_n - J_{n-1} : n \ge 1\}$ are mutually independent exponential randoms, $J_n - J_{n-1}$ being infinite if $R_{X(J_{n-1})} = 0$, and $J_n - J_{n-1}$ being a exponential random variable with mean value $R_{X(J_{n-1})}^{-1}$ if $R_{X(J_{n-1})} > 0$.

5.2.4 Kolmogorov's Forward and Backward Equations

As we saw at the end of the Sect. 3.2.1, apart from its initial distribution, the distribution of a Markov process is completely determined by the transition probability function $t \rightsquigarrow \mathbf{P}(t)$ in (5.2.2). Although the expression in (5.2.2) is useful from a theoretical perspective, for many applications it is not practical. For this reason it is desirable to have other ways of doing computations involving $\mathbf{P}(t)$, and we will develop some in this subsection.

It should be clear both from (5.1.9) and (5.2.2) that $\mathbf{P}(t)$ is some sort of exponential function. Indeed, (5.2.2) makes it reasonable to think that $\mathbf{P}(t)$ ought to be written as $e^{t\mathbf{Q}}$, and that is how a functional analyst would write it. Taking a hint

from the usual exponential function for real numbers, one should guess that another way to describe $\mathbf{P}(t)$ is as the solution to the equation

$$\dot{\mathbf{P}}(t) \equiv \frac{d}{dt}\mathbf{P}(t) = \mathbf{P}(t)\mathbf{Q}$$
 with $\mathbf{P}(0) = \mathbf{I}$. (5.2.14)

In fact, we made implicit use of (5.2.14), which is called *Kolmogorov's forward equation*, in Lemma 5.2.5. The term "forward" can be understood by writing (5.2.14) in coordinates:

$$(\dot{\mathbf{P}}(t))_{ij} = \sum_{k \in \mathbb{S}} (\mathbf{P}(t))_{ik} (\mathbf{Q})_{kj}.$$

Thinking in terms of someone traveling along the paths of the process, in this expression *i* is the *backward variable* since it is the variable that specifies where the traveler was at time 0 and is therefore the one the traveler sees when he looks backward in time. Similarly, *j* is the *forward variable* because it is the one that says where he will be at time *t*. Thus (5.2.14) is an equation that gives the evolution of $t \rightsquigarrow \mathbf{P}(t)$ in terms of its forward variable when the backward variable is fixed. Of course, it is equally reasonable to expect that

$$\dot{\mathbf{P}}(t) = \mathbf{Q}\mathbf{P}(t) \quad \text{with } \mathbf{P}(0) = \mathbf{I}, \tag{5.2.15}$$

which is known as *Kolmogorov's backward equation* since its coordinate version expresses the evolution of $t \rightsquigarrow \mathbf{P}(t)$ as a function of its backward variable.

The proof that $t \rightsquigarrow \mathbf{P}(t)$ solves (5.2.14) and (5.2.15) is easy. By (5.1.9),

$$\mathbf{P}(t+h) - \mathbf{P}(t) = \mathbf{P}(t) \big(\mathbf{P}(h) - \mathbf{I} \big) = \big(\mathbf{P}(h) - \mathbf{I} \big) \mathbf{P}(t),$$

and so, for $h \in (0, 1]$,

$$\begin{split} \left\| \frac{\mathbf{P}(t+h) - \mathbf{P}(t)}{h} - \mathbf{P}(t) \mathbf{Q} \right\|_{\mathbf{u},\mathbf{v}} & \vee \left\| \mathbf{Q} \mathbf{P}(t) - \frac{\mathbf{P}(t+h) - \mathbf{P}(t)}{h} \right\|_{\mathbf{u},\mathbf{v}} \\ & \leq \left\| \sum_{n=2}^{\infty} \frac{h^{n-1} \mathbf{Q}^n}{n!} \right\|_{\mathbf{u},\mathbf{v}} \leq h e^{2R}, \end{split}$$

where $R = \sup_{i \in \mathbb{S}} (\mathbf{Q})_{ii}$. In addition, one can show that $t \rightsquigarrow \mathbf{P}(t)$ is uniquely determined by either (5.2.14) or (5.2.15). In fact, suppose that $t \in [0, \infty) \longmapsto \mathbf{M}(t) \in M_{u,v}(\mathbb{S})$ is a continuously differentiable and

$$\dot{\mathbf{M}}(t) = \mathbf{M}(t)\mathbf{Q}$$
 or $\dot{\mathbf{M}}(t) = \mathbf{Q}\mathbf{M}(t)$.

Then $\mathbf{M}(t) = \mathbf{0}$ for $t \ge 0$ if $\mathbf{M}(0) = \mathbf{0}$. To see this, note that in either case,

$$\left\|\mathbf{M}(t)\right\|_{\mathbf{u},\mathbf{v}} \le \|\mathbf{Q}\|_{\mathbf{u},\mathbf{v}} \int_0^t \left\|\mathbf{M}(\tau)\right\|_{\mathbf{u},\mathbf{v}} \quad \text{for all } t \ge 0$$

and apply *Gronwall's inequality* (cf. Exercise 5.5.5). Applying this to the difference of two solutions of either (5.2.14) or of (5.2.15) one sees that the difference vanishes.

In applications, one often uses these ideas to compute $t \rightsquigarrow \mu \mathbf{P}(t)$ or $t \rightsquigarrow \mathbf{P}(t)\mathbf{f}$ for probability vectors $\boldsymbol{\mu}$ or bounded column vectors \mathbf{f} . That is, (5.2.14) and (5.2.15) say that

$$\frac{d}{dt}\boldsymbol{\mu}\mathbf{P}(t) = \boldsymbol{\mu}\mathbf{P}(t)\mathbf{Q} \quad \text{with } \boldsymbol{\mu}\mathbf{P}(0) = \boldsymbol{\mu}$$

$$\frac{d}{dt}\mathbf{P}(t)\mathbf{f} = \mathbf{Q}\mathbf{P}(t)\mathbf{f} \quad \text{with } \mathbf{P}(0)\mathbf{f} = \mathbf{f},$$
(5.2.16)

and, just as in the preceding, one can use Gronwall's inequality to check that $t \rightsquigarrow \boldsymbol{\mu} \mathbf{P}(t)$ and $t \rightsquigarrow \mathbf{P}(t) \mathbf{f}$ are uniquely determined by these equations.

5.3 Unbounded Rates

Thus far we have been assuming that the rates $\Re = \{R_i : i \in \mathbb{S}\}\$ are bounded. This assumption simplified our considerations in several ways. For one thing, it allowed us to use (5.2.2) to define and derive properties of $t \rightsquigarrow \mathbf{P}(t)$. Secondly, and more important, it allowed us to show that there exists a process that satisfies the conditions (a) and (b) at the beginning of Sect. 5.2.1. The easiest way to understand what might prevent one from constructing such a process is to look at the construction given in Sect. 5.2.2. No matter what the rates are, the process (5.2.9) can be used to define $\tilde{X}(t)$ as long as $t < \sum_{m=1}^{\infty} R_{\tilde{X}_{m-1}}^{-1} \tilde{E}_m$. When the R_i 's are bounded, and therefore $\sum_{m=1}^{\infty} R_{\tilde{X}_{m-1}}^{-1} \tilde{E}_m$ may be finite with positive probability, with the result that $\tilde{X}(t)$ will not be defined for all t. In this section I will give conditions, other than boundedness, that enable one to overcome this problem.

5.3.1 Explosion

The setting here is the same as the one in Sect. 5.2.2, only now we are no longer assuming that the rates are bounded. Even so, as I just said, $\tilde{X}(t)$ is well-defined by the prescription in (5.2.9) as long as $0 \le t < \tilde{J}_{\infty} \equiv \sum_{m=1}^{\infty} R_{\tilde{X}_{m-1}}^{-1} \tilde{E}_m$.

Our first step is to show that, with probability 1, J_{∞} coincides with the time when the process explodes out of the state space. To be precise, choose an *exhaustion* $\{F_N : N \ge 1\}$ of \mathbb{S} by non-empty, finite subsets. That is, for each N, F_N is a nonempty, finite subset of \mathbb{S} , $F_N \subseteq F_{N+1}$, and $\mathbb{S} = \bigcup_N F_N$. Next, take $\mathfrak{R}^{(N)}$ to be the set of rates given by

$$R_i^{(N)} = \begin{cases} R_i & \text{if } i \in F_N \\ 0 & \text{if } i \notin F_N. \end{cases}$$

Then, for each $N \ge 1$, (5.2.9) with $\mathfrak{R}^{(N)}$ replacing \mathfrak{R} determines a Markov process $\{\tilde{X}^{(N)}(t) : t \ge 0\}$. Moreover, if $\tilde{\zeta}_N \equiv \inf\{t \ge 0 : \tilde{X}^{(N)}(t) \notin F_N\}$, then

$$t \wedge \tilde{\zeta}_N < \tilde{J}_\infty$$
 and $\tilde{X}(t \wedge \tilde{\zeta}_N) = \tilde{X}^{(N)}(t \wedge \tilde{\zeta}_N)$ for all $t \ge 0.$ (5.3.1)

To see this, choose $n \ge 0$ so that

$$\tilde{J}_n^{(N)} \le t \wedge \tilde{\zeta}_N < \tilde{J}_{n+1}^{(N)},$$

where $\tilde{J}_0^{(N)} = 0$ and

$$\tilde{J}_n^{(N)} = \sum_{m=1}^n \frac{\tilde{E}_m}{R_{\tilde{X}_{m-1}}^{(N)}} \text{ for } n \ge 1.$$

Because $\tilde{X}_{m-1} \in F_N$ for $1 \le m \le n+1$, $\tilde{J}_{\ell}^{(N)} = \tilde{J}_{\ell} < \tilde{J}_{\infty}$ for all $0 \le \ell \le n+1$, and so $t \land \tilde{\zeta}_N < J_{\infty}$ and $\tilde{J}_n \le t \land \tilde{\zeta}_N < \tilde{J}_{n+1}$. At the same time, it also means that

$$\tilde{X}(t \wedge \tilde{\zeta}_N) = \Phi^{\mathfrak{R}} \left(t \wedge \tilde{\zeta}_N; (\tilde{E}_1, \dots, \tilde{E}_n, \dots), (\tilde{X}_0, \dots, \tilde{X}_n, \dots) \right)$$
$$= \Phi^{\mathfrak{R}^{(N)}} \left(t \wedge \tilde{\zeta}_N; (\tilde{E}_1, \dots, \tilde{E}_n, \dots), (\tilde{X}_0, \dots, \tilde{X}_n, \dots) \right) = \tilde{X}^{(N)}(t \wedge \tilde{\zeta}_N).$$

As a consequence of (5.3.1), we know that $\tilde{X}^{(N+1)}(t \wedge \tilde{\zeta}_N) = X^{(N)}(t \wedge \tilde{\zeta}_N)$ for all $t \ge 0$ and therefore that $\tilde{\zeta}_N \le \tilde{\zeta}_{N+1}$. Hence the *explosion time* $\tilde{\mathfrak{e}} \equiv \lim_{N \to \infty} \tilde{\zeta}_N$ exists (in $[0, \infty]$).

Lemma 5.3.2 $\mathbb{P}(\tilde{\mathfrak{e}} = \tilde{J}_{\infty}) = 1$ and so, with probability 1, $\tilde{X}(t)$ is well-defined for $t \in [0, \tilde{\mathfrak{e}})$.

Proof Because $\{\tilde{e} \neq \tilde{J}_{\infty}\}$ can be written as the union of the sets

$$\{\tilde{\mathfrak{e}} > T \geq \tilde{J}_{\infty}\} \cup \{\tilde{J}_{\infty} > T \geq \tilde{\mathfrak{e}}\}$$

as T runs over the positive rational numbers, in order to prove that $\tilde{\mathfrak{e}} = \tilde{J}_{\infty}$ with probability 1, (7.1.5) says that it suffices for us to show that

$$\mathbb{P}(\tilde{\mathfrak{e}} > T \ge \tilde{J}_{\infty}) = 0 = \mathbb{P}(\tilde{J}_{\infty} > T \ge \tilde{\mathfrak{e}}) \quad \text{for each } T > 0.$$

To this end, first suppose that $\mathbb{P}(\tilde{\mathfrak{e}} > T \ge \tilde{J}_{\infty}) > 0$ for some *T*. Then there exists an *N* such that $\mathbb{P}(\tilde{\zeta}_N > T \ge \tilde{J}_{\infty}) > 0$. On the other hand,

$$ilde{\zeta}_N > T \geq ilde{J}_\infty \quad \Longrightarrow \quad ilde{J}_\infty \geq rac{1}{R^{(N)}} \sum_1^\infty ilde{E}_m,$$

where $R^{(N)} = \sup_{i \in \mathbb{S}} R_j^{(N)}$, and therefore we are led to the contradiction

$$0 < \mathbb{P}(\tilde{\zeta}_N > T \ge \tilde{J}_{\infty}) \le \mathbb{P}\left(\sum_{m=1}^{\infty} E_m \le R^{(N)}T\right) = 0.$$

Next suppose that $\mathbb{P}(\tilde{J}_{\infty} > T \ge \tilde{\mathfrak{e}}) > 0$. Then there exists an $n \ge 1$ such that $\mathbb{P}(\tilde{J}_n > T \ge \tilde{\mathfrak{e}}) > 0$. On the other hand, as $N \to \infty$,

$$\mathbb{P}(\tilde{J}_n > T \ge \tilde{\mathfrak{e}}) \le \mathbb{P}(\tilde{J}_n > T \ge \tilde{\zeta}_N)$$
$$\le \mathbb{P}(\exists 0 \le m \le n \tilde{X}_m \notin F_N) \longrightarrow 0.$$

That is, we know that $\mathbb{P}(\tilde{J}_n > T \ge \tilde{\mathfrak{e}}) = 0.$

Lemma 5.3.3 If $\mathbb{P}(\tilde{\mathfrak{e}} = \infty | \tilde{X}(0) = i) = 1$ for all $i \in \mathbb{S}$ and

$$(\mathbf{P}(t))_{ij} \equiv \mathbb{P}(\tilde{X}(t) = j | \tilde{X}(0) = i),$$

then, for each initial distribution μ and T > 0,

$$\lim_{N \to \infty} \sup_{t \in (0,T]} \| \boldsymbol{\mu} \mathbf{P}^{(N)}(t) - \boldsymbol{\mu} \mathbf{P}(t) \|_{v} = 0,$$
(5.3.4)

where $t \rightsquigarrow \mathbf{P}^{(N)}(t)$ is the transition probability function in (5.2.2) when the *Q*-matrix there is given by $(\mathbf{Q}^{(N)})_{ij} = \mathbf{1}_{F_N}(i)R_i(\mathbf{P} - \mathbf{I})_{ij}$. Moreover, $t \rightsquigarrow \mathbf{P}(t)$ satisfies (5.1.9) and $\{\tilde{X}(t) : t \ge 0\}$ is a Markov process for which it is the transition probability function. Finally, $t \rightsquigarrow \mathbf{P}(t)$ satisfies Kolmogorov's backward equation in the sense that, for each $(i, j) \in \mathbb{S}^2$,

$$\left(\mathbf{P}(t)\right)_{ij} = \delta_{i,j} + \int_0^t \left(\mathbf{Q}\mathbf{P}(\tau)\right)_{ij} d\tau, \qquad (5.3.5)$$

where $\mathbf{Q} = \mathbf{R}(\mathbf{P} - \mathbf{I})$ with $(\mathbf{R})_{ij} = R_i \delta_{i,j}$.

Proof First note that

$$\mathbb{P}(\mathfrak{e}=\infty) = \sum_{i \in \mathbb{S}} (\boldsymbol{\mu})_i \mathbb{P}\big(\mathfrak{e}=\infty \,|\, X(0)=i\big) = 1,$$

and therefore that $\lim_{N\to\infty} \mathbb{P}(\zeta_N \leq T) = 0$ for each $T \in [0, \infty)$. At the same time, by Theorem 5.3.2,

$$\left|\left(\boldsymbol{\mu}\mathbf{P}^{(N)}(t)\right)_{j}-\left(\boldsymbol{\mu}\mathbf{P}(t)\right)_{j}\right|\leq\mathbb{P}\left(\tilde{X}(t)=j\ \&\ \zeta_{N}\leq T\right)$$

for all $0 \le t \le T$ and $j \in \mathbb{S}$. Thus,

$$\sup_{t \in (0,T]} \left\| \boldsymbol{\mu} \mathbf{P}^{(N)}(t) - \boldsymbol{\mu} \mathbf{P}(t) \right\|_{v} \le \mathbb{P}(\zeta_{N} \le T) \longrightarrow 0 \quad \text{as } N \to \infty.$$

Given (5.3.4), we have that

$$\left(\mathbf{P}(s+t)\right)_{ij} = \lim_{N \to \infty} \left(\mathbf{P}^{(N)}(s+t)\right)_{ij} = \lim_{N \to \infty} \left(\mathbf{P}^{(N)}(s)\mathbf{P}^{(N)}(t)\right)_{ij}$$

for each $(i, j) \in \mathbb{S}^2$. In addition, by applying (5.3.4) with $\mu = (\mathbf{P}(s))_i$, we know that

$$\left(\mathbf{P}(s)\mathbf{P}(t)\right)_{ij} = \lim_{N \to \infty} \left(\mathbf{P}(s)\mathbf{P}^{(N)}(t)\right)_{ij}$$

for each $(i, j) \in \mathbb{S}^2$. Finally, again by (5.3.4), for each $M \ge 1$

$$\sum_{k\in F_M} \left(\mathbf{P}(s) \right)_{ik} = \lim_{N \to \infty} \sum_{k\in F_M} \left(\mathbf{P}^{(N)}(s) \right)_{ik},$$

and therefore, since

$$\sum_{k \notin F_M} \left(\mathbf{P}(s) \right)_{ik} = 1 - \sum_{k \in F_M} \left(\mathbf{P}(s) \right)_{ik},$$

we know that

$$\sum_{k \notin F_M} \left(\mathbf{P}(s) \right)_{ik} = \lim_{N \to \infty} \sum_{k \notin F_M} \left(\mathbf{P}^{(N)}(s) \right)_{ik}$$

Now let $\epsilon > 0$ be given, and use the preceding to choose an M such that

$$\sum_{k \notin F_M} \left(\mathbf{P}^{(N)}(s) \right)_{ik} \le \epsilon \quad \text{for all } N \ge M.$$

Then, for $N \ge M$,

$$\left| \left(\mathbf{P}(s)\mathbf{P}(t) \right)_{ij} - \left(\mathbf{P}(s)\mathbf{P}^{(N)}(t) \right)_{ij} \right| \leq \sum_{k \in F_M} \left| \left(\mathbf{P}(s) \right)_{ik} - \left(\mathbf{P}^{(N)}(s) \right)_{ik} \right| + \epsilon.$$

Since, by (5.3.4), the first term on the right tends to 0 as $N \to \infty$, we have now shown that $\mathbf{P}(t + s) = \mathbf{P}(s)\mathbf{P}(t)$.

To prove the Markov property, let $A \in \sigma({\tilde{X}(\tau) : \tau \in [0, s]})$ be given, and assume that $\tilde{X}(s) = i$ on A. Then, since, by (5.3.1),

$$A \cap \{\tilde{\zeta}_N > s\} \in \sigma\left(\left\{\tilde{X}^{(N)}(\tau) : \tau \in [0, s]\right\}\right),$$

the Markov property for $\{\tilde{X}^{(N)}(t) : t \ge 0\}$ plus (5.3.1) and the result in Lemma 5.3.2 justify

$$\mathbb{P}\left(\left\{\tilde{X}(s+t)=j\right\}\cap A\right) = \lim_{N\to\infty} \mathbb{P}\left(\left\{\tilde{X}^{(N)}(s+t)=j\right\}\cap A\cap\{\tilde{\zeta}_N>s\}\right)$$
$$= \lim_{N\to\infty} \left(\mathbf{P}^{(N)}(t)\right)_{ij} \mathbb{P}\left(A\cap\{\zeta_N>s\}\right) = \left(\mathbf{P}(t)\right)_{ij} \mathbb{P}(A).$$

Finally, to check (5.3.5), note that, by (5.2.15) applied to $\{\mathbf{P}^{(N)}(t) : t \ge 0\}$,

$$\left(\mathbf{P}^{(N)}(t)\right)_{ij} = \delta_{i,j} + \int_0^t \left(\mathbf{Q}\mathbf{P}^{(N)}(\tau)\right)_{ij} d\tau$$

as soon as N is large enough that $i \in F_N$. Hence, since

$$1 \ge \left(\mathbf{P}^{(N)}(\tau)\right)_{kj} \longrightarrow \left(\mathbf{P}(\tau)\right)_{kj}$$

while $\sum_{k} |(\mathbf{Q})_{ik}| = 2R_i < \infty$, (5.3.5) follows by Lebesgue's dominated convergence theorem.

Theorem 5.3.6 Let \mathbf{Q} be a Q-matrix and, for $N \ge 1$, determine the Q-matrix $\mathbf{Q}^{(N)}$ by $(\mathbf{Q}^{(N)})_{ij} = \mathbf{1}_{F_N}(i)(\mathbf{Q})_{ij}$. Given a point Δ not in \mathbb{S} and an $i \in \mathbb{S}$, there is a $\mathbb{S} \cup \{\Delta\}$ -valued, right continuous, piecewise constant stochastic process $\{X(t) : t \ge 0\}$ with the properties that

- (1) $\mathbb{P}(X(0) = i) = 1$,
- (2) for each $N \ge 1$ with $i \in F_N$, $\{X(t \land \zeta^{F_N}) : t \ge 0\}$, where $\zeta^{F_N} = \inf\{t \ge 0 : X(t) \notin F_N\}$, is a Markov process generated by $\mathbf{Q}^{(N)}$,
- (3) if $\mathfrak{e} \equiv \lim_{N \to \infty} \zeta^{F_N}$, then $\mathbb{P}(X(t) = \Delta \text{ for } t \in [\mathfrak{e}, \infty)) = 1$.

Moreover, the distribution of $\{X(t) : t \ge 0\}$ is uniquely determined by these properties. Finally, if $\mathbb{P}(\mathfrak{e} = \infty) = 1$, then $\{X(t) : t \ge 0\}$ is a Markov process with the transition probability function $t \rightsquigarrow \mathbf{P}(t)$ described in Lemma 5.3.3.

Proof Let \Re and \mathbf{P} be the canonical rates and transition probability determined by \mathbf{Q} , and define $\{\tilde{X}(t) : t \in [0, \tilde{J}_{\infty})\}$ accordingly for \Re and \mathbf{P} with $\tilde{X}(0) \equiv i$. Next, set $X(t) = \tilde{X}(t)$ if $t \in [0, \tilde{J}_{\infty} \land \tilde{\mathfrak{e}})$ and $X(t) = \Delta$ if $t \in [\tilde{J}_{\infty} \land \tilde{e}, \infty)$. Obviously, $\mathbb{P}(X(0) = i) = 1, \zeta^{F_N} = \tilde{\zeta}_N$ for all $N \ge 1$, and, by Lemma 5.3.2, $\mathbb{P}(\mathfrak{e} = \tilde{\mathfrak{e}}) = 1$. Hence, for each $N \ge 1$, $X(t \land \zeta^{F_N}) = \tilde{X}^{(N)}(t)$ for all $t \ge 0$, and so, if $i \in F_N$, then $\{X(t \land \zeta^{F_N}) : t \ge 0\}$ is a Markov process generated by $\mathbf{Q}^{(N)}$ starting at i, and $\mathbb{P}(X(t) = \Delta, t \in [\mathfrak{e}, \infty)) = 1$.

To see that the distribution of $\{X(t) : t \ge 0\}$ is uniquely determined, let $0 = t_0 < \cdots < t_K$ and $(j_0, \ldots, j_K) \in (\mathbb{S} \cup \{\Delta\})^{K+1}$ be given, and set $A = \{X(t_k) = j_k, 0 \le k \le K\}$. If $j_0 \ne i$, then $\mathbb{P}(A) = 0$. Thus, assume that $j_0 = i$, in which case,

$$\mathbb{P}(A) = \sum_{\ell=1}^{K} \mathbb{P}(A \cap \{t_{\ell-1} < \mathfrak{e} \le t_{\ell}\}) + \mathbb{P}(A \cap \{\mathfrak{e} > t_{K}\})$$

For each $1 \le \ell \le K + 1$, set $A_{\ell} = \{X(t_k) = j_k, 0 \le k < \ell\}$. Then $A_{K+1} = A$ and

$$\mathbb{P}(A \cap \{t_{\ell-1} < \mathfrak{e} \le t_{\ell}\}) = \begin{cases} 0 & \text{if } j_k = \Delta \text{ for some } 0 \le k < \ell \\ \text{or } j_k \neq \Delta \text{ for some } \ell \le k \le K \\ \mathbb{P}(A_{\ell} \cap \{\mathfrak{e} > t_{\ell-1}\}) & \text{otherwise.} \end{cases}$$

Assuming that $j_k \neq \Delta$ for $0 \leq k < \ell$ and that $j_k = \Delta$ for $\ell \leq k \leq K$,

$$\mathbb{P}(A_{\ell} \cap \{\mathfrak{e} > t_{\ell-1}\}) = \lim_{N \to \infty} \mathbb{P}(A_{\ell} \cap \{\zeta^{(N)} > t_{\ell-1}\})$$
$$= \lim_{N \to \infty} \mathbb{P}(X(t_k \wedge \zeta^{F_N}) = j_k, 0 \le k < \ell)$$
$$= \lim_{N \to \infty} \mathbb{P}(\tilde{X}^{(N)}(t_k) = j_k, 0 \le k < \ell).$$

Hence, in all cases $\mathbb{P}(A)$ is uniquely determined by (1), (2), and (3).

Finally, when $\mathbb{P}(\mathfrak{e} = \infty) = 1$, then the last assertion is covered by Lemma 5.3.3.

As Theorem 5.3.6 makes clear, aside from the choice of Δ and the initial state *i*, the distribution of the process $\{X(t) : t \ge 0\}$ and the random variable \mathfrak{e} are completely determined by **Q**. For this reason, when $\mathbb{P}(\mathfrak{e} = \infty) = 1$, we say that the process *does not explode* and will call any right continuous, piecewise constant process with this distribution a *Markov process generated by* **Q** *starting from i*.

Corollary 5.3.7 Assume that explosion does not occur for any $i \in S$. Then, for each probability vector $\boldsymbol{\mu}$ on S there exists a right continuous, piecewise constant Markov process $\{X(t) : t \ge 0\}$ with initial distribution $\boldsymbol{\mu}$ such that, for each $i \in S$ with $\boldsymbol{\mu}_i > 0$, the conditional distribution of $\{X(t) : t \ge 0\}$ given that X(0) = i is a Markov process generated by \mathbf{Q} starting from i. In particular, the distribution of such a process is uniquely determined by \mathbf{Q} and $\boldsymbol{\mu}$, and the function $t \rightsquigarrow \mathbf{P}(t)$ in Lemma 5.3.3 is the transition probability function for $\{X(t) : t \ge 0\}$.

Proof For each $i \in S$ let $\{X_i(t) : t \ge 0\}$ be a Markov process generated by \mathbb{Q} starting from i, let X_0 be a random variable with distribution μ which is independent of $\sigma(\{X_i(t) : i \in S \& t \ge 0\})$, and set $X(t) = X_{X_0}(t)$. Then it is easy to check that $\{X(t) : t \ge 0\}$ has all the required properties. Furthermore, since the distribution of X(0) is μ and the conditional distribution of $\{X(t) : t \ge 0\}$ given that X(0) = i is uniquely determined by i and \mathbb{Q} when $\mu_i > 0$, the distribution of $\{X(t) : t \ge 0\}$ is uniquely determined by \mathbb{Q} and μ .

We will call a right continuous, piecewise constant process whose distribution is the one described in Corollary 5.3.7 a *Markov process generated by* \mathbf{Q} *with initial distribution* $\boldsymbol{\mu}$.

Remark Unfortunately, the relationship between \mathbf{Q} and the process in Corollary 5.3.7 is a quite unsatisfactory because it fails to give a direct connection. What one would like is a description like that given in the conditions (a) and (b) at the beginning of Sect. 5.2.1. Indeed, we know that our process is right continuous and piecewise constant and has the specified initial distribution. Furthermore, we know that it satisfies a weak form of (5.2.1). To be precise, because

$$\mathbb{P}(X(t+h) = j \mid X(\tau), \tau \in [0, t]) = (\mathbf{P}(h))_{X(t)j},$$

and, by (5.3.5),

$$\left(\mathbf{P}(h)\right)_{ij} = \delta_{ij} + h\mathbf{Q} + o(h) = (1 - hR_i)\delta_{ij} + h\left(\mathbf{P}(h)\right)_{ij} + o_{ij}(h),$$

where $o_{ij}(h)$ is a function that, as $h \searrow 0$, tends to 0 faster than h, we have that, for each $(i, j) \in \mathbb{S}^2$,

$$\mathbb{P}(X(t+h) = j \mid X(\tau), \tau \in [0, t]) = (1 - hR_i)\delta_{ij} + h(\mathbf{P}(h))_{ij} + o_{ij}(h). \quad (5.3.8)$$

What makes (5.3.8) significantly weaker than (5.2.1) is that the function $o_{ij}(h)$ depends on *i* and *j* and, when \Re is are unbounded, the rate at which it tends to 0 will depend on them also. Using more sophisticated techniques, one can show that, under the condition in Corollary 5.3.7, (a) together with the initial distribution and (5.3.8) uniquely determine the distribution of { $X(t) : t \ge 0$ }, but these techniques require material that we have not covered.

5.3.2 Criteria for Non-explosion or Explosion

In this subsection we will first develop two criteria that guarantee *non-explosion*: $e = \infty$ with probability 1. We will also derive a condition that guarantees that explosion occurs with probability 1.

Theorem 5.3.9 If **P** is a transition probability matrix and if $i \in S$ is **P**-recurrent, then for every choice of rates \mathfrak{R} , $\mathbb{P}(\mathfrak{e} = \infty) = 1$ for the process in Theorem 5.3.6 corresponding to the *Q*-matrix determined by \mathfrak{R} and **P** starting at *i*.

Proof By Lemma 5.3.2, what we must show is that $\mathbb{P}(\tilde{J}_{\infty} = \infty) = 1$. Equivalently, if $\{\tilde{X}_n : n \ge 0\}$ is a Markov chain with transition probability **P** with $\tilde{X}_0 = i$ and if $\{\tilde{E}_n : n \ge 1\}$ is a sequence to mutually independent independent unit exponential random variables that are independent of $\sigma(\{\tilde{X}_n : n \ge 0\})$, then we must show that $\sum_{n=0}^{\infty} (R_{\tilde{X}_{n-1}})^{-1}\tilde{E}_n = \infty$ with probability 1. To this end, let $\tilde{\rho}_i^{(m)}$ be the time of the *m*th return of $\{\tilde{X}_n : n \ge 0\}$ to *i*. Then the set *A* on which $\tilde{\rho}_i^{(m)} < \infty$ for all $m \ge 0$ is a $\sigma(\{\tilde{X}_n : n \ge 0\})$ -measurable set having probability 1, and

$$\sum_{n=0}^{\infty} (R_{\tilde{X}_{n-1}})^{-1} \tilde{E}_n \ge (R_i)^{-1} \sum_{m=1}^{\infty} \tilde{E}_{\tilde{\rho}_i^{(m)} + 1} \quad \text{on } A.$$

Since, given $\sigma({\tilde{X}_n : n \ge 0}), {\tilde{E}_{\tilde{\rho}_i^{(m)}+1} : m \ge 1}$ is a sequence of mutually independent, unit exponential random variables and therefore the subset of *A* for which $\sum_{m=1}^{\infty} \tilde{E}_{\tilde{\rho}_i^{(m)}+1} < \infty$ has probability 0, the proof is complete.

I will group our second non-explosion criterion together with our criterion for explosion as two parts of the same theorem. In this theorem, the process is the one in Theorem 5.3.6 determined by the rates \Re and transition probability **P**.

Theorem 5.3.10 If there exists a non-negative function u on \mathbb{S} with the properties that $U_N \equiv \inf_{j \notin F_N} u(j) \longrightarrow \infty$ as $N \to \infty$ and, for some $\alpha \in [0, \infty)$,

$$\sum_{j\in\mathbb{S}} (\mathbf{P})_{ij} u(j) \le \left(1 + \frac{\alpha}{R_i}\right) u(i) \quad \text{whenever } i \in \mathbb{S} \text{ and } R_i > 0,$$

then $\mathbb{P}(\mathfrak{e} = \infty | X(0) = i) = 1$ for all $i \in \mathbb{S}$. On the other hand, if, for some $i \in \mathbb{S}$, $R_j > 0$ whenever $i \rightarrow j$ and there exists a non-negative function u on \mathbb{S} with the property that, for some $\epsilon > 0$,

$$\sum_{\{k:i\to k\}} (\mathbf{P})_{jk} u(k) \le u(j) - \frac{\epsilon}{R_j} \quad \text{whenever } i \to j,$$

then $\mathbb{P}(\mathfrak{e} = \infty | X(0) = i) = 0.$

Proof To prove the first part, for each $N \ge 1$, set $u^{(N)}(j) = u(j)$ when $j \in F_N$ and $u^{(N)}(j) = U_N$ when $j \notin F_N$. It is an easy matter to check that if $\mathbf{Q}^{(N)} = \mathbf{R}^{(N)}(\mathbf{P} - \mathbf{I})$, where $(\mathbf{R}^{(N)})_{ij} = R_i^{(N)} \delta_{i,j}$, and $\mathbf{u}^{(N)}$ is the column vector determined by $(\mathbf{u}^{(N)})_i = u^{(N)}(i)$, then, for all $i \in \mathbb{S}$, $(\mathbf{Q}^{(N)}\mathbf{u}^{(N)})_i \le \alpha u^{(N)}(i)$. Hence, by Kolmogorov's forward equation (5.2.14) applied to the transition probability function $t \rightsquigarrow \mathbf{P}^{(N)}(t)$ corresponding to $\mathbf{Q}^{(N)}$,

$$\frac{d}{dt} \left(\mathbf{P}^{(N)}(t) \mathbf{u}^{(N)} \right)_i \le \alpha \left(\mathbf{P}^{(N)}(t) \mathbf{u}^{(N)} \right)_i,$$

and so, by Gronwall's inequality, $(\mathbf{P}^{(N)}(T)\mathbf{u}^{(N)})_i \leq e^{\alpha T} u^{(N)}(i)$. But, since

$$u^{(N)}(X^{(N)}(T)) = u^{(N)}(X^{(N)}(\zeta_N)) \ge U_N \text{ if } \zeta_N \le T,$$

this means that

$$\mathbf{P}(\zeta_N \le T \mid X(0) = i) \le \frac{1}{U_N} \mathbb{E}\left[u^{(N)}(X^{(N)}(T)) \mid X(0) = i\right] \\ = \frac{(\mathbf{P}^{(N)}(T)\mathbf{u}^{(N)})_i}{U_N} \le \frac{e^{\alpha T}u^{(N)}(i)}{U_N} \le \frac{e^{\alpha T}u(i)}{U_N} \quad \text{if } i \in F_N,$$

and so, by the monotone convergence theorem, we have now proved that $\mathbf{P}(\mathfrak{e} \leq T | X(0) = i) \leq \lim_{N \to \infty} \mathbb{P}(\zeta_N \leq T | X(0) = i) = 0.$

In proving the second assertion, we may and will assume that $i \rightarrow j$ for all $j \in S$. Now take $u^{(N)}(j) = u(j)$ if $j \in F_N$ and $u^{(N)}(j) = 0$ if $j \notin F_N$, and use $\mathbf{u}^{(N)}$ to denote the associated column vector. Clearly $(\mathbf{Q}^{(N)}\mathbf{u}^{(N)})_j$ is either less than or equal to $-\epsilon$ or equal to 0 according to whether $j \in F_N$ or $j \notin F_N$. Hence, again by Kolmogorov's forward equation,

$$\frac{d}{dt} \left(\mathbf{P}^{(N)}(t) \mathbf{u}^{(N)} \right)_j \leq -\epsilon \mathbf{1}_{F_N}(j) \sum_{k \in F_N} \left(\mathbf{P}^{(N)}(t) \right)_{jk},$$

and so

$$\mathbb{E}\left[u^{(N)}(X^{(N)}(T)) \mid X^{(N)}(0) = i\right] - u^{(N)}(i)$$

$$\leq -\epsilon \mathbb{E}\left[\int_{0}^{T} \mathbf{1}_{F_{N}}(X^{(N)}(t)) dt \mid X^{(N)}(0) = i\right] = -\epsilon \mathbb{E}\left[T \wedge \zeta_{N} \mid X^{(N)}(0) = i\right].$$

Since $u^{(N)} \ge 0$, this means that $\mathbb{E}[\zeta_N | X^{(N)}(0) = i] \le \frac{u(i)}{\epsilon}$ for all N, and so $\mathbb{E}[\mathfrak{e}|X(0) = i] \le \frac{u(i)}{\epsilon} < \infty$.

5.3.3 What to Do when Explosion Occurs

Although I do not intend to repent, I feel compelled to admit that I have ignored what, from the mathematical standpoint, is the most interesting aspect of the theory under consideration. Namely, so far I have said nothing about the myriad options one has when explosion occurs with positive probability, and in this subsection I will discuss only the most banal of the many choices available.

In Theorem 5.3.6, when explosion occurred I banished the process to a state Δ outside of S. This is the least imaginative option and, for reasons that I will not attempt to explain, is called the minimal extension. It always works, but it has the disadvantage that it completely ignores the manner in which the explosion took place. For example, when $S = \mathbb{Z}$, explosion can occur because, given that $\mathfrak{e} < \infty$, $\lim_{t \neq \mathfrak{e}} X(t) = +\infty$ with probability 1 or because, although $\lim_{t \neq \mathfrak{e}} |X(t)| = \infty$ with probability 1, both $\lim_{t \neq \mathfrak{e}} X(t) = +\infty$ and $\lim_{t \neq \mathfrak{e}} X(t) = -\infty$ occur with positive probability. In the latter case, one might want to record which of the two possibilities occurred, and this could be done by introducing two new absorbing states, Δ_+ for those paths that escape via $+\infty$ and Δ_- for those that escape via $-\infty$.

Alternatively, rather than thinking of the explosion time as a time to banish the process from S, one can turn it into a time of reincarnation by redistributing the process over S at time \mathfrak{e} and running it again until it explodes, etc. Obviously, there is an infinity of possibilities. Suffice it to say that the preceding discussion barely scratches the surface.

5.4 Ergodic Properties

In this section we will examine the ergodic behavior of the Markov processes which we have been discussing in this chapter. Our running assumption will be that the process with which we are dealing is a continuous time Markov process of the sort described in Sect. 5.2 under the condition that there is no explosion.

5.4.1 Classification of States

We begin by classifying the states of S.

Given a *Q*-matrix **Q**, let (cf. (5.2.12)) \Re and **P** be the canonical rates and transition probability that it determines. Then from Theorems 5.2.10 and 5.3.6, it is clear that the states visited by a Markov process generated by **Q** starting from a state *i* are exactly the same as those visited by the Markov chain with transition probability **P** starting from *i*. Hence, we will say that *i* is **Q**-recurrent or **Q**-transient if and only if it is recurrent or transient for the Markov chain with transition probability **P**. In particular, *i* is **Q**-recurrent if $R_i = 0$. Similarly, we will say that *j* is **Q**-accessible from *i* and will write $i \stackrel{\mathbf{Q}}{\rightarrow} j$ if and only if $\mathbf{P}_{ij}^n > 0$ for some $n \ge 0$, and we will say that *i* **Q**-communicates with *j* and will write $i \stackrel{\mathbf{Q}}{\Rightarrow} j$ if and only if $i \stackrel{\mathbf{Q}}{\Rightarrow} j$ and $j \stackrel{\mathbf{Q}}{\Rightarrow} i$.

From the results in Sect. 3.1, we know that Q-accessibility is a transitive relation and that Q-communication is an equivalence relation. In addition, Q-recurrence and Q-transitivity are Q-communicating class properties. Thus, if the state space S is Q*irreducible* in the sense that all states Q-communicate with one another, then either all its states are Q-recurrent or all are Q-transient.

I next want to describe these results in terms of the paths of the associate Markov process. Thus, let $\{X(t) : t \ge 0\}$ be a Markov process generated by **Q**. In order to describe the analog of the first return time to a state *i*, let J_n be the time when the process makes its *n*th jump. That is,

$$J_n = \begin{cases} \inf\{t > 0 : X(t) \neq X(t-)\} & \text{if } n = 1\\ \inf\{t > J_{n-1} : X(t) \neq X(t-)\} & \text{if } n > 1, \end{cases}$$

where $X(t-) \equiv \lim_{s \nearrow t} X(s)$. Then the first return time to *i* is

$$\sigma_i \equiv \inf \{ t \ge J_1 : X(t) = i \}.$$

More generally, define $\sigma_i^{(0)} = 0$, $\sigma_i^{(1)} = \sigma_i$, and, for m > 1,

$$\sigma_i^{(m)} = \begin{cases} \inf\{t \ge J_{\ell+1} : X(t) = i\} & \text{if } \sigma_i^{(m-1)} = J_{\ell} < \infty \\ \infty & \text{if } \sigma_i^{(m-1)} = \infty. \end{cases}$$

We will call $\sigma_i^{(m)}$ the *m*th return time to *i*.

Clearly, *i* is **Q**-recurrent if and only if $\mathbb{P}(\sigma_i < \infty | X(0) = i) = 1$. Furthermore, if *i* is **Q**-recurrent, then, with probability 1, a process started at *i* returns to *i* infinitely often. To see this, let $\{X_n : n \ge 0\}$ be a Markov chain with transition probability **P**, and observe that, by Theorems 5.2.10 and 5.3.6, $\mathbb{P}(\sigma_i^{(m)} < \infty) = \mathbb{P}(\rho_i^{(m)} < \infty)$, where $\rho_i^{(m)}$ is the *m*th return time of $\{X_n : n \ge 0\}$ to *i*. Thus, by Theorem 2.3.6, for

any $(i, j) \in \mathbb{S}^2$ and $m \ge 1$

$$\mathbb{P}\left(\sigma_{j}^{(m)} < \infty \mid X(0) = i\right)$$
$$= \mathbb{P}\left(\sigma_{j} < \infty \mid X(0) = i\right) \mathbb{P}\left(\sigma_{j}^{(m)} < \infty \mid X(0) = j\right)^{m-1}.$$
(5.4.1)

Just as in the case of chains, these considerations are intimately related to the amount of time that the process spends in a state. Indeed,

$$\mathbb{E}\left[\int_0^\infty \mathbf{1}_{\{j\}} (X(t)) dt \, \Big| \, X(0) = i\right]$$
$$= \sum_{m=0}^\infty \mathbb{E}\left[\int_{\sigma_j^{(m)}}^{\sigma_j^{(m+1)}} \mathbf{1}_{\{j\}} (X(t)) dt, \sigma_j^{(m)} < \infty \,\Big| \, X(0) = 0\right].$$

Furthermore, using the notation in Theorem 5.2.10 and letting $\tilde{\sigma}_j^{(m)}$ denote the *m*th return time of $\{\tilde{X}(t) : t \ge 0\}$ to *j*, we have that

$$\begin{split} & \mathbb{E}\bigg[\int_{\sigma_{j}^{(m+1)}}^{\sigma_{j}^{(m+1)}} \mathbf{1}_{\{j\}}\big(X(t)\big) \, dt, \, \sigma_{j}^{(m)} < \infty \, \Big| \, X(0) = 0 \bigg] \\ &= \sum_{\ell=0}^{\infty} \mathbb{E}\big[J_{\ell+1} - J_{\ell}, \, \sigma_{j}^{(m)} = J_{\ell+1} < \infty \, \Big| \, X(0) = i \big] \\ &= \sum_{\ell=0}^{\infty} \mathbb{E}\big[\tilde{E}_{\ell+1}, \, \tilde{\sigma}_{j}^{(m)} = \tilde{J}_{\ell+1} < \infty \, \Big| \, \tilde{X}(0) = i \big] \\ &= \frac{1}{R_{j}} \sum_{\ell=0}^{\infty} \mathbb{P}\big(\tilde{\sigma}_{j}^{(m)} = \tilde{J}_{\ell+1} < \infty \, \Big| \, \tilde{X}(0) = i \big) \\ &= \frac{\mathbb{P}(\tilde{\sigma}_{j}^{(m)} < \infty \, | \, \tilde{X}(0) = i)}{R_{j}} = \frac{\mathbb{P}(\sigma_{j}^{(m)} < \infty \, | \, X(0) = i)}{R_{j}} \\ &= \begin{cases} 0 \quad \text{if } \, \mathbb{P}(\sigma_{j} < \infty \, | \, X(0) = i) = 0 \\ \infty \quad \text{if } R_{j} = 0 \text{ and } \, \mathbb{P}(\sigma_{j} < \infty \, | \, X(0) = i) > 0 \\ &= \begin{cases} 0 \quad \text{if } R_{j} = 0 \text{ and } \, \mathbb{P}(\sigma_{j} < \infty \, | \, X(0) = i) > 0 \\ &= \frac{\mathbb{P}(\sigma_{j} < \infty \, | \, X(0) = i) \mathbb{P}(\sigma_{j} < \infty \, | \, X(0) = j)^{(m-1)^{+}}}{R_{j}} \quad \text{otherwise.} \end{cases} \end{split}$$

Hence, by exactly the argument with which we passed from Theorem 2.3.6 to (2.3.7), we now know that

$$\mathbb{E}\left[\int_{0}^{\infty} \mathbf{1}_{\{j\}}(X(t)) dt \left| X(0) = i \right] = \frac{1}{R_{j}} \left(\delta_{i,j} + \frac{\mathbb{P}(\sigma_{j} < \infty | X(0) = i)}{\mathbb{P}(\sigma_{j} = \infty | X(0) = j)} \right) \\ \mathbb{E}\left[\int_{0}^{\infty} \mathbf{1}_{\{i\}}(X(t)) dt \left| X(0) = i \right] = \infty \\ \iff \mathbb{P}\left(\int_{0}^{\infty} \mathbf{1}_{\{i\}}(X(t)) dt = \infty \left| X(0) = i \right] = 1$$
(5.4.2)
$$\mathbb{E}\left[\int_{0}^{\infty} \mathbf{1}_{\{i\}}(X(t)) dt \left| X(0) = i \right] < \infty \\ \iff \mathbb{P}\left(\int_{0}^{\infty} \mathbf{1}_{\{i\}}(X(t)) dt < \infty \left| X(0) = i \right] = 1,$$

where, in the first line when $R_j = 0$, the right hand side is ∞ if either i = j or $\mathbb{P}(\sigma_j | X(0) = i) > 0$ and is 0 if $i \neq j$ and $\mathbb{P}(\sigma_j | X(0) = i) = 0$. In particular, when *i* is **Q**-transient, the expected amount of time that the process spends at *j* is uniformly bounded independent of its initial distribution.

Our final goal in this subsection is to prove the following statement.

Theorem 5.4.3 For any given state $i \in S$, the following are equivalent:

- (1) *i* is **Q**-recurrent.
- (2) There is a $t \in (0, \infty)$ such that *i* is recurrent relative to the transition probability $\mathbf{P}(t)$.
- (3) *i* is recurrent relative to $\mathbf{P}(t)$ for all $t \in (0, \infty)$.

Proof We will prove this equivalence by checking that the same equivalence holds when "recurrent" is replaced throughout by "transient." To this end, first observe that

$$\mathbb{E}\left[\int_0^\infty \mathbf{1}_{\{i\}} (X(t)) dt \, \Big| \, X(0) = i\right] = \int_0^\infty (\mathbf{P}(t))_{ii} dt,$$

and therefore, from the first line of (5.4.2), that

i is **Q**-transient
$$\iff \int_0^\infty (\mathbf{P}(t))_{ii} dt < \infty.$$
 (5.4.4)

Next, notice that for $0 \le s < t$,

$$\left(\mathbf{P}(t)\right)_{ii} \ge \left(\mathbf{P}(t-s)\right)_{ii} \left(\mathbf{P}(s)\right)_{ii} \ge e^{-(t-s)R_i} \left(\mathbf{P}(s)\right)_{ii} \tag{5.4.5}$$

since $(\mathbf{P}(h))_{ii} \ge \mathbb{P}(J_1 > h | X(0) = i) = e^{-hR_i}$. Hence, for any t > 0 and $n \in \mathbb{N}$,

$$\left(\mathbf{P}(t)^{n+1}\right)_{ii} = \left(\mathbf{P}\left((n+1)t\right)\right)_{ii} \ge e^{-tR_i}\left(\mathbf{P}(\tau)\right)_{ii} \quad \text{for all } \tau \in [nt, (n+1)t]$$

and

$$e^{-tR_i} \left(\mathbf{P}(t)^n \right)_{ii} = e^{-tR_i} \left(\mathbf{P}(nt) \right)_{ii} \le \left(\mathbf{P}(\tau) \right)_{ii} \quad \text{for all } \tau \in \left[nt, (n+1)t \right].$$

Since this means that

$$te^{-tR_i}\sum_{n=0}^{\infty} \left(\mathbf{P}(t)^n\right)_{ii} \leq \int_0^{\infty} \left(\mathbf{P}(\tau)\right)_{ii} d\tau \leq te^{tR_i}\sum_{n=0}^{\infty} \left(\mathbf{P}(t)^{n+1}\right)_{ii},$$

the asserted implications are now immediate from (5.4.4).

5.4.2 Stationary Measures and Limit Theorems

In this subsection I will prove the analogs of the results in Sects. 4.1.6 and 4.1.7, and, because periodicity plays no role here, the statements are simpler.

Theorem 5.4.6 *For each* $j \in \mathbb{S}$

$$\hat{\pi}_{jj} \equiv \lim_{t \to \infty} \left(\mathbf{P}(t) \right)_{jj}$$
 exists

and

$$\lim_{t \to \infty} (\mathbf{P}(t))_{ij} = \hat{\pi}_{ij} \equiv \mathbb{P}(\sigma_j < \infty | X(0) = i) \hat{\pi}_{jj} \quad \text{for } i \neq j.$$

Moreover, if $\hat{\pi}_{jj} > 0$, then $\hat{\pi}_{ii} > 0$ for all $i \in C \equiv \{i : i \stackrel{\mathbf{Q}}{\Leftrightarrow} j\}$, and when the row vector $\hat{\pi}^{C}$ is determined by $(\hat{\pi}^{C})_{i} = \mathbf{1}_{C}(i)\hat{\pi}_{ii}, \hat{\pi}^{C}$ is, for each s > 0, the unique probability vector $\boldsymbol{\mu} \in \text{Stat}(\mathbf{P}(s))$ for which $(\boldsymbol{\mu})_{k} = 0$ if $k \notin C$. In fact, if $\boldsymbol{\mu} \in \text{Stat}(\mathbf{P}(s))$ for some s > 0, then, for each $j \in C$,

$$(\boldsymbol{\mu})_j = \left(\sum_{\substack{i \in j \\ i \leftrightarrow j}} (\boldsymbol{\mu})_i\right) \hat{\pi}_{jj},$$

Proof We begin with the following continuous-time version of the renewal equation (cf. (4.1.6)):

$$\left(\mathbf{P}(t)\right)_{ij} = e^{-tR_i}\delta_{i,j} + \mathbb{E}\left[\left(\mathbf{P}(t-\sigma_j)\right)_{jj}, \sigma_j \le t \mid X(0) = i\right].$$
(5.4.7)

The proof of (5.4.7) runs as follows. First, write $(\mathbf{P}(t))_{ij}$ as

$$\mathbb{P}(X(t) = j \& J_1 > t \mid X(0) = i) + \mathbb{P}(X(t) = j \& J_1 \le t \mid X(0) = i).$$

Clearly, the first term is 0 unless i = j, in which case it is equal e^{-tR_i} . To handle the second term, write it as

$$\sum_{m=1}^{\infty} \mathbb{P}(X(t) = j \& \sigma_j = J_m \le t \mid X(0) = i),$$

and observe that (cf. (5.2.9) and (5.2.11))

$$\mathbb{P}(X(t) = j \& \sigma_j = J_m \le t \mid X(0) = i)$$

= $\mathbb{P}(\Phi^{\mathfrak{R}}(t - \tilde{J}_m; (\tilde{E}_{m+1}, \dots, \tilde{E}_{m+n}, \dots), (j, \tilde{X}_{m+1}, \dots, \tilde{X}_{m+n}, \dots)) = j$
& $\tilde{\sigma}_j = \tilde{J}_m \le t \mid \tilde{X}_0 = i)$
= $\mathbb{E}[(\mathbf{P}(t - J_m))_{jj}, \sigma_j = J_m \le t \mid X(0) = i].$

Hence, after summing this over $m \ge 1$ and combining this result with the preceding, one arrives at (5.4.7).

Knowing (5.4.7), we see that the first part of the theorem will be proved once we treat the case when i = j. To this end, we begin by observing that, because, by (5.4.5), $(\mathbf{P}(s))_{ii} \ge e^{-sR_i} > 0$ for all s > 0 and $i \in \mathbb{S}$, each *i* is aperiodic relative to $\mathbf{P}(s)$, and so, by (4.1.15), we know that $\pi(s)_{ii} \equiv \lim_{n\to\infty} (\mathbf{P}(s)^n)_{ii}$ exists for all s > 0 and $i \in \mathbb{S}$. We need to show that $\pi(1)_{ii} = \lim_{t\to\infty} (\mathbf{P}(t))_{ii}$ for all t > 0, and when $\pi(1)_{ii} = 0$, this is easy. Indeed, by (5.4.5),

$$\overline{\lim_{t\to\infty}} (\mathbf{P}(t))_{ii} \le e^{R_i} \overline{\lim_{t\to\infty}} (\mathbf{P}(\lfloor t \rfloor + 1))_{ii} = e^{R_j} \pi(1)_{ii},$$

where $\lfloor t \rfloor$ denotes the integer part of *t*.

In view of the preceding, what remains to be proved in the first assertion is that $\lim_{t\to\infty} (\mathbf{P}(t))_{ii} = \pi(1)_{ii}$ when $\pi(1)_{ii} > 0$, and the key step in our proof will be the demonstration that $\pi(s)_{ii} = \pi(1)_{ii}$ for all s > 0, a fact which we already know when $\pi(1)_{ii} = 0$. Thus, assume that $\pi(1)_{ii} > 0$, and let *C* be the communicating class of *i* relative to $\mathbf{P}(1)$. By Theorem 5.4.3, *C* is also the communicating class of *i* relative to $\mathbf{P}(s)$ for every s > 0. In addition, because $\pi(1)_{ii} > 0$, *i* is recurrent, in fact positive recurrent, relative to $\mathbf{P}(1)$, and therefore, by Theorem 5.4.3, it is also recurrent relative to $\mathbf{P}(s)$ for all s > 0. Now determine the row vector $\pi(1)^C$ so that $(\pi(1)^C)_j = \mathbf{1}_C(j)\pi(1)_{jj}$. Then, because $\pi(1)_{ii} > 0$, we know, by Theorem 4.1.10, that $\pi(1)^C$ is the one and only $\boldsymbol{\mu} \in \text{Stat}(\mathbf{P}(1))$ which vanishes off of *C*. Next, given s > 0, consider $\boldsymbol{\mu} = \pi(1)^C \mathbf{P}(s)$. Then $\boldsymbol{\mu}$ is a probability vector and, because $(\mathbf{P}(1))_{jk} = 0$ when $j \in C$ and $k \notin C$, $\boldsymbol{\mu}$ vanishes off of *C*. Also, because

$$\mu \mathbf{P}(1) = \pi (1)^C \mathbf{P}(s) \mathbf{P}(1) = \pi (1)^C \mathbf{P}(s+1)$$
$$= \pi (1)^C \mathbf{P}(1) \mathbf{P}(s) = \pi (1)^C \mathbf{P}(s) = \mu,$$

 μ is stationary for **P**(1). Hence, by uniqueness, we conclude that $\pi(1)^C \mathbf{P}(s) = \pi(1)^C$, and therefore that $\pi(1)^C$ is a stationary probability for **P**(*s*) which vanishes off of *C*. But, by (4.1.8) and the fact that *C* is the communicating class of *i* relative to **P**(*s*), this means that

$$\pi(1)_{ii} = \left(\pi(1)^{C}\right)_{i} = \left(\sum_{j \in C} \left(\pi(1)^{C}\right)_{j}\right) \pi(s)_{ii} = \pi(s)_{ii}.$$

To complete the proof of the first part from here, note that, by (5.4.5),

$$e^{-sR_j} (\mathbf{P}(ns))_{jj} \le (\mathbf{P}(t))_{jj} \le e^{sR_j} (\mathbf{P}((n+1)s))_{jj}$$
 when $ns \le t \le (n+1)s$,

and so, since $\pi(s)_{jj} = \pi(1)_{jj}$ and $\mathbf{P}(nt) = \mathbf{P}(t)^n$,

$$e^{-sR_j}\pi(1)_{jj} \leq \underline{\lim}_{t\to\infty} (\mathbf{P}(t))_{jj} \leq \overline{\lim}_{t\to\infty} (\mathbf{P}(t))_{jj} \leq e^{sR_j}\pi(1)_{jj}.$$

Now let $s \searrow 0$, and simply define $\hat{\pi}_{jj} = \pi(1)_{jj}$.

Given the first part, the proof of the second part is easy. Namely, if $\hat{\pi}_{jj} > 0$ and $C = \{i : i \leftrightarrow j\}$, then, for each s > 0, we know that *C* is the communicating class of *j* relative to $\mathbf{P}(s)$ and that $\hat{\pi}^C = \pi(s)^C \in \text{Stat}(\mathbf{P}(s))$. Conversely, by (4.1.8), we know that if $\boldsymbol{\mu} \in \text{Stat}(\mathbf{P}(s))$ for some s > 0, then

$$(\boldsymbol{\mu})_{j} = \left(\sum_{\{i:i \stackrel{\boldsymbol{Q}}{\Leftrightarrow} j\}} (\boldsymbol{\mu})_{i}\right) \pi(s)_{jj} = \left(\sum_{\{i:i \stackrel{\boldsymbol{Q}}{\Leftrightarrow} j\}} (\boldsymbol{\mu})_{i}\right) \hat{\pi}_{jj}.$$

With the preceding result in mind, we will say that *i* is **Q**-*positive recurrent* if $\hat{\pi}_{ii} > 0$ and will say that *i* is **Q**-*null recurrent* if *i* is **Q**-recurrent but not **Q**-positive recurrent. From the preceding and Theorem 4.1.10, we already know that *Q*-*positive recurrence is a Q*-communicating class property.

The following corollary comes at very little additional cost.

Corollary 5.4.8 (Mean Ergodic Theorem) Assume that j is **Q**-positive recurrent and that $\mathbb{P}(X(0) \stackrel{\mathbf{Q}}{\leftrightarrow} j) = 1$. Then,

$$\lim_{T\to\infty} \mathbb{E}\left[\left(\frac{1}{T}\int_0^T \mathbf{1}_{\{j\}}(X(t))\,dt - \hat{\pi}_{jj}\right)^2\right] = 0.$$

See Exercise 5.5.11 for the more refined version.

Proof The proof is really just an obvious transcription to the continuous setting of the argument used to prove Theorem 4.1.14. By precisely the argument given there, it suffices to handle the case when $\hat{\pi}^C$ is the initial distribution, where $C = \{i : j \Leftrightarrow i\}$. Thus, if $f = \mathbf{1}_{\{j\}} - \hat{\pi}_{jj}$, what we need to show is that

$$\lim_{T \to \infty} \mathbb{E}\left[\left(\frac{1}{T} \int_0^T f(X(t)) dt\right)^2\right] = 0,$$

when $\hat{\pi}^{C}$ is the distribution of X(0). But, because $\hat{\pi}^{C}$ is $\mathbf{P}(t)$ -stationary,

$$\mathbb{E}\left[\left(\frac{1}{T}\int_0^T f(X(t))dt\right)^2\right] = \frac{2}{T^2}\int_0^T \left(\int_0^t \mathbb{E}\left[f(X(s))f(X(t))\right]ds\right)dt$$
$$= \frac{2}{T^2}\int_0^T \left(\int_0^t \alpha(t-s)ds\right)dt$$
$$= \frac{2}{T}\int_0^T \left(1-\frac{t}{T}\right)\alpha(t)dt,$$

where $\alpha(t) \equiv \hat{\pi}^{C}(f\mathbf{P}(t)\mathbf{f})$, **f** being the column vector corresponding to the function f and $f\mathbf{P}(t)\mathbf{f}$ being the column vector determined by $(f\mathbf{P}(t)\mathbf{f})_{i} = f(i)(\mathbf{P}(t)\mathbf{f})_{i}$. Finally, since, for each $i \in C$, $\lim_{t\to\infty} (\mathbf{P}(t)\mathbf{f})_{i} = 0$, $\lim_{t\to\infty} \alpha(t) = 0$, and therefore

$$\left|\frac{2}{T}\int_0^T \left(1-\frac{t}{T}\right)\alpha(t)\,dt\right| \le \frac{2}{T}\int_0^T \left|\alpha(t)\right|\,dt \longrightarrow 0 \quad \text{as } T \to \infty.$$

5.4.3 Interpreting and Computing $\hat{\pi}_{ii}$

Although we have shown that the limit $\hat{\pi}_{ij} = \lim_{t\to\infty} (\mathbf{P}(t))_{ij}$ exists, we have yet to give an expression, analogous to the one in (4.1.5), for $\hat{\pi}_{ii}$. However, as we are about to see, such an expression is readily available from (5.4.7). Namely, for each $\alpha > 0$, set

$$L(\alpha)_{ij} = \alpha \mathbb{E}\left[\int_0^\infty e^{-\alpha t} \mathbf{1}_{\{j\}} (X(t)) dt \, \Big| \, X(0) = i\right] = \alpha \int_0^\infty e^{-\alpha t} (\mathbf{P}(t))_{ij} dt.$$

Because $\hat{\pi}_{ii} = \lim_{t \to \infty} (\mathbf{P}(t))_{ii}$, the second of these representations makes it is easy to identify $\hat{\pi}_{ii}$ as $\lim_{\alpha \searrow 0} L(\alpha)_{ii}$. At the same time, from (5.4.7), we see that

$$L(\alpha)_{ii} = \frac{\alpha}{\alpha + R_i} + \mathbb{E} \Big[e^{-\alpha \sigma_i} \, \big| \, X(0) = i \Big] L(\alpha)_{ii}.$$

Hence, since

$$\frac{1 - \mathbb{E}[e^{-\alpha\sigma_i} \mid X_0 = i]}{\alpha} = \frac{\mathbb{E}[1 - e^{-\alpha\sigma_i} \mid X_0 = i]}{\alpha} \longrightarrow \mathbb{E}[\sigma_i \mid X_0 = r_i]$$

if $R_i > 0$ and $\mathbb{P}(\sigma_i = \infty | X_0 = i) = 1$ if $R_i = 0$, we have now shown that

$$\hat{\pi}_{ii} = \begin{cases} 1 & \text{if } R_i = 0\\ \frac{1}{R_i \mathbb{E}[\sigma_i | X(0) = i]} & \text{if } R_i > 0, \end{cases}$$

which, in conjunction with the second equality in Theorem 5.4.6, proves that

$$\hat{\pi}_{ij} = \begin{cases} \delta_{i,j} + (1 - \delta_{i,j}) \mathbb{P}(\sigma_j < \infty | X(0) = i) & \text{if } R_j = 0\\ \frac{\mathbb{P}(\sigma_j < \infty | X(0) = i)}{R_j \mathbb{E}[\sigma_j | X(0) = j]} & \text{if } R_j > 0. \end{cases}$$
(5.4.9)

Of course, as an immediate corollary of (5.4.9), we now know that *i* is positive recurrent if and only if either $R_i = 0$ or $\mathbb{E}[\sigma_i | X(0) = i] < \infty$.

Finally, we will apply the results in Sect. 3.2 to the setting here. To this end, first observe that as an immediate consequence of (5.2.2), for any probability vector μ ,

$$\mu \mathbf{Q} = \mathbf{0} \iff \mu \mathbf{P}(t) \text{ for all } t > 0 \text{ when } \Re \text{ is bounded.}$$
 (5.4.10)

Theorem 5.4.11 Assume that S is finite and that there is only one stationary probability $\hat{\pi}$ for $t \rightsquigarrow \mathbf{P}(t)$. Then 0 is a simple eigenvalue of \mathbf{Q} and

$$(\hat{\boldsymbol{\pi}})_i = \frac{\det((-Q)^{\{i\}})}{\Pi_{-Q}} = \frac{\det((-Q)^{\{i\}})}{\sum_{j \in \mathbb{S}} \det((-Q)^{\{j\}})}$$

for all $i \in \mathbb{S}$.

Proof The argument is essentially the same as the one for chains. To see that 0 is a simple eigenvalue, suppose not. Then there exists a column vector $\mathbf{v} \neq \mathbf{0}$ such that $\mathbf{Q}\mathbf{v} = \mathbf{0}$ and $\hat{\pi}\mathbf{v} = 0$. By (5.4.10), $\mathbf{v} = \mathbf{P}(t)\mathbf{v}$ for all t > 0, and so Theorem 5.4.6 leads to the contradiction

$$\mathbf{v} = \lim_{t \to \infty} \mathbf{P}(t) \mathbf{v} = (\mathbf{v}, \hat{\boldsymbol{\pi}})_{\mathbb{C}^{\mathbb{S}}} \mathbf{1} = \mathbf{0}.$$

Knowing that 0 is a simple eigenvalue and that Q1 = 0, the proof of the second assertion is precisely the same as the proof of the analogous result in Theorem 3.2.6.

The reason why I wrote the equation for $\hat{\pi}$ in terms of $-\mathbf{Q}$ instead of \mathbf{Q} is that, by arguments like those in the chain case, $(-\mathbf{Q})^{\{i\}}$ is non-negative for all *i* and is strictly positive for \mathbf{Q} -recurrent *i*'s.

5.5 Exercises

Exercise 5.5.1 The purpose of this exercise is to give another derivation of (5.1.5). Let $\{E_n : n \ge 1\}$ be a sequence of mutually independent, unit exponential random variables, and define $\{J_n : n \ge 0\}$ and $\{N(t) : t \ge 0\}$ accordingly, as in Sect. 5.1.1. Given $0 = t_0 < \cdots < t_L$ and $0 = n_0 < \cdots < n_L$, use the change of variables formula for multi-dimensional integrals to justify:

$$\begin{aligned} \mathbf{P} \Big(N(t_1) &= n_1, \dots, N(t_L) = n_L \Big) \\ &= \mathbf{P} (J_{n_1} \leq t_1 < J_{n_1+1}, \dots, J_{n_L} \leq t_\ell < J_{n_L+1}) \\ &= \int \dots \int \exp \left(-\sum_{\ell=1}^{n_L+1} \xi_\ell \right) d\xi_1 \dots d\xi_{n_L+1} \\ &= \int \dots \int e^{-\eta_{n_L+1}} d\eta_1 \dots d\eta_{n_L+1} \\ &= e^{-t_L} \prod_{\ell=1}^L \operatorname{vol}(\Delta_\ell) = \prod_{\ell=1}^L e^{-(t_\ell - t_{\ell-1})} \frac{(t_\ell - t_{\ell-1})^{n_\ell - n_{\ell-1}}}{(n_\ell - n_{\ell-1})!}, \end{aligned}$$

where

$$A = \left\{ (\xi_1, \dots, \xi_{n_L+1}) \in (0, \infty)^{n_L+1} : \sum_{1}^{n_\ell} \xi_j \le t_\ell < \sum_{1}^{n_\ell+1} \xi_j \text{ for } 1 \le \ell \le L \right\},\$$

$$B = \left\{ (\eta_1, \dots, \eta_{n_L+1}) \in (0, \infty)^{n_L+1} : \\\eta_i < \eta_{i+1} \text{ for } 1 \le i \le n_L \& \eta_{n_\ell} \le t_\ell < \eta_{n_\ell+1} \text{ for } 1 \le \ell \le L \right\},\$$

$$\Delta_\ell = \left\{ (u_1, \dots, u_{n_\ell - n_{\ell-1}}) \in \mathbb{R}^{n_\ell - n_{\ell-1}} : t_{\ell-1} \le u_1 < \dots < u_{n_\ell - n_{\ell-1}} \le t_\ell \right\}.$$

When $n_{\ell} = n_{\ell-1}$ for some ℓ 's, begin by choosing $0 = \ell_0 < \cdots < \ell_K = L$ so that $n_{\ell_k} = \cdots = n_{\ell_{k+1}-1} < n_{\ell_{k+1}}$ for $0 \le k < L$. Show that

$$\mathbb{P}(N(t_{\ell}) = t_{\ell}, 1 \le \ell \le L)$$

= $\mathbb{P}(J_{n_{\ell_k}} \le t_{\ell_k} \& J_{n_{\ell_k-1}+1} < t_{\ell_k-1}, 1 \le k \le K \text{ and } J_{n_L+1} > t_L),$

and, proceeding as above, conclude that

$$\mathbb{P}\big(N(t_{\ell}) = t_{\ell}, 1 \leq \ell \leq L\big)$$

= $e^{-t_{L}} \prod_{k=1}^{K} \frac{(t_{\ell_{k}} - t_{\ell_{k}-1})^{n_{\ell_{k}} - n_{\ell_{k}-1}}}{(n_{\ell_{k}} - n_{\ell_{k}-1})!} = e^{-t_{L}} \prod_{\ell=1}^{L} \frac{(t_{\ell} - t_{\ell-1})^{n_{\ell} - n_{\ell-1}}}{(n_{\ell} - n_{\ell-1})!}.$

Exercise 5.5.2 Let \mathbf{M}_1 and \mathbf{M}_2 be commuting elements of $M_{u,v}(\mathbb{S})$. After checking that

$$(\mathbf{M}_1 + \mathbf{M}_2)^m = \sum_{\ell=0}^m \binom{m}{\ell} \mathbf{M}_1^\ell \mathbf{M}_2^{m-\ell}$$

for all $m \in \mathbb{N}$, verify (5.1.9).

Exercise 5.5.3 Given a **Q**-recurrent state *i*, set $C = \{j : i \stackrel{\mathbf{Q}}{\leftrightarrow} j\}$, and show that $R_i > 0 \implies R_j > 0$ for all $j \in C$.

Exercise 5.5.4 Gronwall's is an inequality which has many forms, the most elementary of which states that if $u : [0, T] \longrightarrow [0, \infty)$ is a continuous function that satisfies

$$u(t) \le A + B \int_0^t u(\tau) d\tau \quad \text{for } t \in [0, T],$$

then $u(t) \le Ae^{Bt}$ for $t \in [0, T]$. Prove this form of *Gronwall's inequality*.

Hint: Assume that $B \neq 0$, set $U(t) = \int_0^t u(\tau) d\tau$, show that $\dot{U}(t) \leq A + BU(t)$, and conclude that $U(t) \leq \frac{A}{B}(e^{Bt} - 1)$.

Exercise 5.5.5 In this exercise we will give the continuous-time version of the ideas in Exercise 2.4.1. For this purpose, assume that S is irreducible and positive recurrent with respect to \mathbf{Q} , and use $\hat{\pi}$ to denote the unique probability vector which is stationary for each $\mathbf{P}(t)$, t > 0. Next, determine the *adjoint* semigroup { $\mathbf{P}(t)^{\top}$: $t \ge 0$ } so that

$$\left(\mathbf{P}(t)^{\top}\right)_{ij} = \frac{(\hat{\boldsymbol{\pi}})_j (\mathbf{P}(t))_{ji}}{(\hat{\boldsymbol{\pi}})_i}$$

(a) Show that $\mathbf{P}(t)^{\top}$ is a transition probability matrix and that $\hat{\boldsymbol{\pi}} \mathbf{P}(t)^{\top} = \hat{\boldsymbol{\pi}}$ for each $t \ge 0$. In addition, check that $\{\mathbf{P}(t)^{\top} : t \ge 0\}$ is the semigroup determined by the *Q*-matrix \mathbf{Q}^{\top} , where

$$\left(\mathbf{Q}^{\top}\right)_{ij} = \frac{(\hat{\boldsymbol{\pi}})_j(\mathbf{Q})_{ji}}{(\hat{\boldsymbol{\pi}})_i}.$$

(b) Let \mathbb{P} and \mathbb{P}^{\top} denote the probabilities computed for the Markov processes corresponding, respectively, to \mathbf{Q} and \mathbf{Q}^{\top} with initial distribution $\hat{\boldsymbol{\pi}}$. Show that \mathbb{P}^{\top} is the *reverse* of \mathbb{P} in the sense that, for each $n \in \mathbb{N}$, $0 = t_0 < t_1 < \cdots < t_n$, and $(j_0, \ldots, j_n) \in \mathbb{S}^{n+1}$,

$$\mathbb{P}^{\top}(X(t_m) = j_m \text{ for } 0 \le m \le n) = \mathbb{P}(X(t_n - t_m) = j_m \text{ for } 0 \le m \le n).$$

(c) Assume that $R_i \equiv -(\mathbf{Q})_{ii} > 0$ for all $i \in \mathbb{S}$, and define \mathbf{P}^\top by

$$\left(\mathbf{P}^{\top}\right)_{ij} = (1 - \delta_{i,j}) \frac{(\hat{\boldsymbol{\pi}})_j(\mathbf{Q})_{ji}}{R_i(\hat{\boldsymbol{\pi}})_i}.$$

Show that \mathbf{P}^{\top} is again a transition probability matrix, and check that $\mathbf{Q}^{\top} = \mathbf{R}(\mathbf{P}^{\top} - \mathbf{I})$ where $(\mathbf{R})_{ij} = \delta_{i,j} R_i$.

Exercise 5.5.6 Take $\mathbb{S} = \mathbb{N}$ and $(\mathbf{P})_{ij}$ equal to 1 or 0 according to whether j = i + 1 or not. Given a set of strictly positive rates \mathfrak{R} , show that, no matter what its initial distribution, the Markov process determined by \mathfrak{R} and \mathbf{P} explodes with probability 1 if $\sum_{i \in \mathbb{N}} R_i^{-1} < \infty$ and does not explode if $\sum_{i \in \mathbb{N}} R_i^{-1} = \infty$.

Exercise 5.5.7 Here is a more interesting example of explosion. Take $\mathbb{S} = \mathbb{Z}^3$, and let (cf. Exercise 3.4.1) **P** be the transition probability matrix for the symmetric, nearest neighbor random walk on \mathbb{Z}^3 . Given a set of positive rates \mathfrak{R} with the property that $\sum_{\mathbf{k}\in\mathbb{Z}^3} R_{\mathbf{k}}^{-1} < \infty$, show that, starting at every $\mathbf{k} \in \mathbb{Z}^3$, explosion occurs with probability 1 when $(\mathbf{Q})_{\mathbf{k}\ell} = R_{\mathbf{k}}((\mathbf{P})_{\mathbf{k}\ell} - \delta_{\mathbf{k},\ell})$.

Hint: Apply the criterion in the second half of Theorem 5.3.10 to a function of the form

$$\sum_{\boldsymbol{\ell}\in\mathbb{Z}^3}\frac{1}{R_{\boldsymbol{\ell}}}\left(\alpha^2+\sum_{i=1}^3((\mathbf{k})_i-(\boldsymbol{\ell})_i)^2\right)^{-\frac{1}{2}},$$

and use the computation in Exercise 3.4.1.

Exercise 5.5.8 Even when S is finite, writing down a reasonably explicit expression for the solution to (5.2.15) is seldom easy and often impossible. Nonetheless, a little linear algebra often does quite a lot of good. Throughout this exercise, **Q** is a Q-matrix on the state space S, and it is assumed that the associated Markov process exists (i.e., does not explode) starting from any point.

- (a) If $\mathbf{u} \in \mathbb{C}^{\mathbb{S}}$ is a bounded, non-zero, right eigenvector for \mathbf{Q} with eigenvalue $\alpha \in \mathbb{C}$, show that the real part of α must be less than or equal to 0.
- (b) Assume that $N = \#S < \infty$ and that **Q** admits a complete set of linearly independent, right eigenvectors $\mathbf{u}_1, \ldots, \mathbf{u}_N \in \mathbb{C}^S$ with associated eigenvalues $\alpha_1, \ldots, \alpha_N$. Let **U** be the matrix whose *m*th column is \mathbf{u}_m , and show that $e^{t\mathbf{Q}} = \mathbf{U}\Lambda(t)\mathbf{U}^{-1}$, where $\Lambda(t)$ is the diagonal matrix whose *m* diagonal entry is $e^{t\alpha_m}$.

Exercise 5.5.9 Here is the continuous analog of the result in Exercise 4.2.7. Assume that *i* is *Q*-recurrent and $R_i \equiv -(\mathbf{Q})_{ii} > 0$, set

$$\hat{\mu}_j = \mathbb{E}\left[\int_0^{\sigma_i} \mathbf{1}_{\{j\}} (X(t)) dt \, \middle| \, X(0) = i\right] \in [0, \infty] \quad \text{for } j \in \mathbb{S},$$

and let $\hat{\mu}$ be the row vector given by $(\hat{\mu})_j = \hat{\mu}_j$ for each $j \in \mathbb{S}$.

(a) Show that $\hat{\mu}_i = \frac{1}{R_i}$, $\hat{\mu}_j < \infty$ for all $j \in \mathbb{S}$, and that $\hat{\mu}_j > 0$ if and only if $i \stackrel{\mathbf{Q}}{\leftrightarrow} j$. (b) Show that $\hat{\mu} = \hat{\mu} \mathbf{P}(t)$ for all t > 0.

Hint: Check that

$$\left(\hat{\boldsymbol{\mu}}\mathbf{P}(s)\right)_{j} = \mathbb{E}\left[\int_{s}^{s+\sigma_{i}} \mathbf{1}_{\{j\}} \left(X(t)\right) dt \mid X(0) = i\right]$$

and that

$$\mathbb{E}\left[\int_{\sigma_i}^{s+\sigma_i} \mathbf{1}_{\{j\}}(X(t)) dt \, \Big| \, X(0) = i\right] = \mathbb{E}\left[\int_0^s \mathbf{1}_{\{j\}}(X(t)) dt \, \Big| \, X(0) = i\right].$$

(c) In particular, if *i* is **Q**-positive recurrent and $C = \{j : i \stackrel{\mathbf{Q}}{\leftrightarrow} j\}$, show that $\hat{\boldsymbol{\mu}} = (\sum_{i} \hat{\mu}_{j}) \hat{\boldsymbol{\pi}}^{C}$. Equivalently,

$$\left(\hat{\pi}^{C}\right)_{j} = \frac{\mathbb{E}[\int_{0}^{\sigma_{i}} \mathbf{1}_{\{j\}}(X(t)) \, dt \, | \, X(0) = i]}{\mathbb{E}[\sigma_{i} \, | \, X(0) = i]}$$

Exercise 5.5.10 Having given the continuous-time analog of Exercise 4.2.7, we will now give the continuous-time analog of Exercise 4.2.9. For this purpose, again let *i* be a **Q**-recurrent element of \mathbb{S} with $R_i \equiv -(\mathbf{Q})_{ii} > 0$, and define the measure $\hat{\mu}$ accordingly, as in Exercise 5.5.9. Next, assume that $\hat{\nu} \in [0, \infty)^{\mathbb{S}}$ satisfies the conditions that:

$$(\hat{\boldsymbol{\nu}})_i > 0,$$
 $(\hat{\boldsymbol{\nu}})_j = 0$ if $\mathbb{P}(\sigma_j < \infty \mid X(0) = i) < 1,$ and $\hat{\boldsymbol{\nu}} \mathbf{Q} = \mathbf{0}$

in the sense that $R_j(\hat{\boldsymbol{v}})_j = \sum_{k \neq j} (\hat{\boldsymbol{v}})_k \mathbf{Q}_{kj}$ for all $j \in \mathbb{S}^4$. The goal is to show that $\hat{\boldsymbol{v}} = R_i(\hat{\boldsymbol{v}})_i \hat{\boldsymbol{\mu}}$. Equivalently,

$$(\hat{\boldsymbol{v}})_j = R_i(\hat{\boldsymbol{v}})_i \mathbb{E}\left[\int_0^{\sigma_i} \mathbf{1}_{\{j\}} (X(t)) dt \, \Big| \, X(0) = i\right] \text{ for all } j \in \mathbb{S}.$$

In particular, by Exercise 5.5.9, this will mean that $\hat{\mathbf{v}} = \hat{\mathbf{v}} \mathbf{P}(t)$ for all t > 0.

In the following, **P** and \mathfrak{R} are the canonical transition probability and rates determined by **Q**.

(a) Define the row vector \mathbf{v} so that $(\mathbf{v})_j = \frac{R_j(\hat{\mathbf{v}})_j}{R_i(\hat{\mathbf{v}})_i}$, and show that $\mathbf{v} = \mathbf{v}\mathbf{P}$. (b) By combining (a) with Exercise 4.2.7, show that

$$(\mathbf{v})_j = \mathbb{E}\left[\sum_{m=0}^{\rho_i - 1} \mathbf{1}_{\{j\}}(X_m) \, \middle| \, X_0 = i\right],$$

where $\{X_n : n \ge 0\}$ is the Markov chain with transition probability matrix **P**. (c) Show that

$$R_{j}\mathbb{E}\left[\int_{0}^{\sigma_{i}}\mathbf{1}_{\{j\}}(X(t))\,dt\,\Big|\,X(0)=i\right]=\mathbb{E}\left[\sum_{m=0}^{\rho_{i}-1}\mathbf{1}_{\{j\}}(X_{m})\,\Big|\,X_{0}=i\right],$$

and, after combining this with (b), arrive at the desired conclusion.

Exercise 5.5.11 Following the strategy used in Exercise 5.5.10, show that, under the hypotheses in Corollary 5.4.8, one has the following continuous version of the individual ergodic theorem:

$$\mathbb{P}\left(\lim_{T\to\infty}\left\|\mathbf{L}_{T}-\hat{\boldsymbol{\pi}}^{C}\right\|_{v}=0\right)=1,$$

⁴The reason why the condition $\hat{\nu} \mathbf{Q} = \mathbf{0}$ needs amplification is that, because **Q** has negative diagonal entries, possible infinities could cause ambiguities.

where $C = \{i : j \stackrel{\mathbf{Q}}{\leftrightarrow} i\}$ and \mathbf{L}_T is the *empirical measure* determined by

$$(\mathbf{L}_T)_i = \frac{1}{T} \int_0^T \mathbf{1}_{\{i\}} \big(X(t) \big) \, dt.$$

In addition, again following the strategy in Exercise 5.5.10, show that, for any initial distribution,

$$\mathbb{P}\left(\lim_{T \to \infty} \frac{1}{T} \int_0^T \mathbf{1}_{\{j\}} (X(t)) dt = 0\right) = 1$$

when j is not **Q**-positive recurrent.

Exercise 5.5.12 Given a Markov process $\{X(t) : t \ge 0\}$ with *Q*-matrix **Q** and an R > 0, one can produce a Markov process with *Q*-matrix $R\mathbf{Q}$ by speeding up the clock of $t \rightsquigarrow X(t)$ by a factor of *R*. That is, $\{X(Rt) : t \ge 0\}$ is a Markov process with *Q*-matrix $R\mathbf{Q}$. The purpose of this exercise is to carry out the analogous procedure, known in the literature as *random time change*, for variable rates. To be precise, let **P** be a transition probability matrix on \mathbb{S} with $(\mathbf{P})_{ii} = 0$ for all *i*. Choose and fix an $i \in \mathbb{S}$, let $\{X_n : n \ge 0\}$ be a Markov chain starting from *i* with transition probability matrix **P** and $\{N(t) : t \ge 0\}$ a simple Poisson process which is independent of $\sigma(\{X_n : n \ge 0\})$, and set $X^0(t) = X^0_{N(t)}$ for $t \ge 0$. That is, $\{X^0(t) : t \ge 0\}$ is a Markov process with *Q*-matrix (**P** - **I**) starting from *i*. Finally, let \Re be a set of positive rates, and take $\mathbf{Q} = \mathbf{R}(\mathbf{P} - \mathbf{I})$ accordingly.

(a) Define

$$A(t) = \int_0^t \frac{1}{R_{X^0(\tau)}} d\tau \quad \text{for } t \in [0, \infty).$$

observe that $t \rightsquigarrow A(t)$ is strictly increasing, set $A(\infty) = \lim_{t \neq \infty} A(t) \in (0, \infty]$, and use $s \in [0, A(\infty)) \longmapsto A^{-1}(s) \in [0, \infty)$ to denote the inverse of $t \rightsquigarrow A(t)$. Show that

$$A^{-1}(s) = \int_0^s R_{X^0(A^{-1}(\sigma))} d\sigma, \quad s \in [0, A(\infty)).$$

(b) Set $X(s) = X^0(A^{-1}(s))$ for $s \in [0, A(\infty))$. Define $J_0^0 = 0 = J_0$ and, for $n \ge 1$, J_n^0 and J_n to be the times of the *n*th jump of $t \rightsquigarrow X^0(t)$ and $s \rightsquigarrow X(s)$, respectively. After noting that $J_n = A(J_n^0)$, conclude that, for each $n \ge 1$, s > 0, and $j \in \mathbb{S}$,

$$\mathbb{P}(J_n - J_{n-1} > s \& X(J_n) = j \mid X(\sigma), \sigma \in [0, J_n))$$

= $e^{-tR_{X(J_{n-1})}}(\mathbf{P})_{X(J_{n-1})j}$ on $\{J_{n-1} < \infty\}$.

(c) Show that the explosion time for the Markov process starting from *i* with Q-matrix **Q** has the same distribution as $A(\infty)$. In particular, if $A(\infty) = \infty$ with probability 1, use (b) to conclude that $\{X(s) : s \ge 0\}$ is a Markov process starting from *i* with Q-matrix **Q**.

(d) As a consequence of these considerations, show that if the Markov process corresponding to **Q** does not explode, then neither does the one corresponding to **Q'**, where **Q'** is related to **Q** by $(\mathbf{Q}')_{ij} = \alpha_i(\mathbf{Q})_{ij}$ and the $\{\alpha_i : i \in \mathbb{S}\}$ is a bounded subset of $(0, \infty)$.

Chapter 6 Reversible Markov Processes

This chapter is devoted to the study of a class of Markov processes that admit an initial distribution with respect to which they are *reversible* in the sense that, on every time interval, the distribution of the process is the same when it is run backwards as when it is run forwards. That is, for any $n \ge 1$ and $(i_0, \ldots, i_n) \in E^{n+1}$, in the discrete time setting,

$$\mathbb{P}(X_m = i_m \text{ for } 0 \le m \le n) = \mathbb{P}(X_{n-m} = i_m \text{ for } 0 \le m \le n)$$
(6.0.1)

and in the continuous time setting,

$$\mathbb{P}(X(t_m) = i_m \text{ for } 0 \le m \le n) = \mathbb{P}(X(t_n - t_m) = i_m \text{ for } 0 \le m \le n)$$
(6.0.2)

whenever $0 = t_0 < \cdots < t_n$. Notice that the initial distribution of such a process is necessarily stationary. Indeed, depending on whether the setting is that of discrete or continuous time, we have

$$\mathbb{P}(X_0 = i \& X_n = j) = \mathbb{P}(X_n = i \& X_0 = j) \quad \text{or}$$
$$\mathbb{P}(X(0) = i \& X(t) = j) = \mathbb{P}(X(t) = i \& X(0) = j),$$

from which stationarity follows after one sums over j. In fact, what the preceding argument reveals is that reversibility says that the joint distribution of, depending on the setting, (X_0, X_n) or (X(0), X(t)) is the same as that of (X_n, X_0) or (X(t), X(0)). This should be contrasted with the stationarity which gives equality only for the marginal distribution of the first components of these.

In view of the preceding, one should suspect that reversible Markov processes have ergodic properties that are better than those of general stationary processes, and in this chapter we will show that this suspicion is justified.

6.1 Reversible Markov Chains

In this section we will discuss irreducible, reversible Markov chains. Because the initial distribution of such a chain is stationary, we know (cf. Theorem 4.1.10) that the chain must be positive recurrent and that the initial distribution must be the probability vector (cf. (4.1.5) and (4.1.9)) $\boldsymbol{\pi} = \boldsymbol{\pi}^{\mathbb{S}}$ whose *i*th component is $(\boldsymbol{\pi})_i \equiv \mathbb{E}[\rho_i|X_0 = i]^{-1}$. Thus, if **P** is the transition probability matrix, then, by taking n = 1 in (6.0.1), we see that

$$(\boldsymbol{\pi})_i(\mathbf{P})_{ij} = \mathbb{P}(X_0 = i \& X_1 = j) = \mathbb{P}(X_0 = j \& X_1 = i) = (\boldsymbol{\pi})_j(\mathbf{P})_{ji}$$

That is, **P** satisfies¹

$$(\boldsymbol{\pi})_i(\mathbf{P})_{ii} = (\boldsymbol{\pi})_i(\mathbf{P})_{ii}, \text{ the detailed balance condition.}$$
 (6.1.1)

Conversely, (6.1.1) implies reversibility. To see this, one works by induction on $n \ge 1$ to check that

$$\boldsymbol{\pi}_{i_0}(\mathbf{P})_{i_0i_1}\cdots(\mathbf{P})_{i_{n-1}i_n}=\boldsymbol{\pi}_{i_n}(\mathbf{P})_{i_ni_{n-1}}\cdots(\mathbf{P})_{i_1i_0},$$

which is equivalent to (6.0.1).

6.1.1 Reversibility from Invariance

As we have already seen, reversibility implies invariance, and it should be clear that the converse is false. On the other hand, there are two canonical ways in which one can pass from an irreducible transition probability **P** with stationary distribution π to a transition probability for which π is reversible. Namely, for such a **P**, $(\pi)_i > 0$ for all $i \in \mathbb{S}$, and so we can define the *adjoint* \mathbf{P}^{\top} of **P** so that

$$\left(\mathbf{P}^{\top}\right)_{ij} = \frac{(\boldsymbol{\pi})_j (\mathbf{P})_{ji}}{(\boldsymbol{\pi})_i}.$$
(6.1.2)

Obviously, π is reversible for **P** if and only if $\mathbf{P} = \mathbf{P}^{\top}$. Moreover, because $\pi \mathbf{P} = \pi$, \mathbf{P}^{\top} is again a transition probability. In addition, one can easily verify that both

$$\frac{\mathbf{P} + \mathbf{P}^{\top}}{2} \quad \text{and} \quad \mathbf{P}^{\top} \mathbf{P} \tag{6.1.3}$$

are transition probabilities that are reversible with respect to π . As is explained in Exercise 6.6.7 below, each of these constructions has its own virtue.

¹The reader who did Exercise 2.4.1 should recognize that the condition below is precisely the same as the statement that $\mathbf{P} = \mathbf{P}^{\top}$. In particular, if one knows the conclusion of that exercise, then one has no need for the discussion which follows.

6.1.2 Measurements in Quadratic Mean

For reasons which will become increasingly clear, it turns out that we will want here to measure the size of functions using a Euclidean norm rather than the uniform norm $||f||_u$. Namely, we will use the norm

$$\|f\|_{2,\pi} \equiv \sqrt{\langle |f|^2 \rangle_{\pi}} \quad \text{where } \langle g \rangle_{\pi} \equiv \sum_{i} g(i)(\pi)_{i} \tag{6.1.4}$$

is the expected value of g with respect to π . Because $(\pi)_i > 0$ for each $i \in \mathbb{S}$, it is clear that $||f||_{2,\pi} = 0 \iff f \equiv 0$. In addition, if we define the *inner product* $\langle f, g \rangle_{\pi}$ to be $\langle fg \rangle_{\pi}$, then, for any t > 0,

$$0 \le \|tf \pm t^{-1}f\|_{2,\pi}^2 = t^2 \|f\|_{2,\pi}^2 \pm 2\langle f, g \rangle_{\pi} + t^{-2} \|g\|_{2,\pi}^2,$$

and so $|\langle f, g \rangle_{\pi}| \le t^2 ||f||_{2,\pi}^2 + t^{-2} ||g||_{2,\pi}^2$ for all t > 0. To get the best estimate, we minimize the right hand side with respect to t > 0. When either f or g is identically 0, then we see that $\langle f, g \rangle_{\pi} = 0$ by letting $t \to \infty$ or $t \to 0$. If neither f nor g vanishes identically, we can do best by taking $t = (\frac{||g||_{2,\pi}}{||f||_{2,\pi}})^{\frac{1}{4}}$. Hence, in any case, we arrive at *Schwarz's inequality*

$$|\langle f, g \rangle_{\pi}| \le ||f||_{2,\pi} ||g||_{2,\pi}.$$
 (6.1.5)

Given Schwarz's inequality, we know

$$\|f + g\|_{2,\pi}^2 = \|f\|_{2,\pi}^2 + 2\langle f, g \rangle_{\pi} + \|g\|_{2,\pi}^2 \le \left(\|f\|_{2,\pi} + \|g\|_{2,\pi}\right)^2.$$

That is, we have the *triangle inequality*:

$$\|f + g\|_{2,\pi} \le \|f\|_{2,\pi} + \|g\|_{2,\pi}.$$
(6.1.6)

Thus, if $L^2(\pi)$ denotes the space of f for which $||f||_{2,\pi} < \infty$, then $L^2(\pi)$ is a linear space for which $(f, g) \rightsquigarrow ||f - g||_{2,\pi}$ is a metric. In fact, this metric space is complete since, if $\lim_{m\to\infty} \sup_{n>m} ||f_n - f_m||_{2,\pi} = 0$, then $\{f_n(i) : n \ge 0\}$ is Cauchy convergent in \mathbb{R} for each $i \in \mathbb{S}$, and therefore there exists a limit function f such that $f_n(i) \longrightarrow f(i)$ for each $i \in \mathbb{S}$. Moreover, by Fatou's lemma,

$$||f - f_m||_{2,\pi} \le \lim_{n \to \infty} ||f_n - f_m||_{2,\pi} \longrightarrow 0 \text{ as } m \to \infty.$$

In this chapter, we will use the notation

$$\mathbf{P}f(i) = \sum_{j \in \mathbb{S}} f(j)(\mathbf{P})_{ij} = (\mathbf{P}\mathbf{f})_i,$$

where **f** is the column vector determined by f. When f is bounded, it is clear that $\mathbf{P}f(i)$ is well-defined and that $\|\mathbf{P}f\|_{u} \le \|f\|_{u}$, where $\|g\|_{u} \equiv \sup_{i \in \mathbb{S}} |g(i)| = \|\mathbf{g}\|_{u}$

when **g** is the column vector determined by *g*. We next show that, even when $f \in L^2(\pi)$, $\mathbf{P}f(i)$ is well-defined for each $i \in \mathbb{S}$ and that **P** is a contraction in $L^2(\pi)$: $\|\mathbf{P}f\|_{2,\pi} \leq \|f\|_{2,\pi}$. To check that $\mathbf{P}f(i)$ is well-defined, we will show that the series $\sum_{j \in \mathbb{S}} f(j)(\mathbf{P})_{ij}$ is absolutely convergent. But, by the form of Schwarz's inequality in Exercise 1.3.1,

$$\sum_{j\in\mathbb{S}} \left| f(j) \right| (\mathbf{P})_{ij} = \sum_{j\in\mathbb{S}} \left| f(y) \right| \sqrt{(\boldsymbol{\pi})_i} \frac{(\mathbf{P})_{ij}}{\sqrt{(\boldsymbol{\pi})_i}} \le \| f \|_{2,\boldsymbol{\pi}} \left(\sum_{j\in\mathbb{S}} \frac{(\mathbf{P})_{ij}^2}{(\boldsymbol{\pi})_j} \right)^{\frac{1}{2}},$$

and, by (6.1.1),

$$\sum_{j\in\mathbb{S}}\frac{(\mathbf{P})_{ij}^2}{(\boldsymbol{\pi})_j} = \frac{1}{(\boldsymbol{\pi})_i}\sum_{j\in\mathbb{S}}(\mathbf{P})_{ij}(\mathbf{P})_{ji} = \frac{(\mathbf{P}^2)_{ii}}{((\boldsymbol{\pi})_i)} < \infty.$$

As for the estimate $\|\mathbf{P}f\|_{2,\pi} \le \|f\|_{2,\pi}$, we use Exercise 6.6.2 below together with $\sum_{j\in\mathbb{S}}(\mathbf{P})_{ij} = 1$ to see that $(\mathbf{P}f(i))^2 \le \mathbf{P}f^2(i)$ for each *i*. Thus, since π is **P**-stationary,

$$\|\mathbf{P}f\|_{2,\pi} \le \|f\|_{2,\pi}. \tag{6.1.7}$$

An important consequence of (6.1.7) is the fact that, in general (cf. (2.2.4)),

$$\lim_{n \to \infty} \left\| \mathbf{A}_n f - \langle f \rangle_{\boldsymbol{\pi}} \right\|_{2, \boldsymbol{\pi}} = 0 \quad \text{for all } f \in L^2(\boldsymbol{\pi})$$
(6.1.8)

and

$$\mathbf{P} \text{ is aperiodic} \implies \lim_{n \to \infty} \left\| \mathbf{P}^n f - \langle f \rangle_{\pi} \right\|_{2,\pi} = 0 \quad \text{for all } f \in L^2(\pi). \quad (6.1.9)$$

To see these, first observe that, by (4.1.11), (4.1.15), and Lebesgue's dominated convergence theorem, there is nothing to do when f vanishes off of a finite set. Thus, if $\{F_N : N \ge 1\}$ is an exhaustion of \mathbb{S} by finite sets and if, for $f \in L^2(\pi)$, $f_N \equiv \mathbf{1}_{F_N} f$, then, for each $N \in \mathbb{Z}^+$,

$$\begin{aligned} \|\mathbf{A}_n f - \langle f \rangle_{\pi} \|_{2,\pi} &\leq \|\mathbf{A}_n (f - f_N)\|_{2,\pi} + \|\mathbf{A}_n f_N - \langle f_N \rangle_{\pi} \|_{2,\pi} + \langle |f_N - f| \rangle_{\pi} \\ &\leq 2 \|f - f_N\|_{2,\pi} + \|\mathbf{A}_n f_N - \langle f_N \rangle_{\pi} \|_{2,\pi}, \end{aligned}$$

where, in the passage to the second line, we have used $\|\mathbf{A}_n g\|_{2,\pi} \le \|g\|_{2,\pi}$, which follows immediately from (6.1.7). Therefore, for each *N*,

$$\overline{\lim_{n\to\infty}} \left\| \mathbf{A}_n f - \langle f \rangle_{\pi} \right\|_{2,\pi} \leq 2 \| f - f_N \|_{2,\pi},$$

which gives (6.1.8) when $N \to \infty$. The argument for (6.1.9) is essentially the same, and, as is shown in Exercise 6.6.4 below, all these results hold even in the non-reversible setting.

Finally, (6.1.1) leads to

$$\langle g, \mathbf{P}f \rangle_{\pi} = \sum_{(i,j)} (\pi)_i g(i)(\mathbf{P})_{ij} f(j) = \sum_{(i,j)} (\pi)_j f(i)(\mathbf{P})_{ji} g(j) = \langle \mathbf{P}g, f \rangle_{\pi}.$$

In other words, **P** is *symmetric* on $L^2(\pi)$ in the sense that

$$\langle g, \mathbf{P}f \rangle_{\pi} = \langle \mathbf{P}g, f \rangle_{\pi} \quad \text{for } (f, g) \in \left(L^2(\pi)\right)^2.$$
 (6.1.10)

6.1.3 The Spectral Gap

Equation (6.1.7) combined with (6.1.10) say that **P** is a *self-adjoint contraction* on the Hilbert space² $L^2(\pi)$. For the reader who is unfamiliar with these concepts at this level of abstraction, think about the case when $\mathbb{S} = \{1, ..., N\}$. Then the space of functions $f : \mathbb{S} \longrightarrow \mathbb{R}$ can be identified with \mathbb{R}^N . (Indeed, we already made this identification when we gave, in Sect. 2.1.3, the relation between functions and column vectors.) After making this identification, the inner product on $L^2(\pi)$ becomes the inner product $\sum_{1}^{N} (\pi)_i(\mathbf{v})_i(\mathbf{w})_i$ for column vectors \mathbf{v} and \mathbf{w} in \mathbb{R}^N . Hence (6.1.7) says that the matrix **P** acts as a symmetric, contraction on \mathbb{R}^N with respect to this inner product. Alternatively, if $\tilde{\mathbf{P}} \equiv \boldsymbol{\Pi}^{\frac{1}{2}} \mathbf{P} \boldsymbol{\Pi}^{-\frac{1}{2}}$, where $\boldsymbol{\Pi}$ is the diagonal matrix whose *i*th diagonal entry is $(\pi)_i$, then, by (6.1.1), $\tilde{\mathbf{P}}$ is symmetric with respect to the standard inner product $(\mathbf{v}, \mathbf{w})_{\mathbb{R}^N} \equiv \sum_{1}^{N} (\mathbf{v})_i(\mathbf{w})_i$ on \mathbb{R}^N . Moreover, because, by (6.1.7),

$$\|\tilde{\mathbf{P}}\mathbf{f}\|_{\mathbb{R}^{N}}^{2} \equiv (\tilde{\mathbf{P}}\mathbf{f}, \tilde{\mathbf{P}}\mathbf{f})_{\mathbb{R}^{N}} = \sum_{i} (\pi)_{i} \left(\mathbf{P}\boldsymbol{\Pi}^{-\frac{1}{2}}\mathbf{f}\right)_{i}^{2} = \|\mathbf{P}g\|_{2,\pi}^{2} \le \|g\|_{2,\pi}^{2} = \|\mathbf{f}\|_{\mathbb{R}^{N}}^{2},$$

where g is the function determined by the column vector $\mathbf{\Pi}^{-\frac{1}{2}}\mathbf{f}$, we see that, as an operator on \mathbb{R}^N , $\tilde{\mathbf{P}}$ is length contracting. Now, by the standard theory of symmetric matrices on \mathbb{R}^N , we know that $\tilde{\mathbf{P}}$ admits eigenvalues $1 \ge \lambda_1 \ge \cdots \ge \lambda_N$ ≥ -1 with associated eigenvectors $(\mathbf{e}_1, \dots, \mathbf{e}_N)$ that are orthonormal for $(\cdot, \cdot)_{\mathbb{R}^N}$: $(\mathbf{e}_k, \mathbf{e}_\ell)_{\mathbb{R}^N} = \delta_{k,\ell}$. Moreover, because $\sqrt{(\pi)_i} = \sum_{j=1}^N (\tilde{\mathbf{P}})_{ij} \sqrt{(\pi)_j}$, we know that $\lambda_1 = 1$ and can take $(\mathbf{e}_1)_i = \sqrt{(\pi)_i}$. Finally, by setting $\mathbf{g}_\ell = (\mathbf{\Pi})^{-\frac{1}{2}} \mathbf{e}_\ell$ and letting g_ℓ be the associated function on \mathbb{S} , we see that $\mathbf{P}g_\ell = \lambda_\ell g_\ell$, $g_1 \equiv 1$, and $\langle g_k, g_\ell \rangle_{\pi} = \delta_{k,\ell}$. To summarize, when \mathbb{S} has N elements, we have shown that \mathbf{P} on $L^2(\pi)$ has eigenvalues $1 = \lambda_1 \ge \cdots \ge \lambda_N \ge -1$ with corresponding eigenfunctions g_1, \dots, g_N which are orthonormal with respect to $\langle \cdot, \cdot \rangle_{\pi}$. Of course, since $L^2(\pi)$ has dimension N and, by ortho-normality, the g_ℓ 's are linearly independent,

²A Hilbert space is a vector space equipped with an inner product which determines a norm for which the associated metric is complete.

 (g_1, \ldots, g_N) is an orthonormal basis in $L^2(\pi)$. In particular,

$$\mathbf{P}^{n} f - \langle f \rangle_{\boldsymbol{\pi}} = \sum_{\ell=2}^{N} \lambda_{\ell}^{n} \langle f, g_{\ell} \rangle_{\boldsymbol{\pi}} g_{\ell} \quad \text{for all } n \ge 0 \text{ and } f \in L^{2}(\boldsymbol{\pi}), \qquad (6.1.11)$$

and so

$$\|\mathbf{P}^{n} f - \langle f \rangle_{\pi} \|_{2,\pi}^{2} = \sum_{\ell=2}^{N} \lambda_{\ell}^{2n} \langle f, g_{\ell} \rangle_{\pi}^{2} \le (1-\beta)^{2n} \|f - \langle f \rangle_{\pi} \|_{2,\pi}^{2},$$

where $\beta = (1 - \lambda_2) \wedge (1 + \lambda_N)$ is the *spectral gap* between $\{-1, 1\}$ and $\{\lambda_2, \dots, \lambda_N\}$. In other words,

$$\|\mathbf{P}^{n} f - \langle f \rangle_{\boldsymbol{\pi}} \|_{2,\boldsymbol{\pi}} \leq (1-\beta)^{n} \| f - \langle f \rangle_{\boldsymbol{\pi}} \|_{L^{2}(\boldsymbol{\pi})}$$

for all $n \geq 0$ and $f \in L^{2}(\boldsymbol{\pi})$. (6.1.12)

When S is not finite, it will not be true in general that one can find an orthonormal basis of eigenfunctions for **P**. Instead, the closest approximation to the preceding development requires a famous result, known as the Spectral Theorem (cf. Sect. 107 in [7]), about bounded, symmetric operators on a Hilbert space. Nonetheless, seeing as it is the estimate (6.1.12) in which we are most interested, we can get away without having to invoke the Spectral Theorem. To be more precise, observe that (6.1.12) holds when

$$\beta \equiv 1 - \sup\{\|\mathbf{P}f - \langle f \rangle_{\pi}\|_{2,\pi} : f \in L^{2}(\pi) \text{ with } \|f\|_{2,\pi} = 1\}.$$
(6.1.13)

To see this, first note that, when β is given by (6.1.13),

$$\begin{aligned} \left\| \mathbf{P}f - \langle f \rangle_{\pi} \right\|_{2,\pi} &= \|f\|_{2,\pi} \left\| \mathbf{P}\left(\frac{f}{\|f\|_{2,\pi}}\right) - \left\langle \frac{f}{\|f\|_{2,\pi}} \right\rangle_{\pi} \right\|_{2,\pi} \\ &\leq (1-\beta) \|f\|_{2,\pi} \quad \text{when } f \in L^{2}(\pi) \setminus \{0\}, \end{aligned}$$

and that $\|\mathbf{P}f - \langle f \rangle_{\pi}\|_{2,\pi} \le (1 - \beta) \|f\|_{2,\pi}$ trivially when f = 0. Next, because $\langle \mathbf{P}^n f \rangle_{\pi} = \langle f \rangle_{\pi}$ for all $n \ge 0$,

$$\begin{aligned} \left\| \mathbf{P}^{n+1} f - \langle f \rangle_{\pi} \right\|_{2,\pi} &= \left\| \mathbf{P} \left(\mathbf{P}^{n} f - \langle \mathbf{P}^{n} f \rangle_{\pi} \right) \right\|_{2,\pi} \\ &\leq (1-\beta) \left\| \mathbf{P}^{n} f - \langle \mathbf{P}^{n} f \rangle_{\pi} \right\|_{2,\pi} = (1-\beta) \left\| \mathbf{P}^{n} f - \langle f \rangle_{\pi} \right\|_{2,\pi}. \end{aligned}$$

By induction on *n*, we get $\|\mathbf{P}^n f - \langle f \rangle_{\pi}\|_{2,\pi} \le (1-\beta)^n \|f\|_{2,\pi}$ for all $n \ge 0$ and $f \in L^2(\pi)$. Hence, if $f \in L^2(\pi)$ and $\bar{f} = f - \langle f \rangle_{\pi}$, then

$$\|\mathbf{P}^{n}f - \langle f \rangle_{\pi}\|_{2,\pi} = \|\mathbf{P}^{n}\bar{f}\|_{2,\pi} \le (1-\beta)^{n}\|\bar{f}\|_{2,\pi} = (1-\beta)^{n}\|f - \langle f \rangle_{\pi}\|_{2,\pi}.$$

That is, (6.1.12) holds with the β in (6.1.13). Observe that when S is finite the β in (6.1.13) coincides with the one in the preceding paragraph. Hence, we have made a

first step toward generalizing the contents of that paragraph to situations in which the analog of (6.1.11) does not exist.

6.1.4 Reversibility and Periodicity

Clearly, the constant β in (6.1.13) can be as small as 0, in which case (6.1.12) tells us nothing. There are three ways in which this might happen. One way is that there exist an $f \in L^2(\pi)$ with the property that $||f||_{2,\pi} = 1$, $\langle f \rangle_{\pi} = 0$, and $\mathbf{P}f = f$. However, irreducibility rules out the existence of such an f. Indeed, because of irreducibility, we would have (cf. (6.1.8)) the contradiction that $0 = \langle f \rangle_{\pi} = \lim_{n \to \infty} (\mathbf{A}_n \mathbf{f})_i = f(i)$ for all $i \in \mathbb{S}$ and that $f(i) \neq 0$ for some $i \in \mathbb{S}$. Thus, we can ignore this possibility because it never occurs. A second possibility is that there exists an $f \in L^2(\pi)$ with $||f||_{2,\pi} = 1$ such that $\mathbf{P}f = -f$. In fact, if f is such a function, then $\langle f \rangle_{\pi} = \langle \mathbf{P}f \rangle_{\pi} = -\langle f \rangle_{\pi}$, and so $\langle f \rangle_{\pi} = 0$. Hence, we would have that $||\mathbf{P}f - \langle f \rangle_{\pi}||_{2,\pi} = 1$ and therefore that $\beta = 0$. The third possibility is that there is no non-zero solution to $\mathbf{P}f = -f$ but that, nonetheless, there exists a sequence $\{f_n\}_1^{\infty} \subseteq L^2(\pi)$ with $||f_n||_{2,\pi} = 1$ and $\langle f_n \rangle_{\pi} = 0$ such that $||\mathbf{P}f_n||_{2,\pi}$ tends to 1.

Because the analysis of this last possibility requires the Spectral Theorem, we will not deal with it. However, as the next theorem shows, the second possibility has a pleasing and simple probabilistic interpretation. See Exercise 6.6.5 below for an extension of these considerations to non-reversible **P**'s.

Theorem 6.1.14 If **P** is an irreducible transition probability for which there is a reversible initial distribution, which is necessarily π , then the period of **P** is either 1 or 2. Moreover, the period is 2 if and only if there exists an $f \in L^2(\pi) \setminus \{0\}$ for which $f = -\mathbf{P}f$.

Proof We begin by showing that the period *d* must be less than or equal to 2. To this end, remember that, because of irreducibility, $(\boldsymbol{\pi})_i > 0$ for all *i*'s. Hence, the detailed balance condition, (6.1.1), implies that $(\mathbf{P})_{ij} > 0 \iff (\mathbf{P})_{ji} > 0$. In particular, since, for each *i*, $(\mathbf{P})_{ij} > 0$ for some *j* and therefore $(\mathbf{P}^2)_{ii} = \sum_j (\mathbf{P})_{ij} (\mathbf{P})_{ji} > 0$, we see that the period must divide 2.

To complete the proof at this point, first suppose that d = 1. If $f \in L^2(\pi)$ satisfies $f = -\mathbf{P}f$, then, as noted before, $\langle f \rangle_{\pi} = 0$, and yet, because of aperiodicity and (6.1.9), $\lim_{n\to\infty} \mathbf{P}^n f(i) = \langle f \rangle_{\pi} = 0$ for each $i \in \mathbb{S}$. Since $f = \mathbf{P}^{2n} f$ for all $n \ge 0$, this means that $f \equiv 0$. Conversely, if d = 2, take \mathbb{S}_0 and \mathbb{S}_1 accordingly, as in Sect. 4.1.7, and consider $f = \mathbf{1}_{\mathbb{S}_0} - \mathbf{1}_{\mathbb{S}_1}$. Because of (4.1.19), $\mathbf{P}f = -f$, and clearly $\|f\|_{2,\pi} = 1$.

As an immediate corollary of the preceding, we can give the following graph theoretic picture of aperiodicity for irreducible, reversible Markov chains. Namely, if we use **P** to define a graph structure in which the elements of S are the "vertices"

and an "edge" between *i* and *j* exists if and only if $(\mathbf{P})_{ij} > 0$, then the first part of Theorem 6.1.14, in combination with the considerations in Sect. 4.1.7, says that the resulting graph is *bipartite* (i.e., splits into two parts in such a way that all edges run from one part to the other) if and only if the chain fails to be aperiodic, and the second part says that this is possible if and only if there exists an $f \in L^2(\pi) \setminus \{0\}$ satisfying $\mathbf{P}f = -f$.

6.1.5 Relation to Convergence in Variation

Before discussing methods for finding or estimating the β in (6.1.13), it might be helpful to compare the sort of convergence result contained in (6.1.12) to the sort of results we have been getting heretofore. To this end, first observe that

$$\left\|\mathbf{P}^{n}f - \langle f \rangle_{\boldsymbol{\pi}}\right\|_{2,\boldsymbol{\pi}} \leq \left\|\mathbf{P}^{n}f - \langle f \rangle_{\boldsymbol{\pi}}\right\|_{\mathbf{u}} \leq \sup_{i} \left\|\boldsymbol{\delta}_{i}\mathbf{P}^{n} - \boldsymbol{\pi}\right\|_{\mathbf{v}} \|f\|_{\mathbf{u}}$$

In particular, if one knows, as one does when Theorem 2.2.1 applies, that

$$\sup_{i} \left\| \boldsymbol{\delta}_{i} \mathbf{P}^{n} - \boldsymbol{\pi} \right\|_{\mathsf{v}} \le C (1 - \epsilon)^{n} \tag{(*)}$$

for some $C < \infty$ and $\epsilon \in (0, 1]$, then one has that

$$\left\|\mathbf{P}^{n}f-\langle f\rangle_{\boldsymbol{\pi}}\right\|_{2,\boldsymbol{\pi}}\leq C(1-\epsilon)^{n}\|f\|_{u},$$

which looks a lot like (6.1.12). Indeed, the only difference is that on the right hand side of (6.1.12), C = 1 and the norm is $||f||_{2,\pi}$ instead of $||f||_u$. Thus, one should suspect that (*) implies that the β in (6.1.12) is at least as large as the ϵ in (*). If, as will always be the case when S is finite, there exists a $g \in L^2(\pi)$ with the properties that $||g||_{2,\pi} = 1$, $\langle g \rangle_{\pi} = 0$, and either $\mathbf{P}g = (1 - \beta)g$ or $\mathbf{P}g = -(1 - \beta)g$, this suspicion is easy to verify. Namely, set $f = g\mathbf{1}_{[-R,R]}(g)$ where R > 0 is chosen so that $a \equiv \langle f, g \rangle_{\pi} \ge \frac{1}{2}$, and set $\bar{f} = f - \langle f \rangle_{\pi}$. Then, after writing $\bar{f} = ag + (\bar{f} - ag)$ and noting that $\langle g, \bar{f} - ag \rangle_{\pi} = 0$, we see that, for any $n \ge 0$,

$$\|\mathbf{P}^{n}\bar{f}\|_{2,\pi}^{2} = a^{2}(1-\beta)^{2n} \pm 2a(1-\beta)^{n} \langle g, \mathbf{P}^{n}(\bar{f}-ag) \rangle_{\pi} + \|\mathbf{P}^{n}(\bar{f}-ag)\|_{2,\pi}^{2}$$
$$\geq \frac{1}{4}(1-\beta)^{2n},$$

since

$$\pm \langle g, \mathbf{P}^n(\bar{f} - ag) \rangle_{\pi} = \langle \mathbf{P}^n g, \bar{f} - ag \rangle_{\pi} = (1 - \beta)^n \langle g, \bar{f} - ag \rangle_{\pi} = 0.$$

On the other hand, $\|\mathbf{P}^n \bar{f}\|_{2,\pi}^2 \leq C^2 (1-\epsilon)^{2n} \|\bar{f}\|_u^2 \leq (CR)^2 (1-\epsilon)^{2n}$. Thus, $\frac{1}{4}(1-\beta)^{2n} \leq (CR)^2 (1-\epsilon)^{2n}$ for all $n \geq 0$, which is possible only if $\beta \geq \epsilon$. When

no such g exists, the same conclusion holds, only one has to invoke the Spectral Theorem in order to arrive at it.

As the preceding shows, uniform estimates on the variation distance between $\mu \mathbf{P}^n$ and π imply estimates like the one in (6.1.12). However, going in the opposite direction is not always possible. To examine what can be done, let a probability vector $\boldsymbol{\mu}$ be given, and define f so that $f(i) = \frac{(\mu)_i}{(\pi)_i}$. Then, since $\langle f \rangle_{\pi} = 1$ and $\langle g \rangle_{\pi}^2 \leq \langle g^2 \rangle_{\pi}$ for all $g \in L^2(\pi)$,

$$\begin{split} \left\|\boldsymbol{\mu}\mathbf{P}^{n}-\boldsymbol{\pi}\right\|_{v} &= \sum_{j} \left|\left(\boldsymbol{\mu}\mathbf{P}^{n}\right)_{j}-(\boldsymbol{\pi})_{j}\right| = \sum_{j} \left|\left\langle f,\mathbf{P}^{n}\mathbf{1}_{\{j\}}\right\rangle_{\boldsymbol{\pi}}-(\boldsymbol{\pi})_{j}\right| \\ &= \sum_{j} \left|\left\langle \mathbf{P}^{n}f-\mathbf{1},\mathbf{1}_{\{j\}}\right\rangle_{\boldsymbol{\pi}}\right| = \sum_{j} \left|\left\langle \mathbf{P}^{n}f-\left\langle f\right\rangle_{\boldsymbol{\pi}},\mathbf{1}_{\{j\}}\right\rangle_{\boldsymbol{\pi}}\right| \\ &= \sum_{j} (\boldsymbol{\pi})_{j} \left|\mathbf{P}^{n}f(j)-\left\langle f\right\rangle_{\boldsymbol{\pi}}\right| \leq \left(\sum_{j} (\boldsymbol{\pi})_{j} \left|\mathbf{P}^{n}f(j)-\left\langle f\right\rangle_{\boldsymbol{\pi}}\right|^{2}\right)^{\frac{1}{2}} \\ &= \left\|\mathbf{P}^{n}f-\left\langle f\right\rangle_{\boldsymbol{\pi}}\right\|_{2,\boldsymbol{\pi}}, \end{split}$$

where, in the second to last line, I used Exercise 6.6.2. Hence, (6.1.12) implies that

$$\|\boldsymbol{\mu}\mathbf{P}^{n} - \boldsymbol{\pi}\|_{v} \leq \left(\sum_{i} \frac{(\boldsymbol{\mu})_{i}^{2}}{(\boldsymbol{\pi})_{i}} - 1\right)^{\frac{1}{2}} (1 - \beta)^{n}.$$
 (6.1.15)

In the case when S is finite, and therefore there exists an $\lambda \in (0, 1)$ for which $(\pi)_i \ge \lambda$, (6.1.15) yields the Doeblin type estimate

$$\|\boldsymbol{\mu}\mathbf{P}^n-\boldsymbol{\pi}\|_{\mathrm{v}}\leq \left(\frac{1-\lambda}{\lambda}\right)^{\frac{1}{2}}(1-\beta)^n.$$

However, when S is infinite, (6.1.15), as distinguished from Doeblin, does not give a rate of convergence which is independent of μ . In fact, unless $\sum_i \frac{\langle \mu \rangle_i^2}{\langle \pi \rangle_i} < \infty$, it gives no information at all. Thus one might ask why we are considering estimates like (6.1.15) when Doeblin does as well and sometimes does better. The answer is that, although Doeblin may do well when it works, it seldom works when S is infinite and, even when S is finite, it usually gives a far less than optimal rate of convergence. See, for example, Exercise 6.6.15 below where an example is given of a situation to which Doeblin's condition does not apply but the considerations here do.

6.2 Dirichlet Forms and Estimation of β

Our purpose in this section will be to find methods for estimating the optimal (i.e., largest) value of β for which (6.1.12) holds. Again, we will assume that the chain is reversible and irreducible. Later, we will add the assumption that it is aperiodic.

6.2.1 The Dirichlet Form and Poincaré's Inequality

Our first step requires us to find other expressions the right hand side of (6.1.13). To this end, begin with the observation that

$$1 - \beta = \sup \{ \|\mathbf{P}f\|_{2,\pi} : f \in L^2_0(\pi) \text{ with } \|f\|_{2,\pi} = 1 \}$$

where $L^2_0(\pi) \equiv \{ f \in L^2(\pi) : \langle f \rangle_{\pi} = 0 \}.$ (6.2.1)

It is obvious that the supremum on the right hand side of (6.1.13) dominates the supremum on the right above. On the other hand, if $f \in L^2(\pi)$ with $||f||_{2,\pi} = 1$, then either $||f - \langle f \rangle_{\pi}||_{2,\pi} = 0$ and therefore $||\mathbf{P}f - \langle f \rangle_{\pi}||_{2,\pi} = 0$, or $1 \ge ||f - \langle f \rangle_{\pi}||_{2,\pi} > 0$, in which case

$$\left\|\mathbf{P}f - \langle f \rangle_{\pi}\right\|_{2,\pi} \le \left\|\mathbf{P}\left(\frac{f - \langle f \rangle_{\pi}}{\|f - \langle f \rangle_{\pi}\|_{2,\pi}}\right)\right\|_{2,\pi}$$

is also dominated by the right hand side of (6.2.1).

To go further, we need to borrow the following simple result from the theory of symmetric operators on a Hilbert space. Namely,

$$\sup \{ \|\mathbf{P}f\|_{2,\pi} : f \in L_0^2(\pi) \& \|f\|_{2,\pi} = 1 \}$$

=
$$\sup \{ |\langle f, \mathbf{P}f \rangle_{\pi} | : f \in L_0^2(\pi) \& \|f\|_{2,\pi} = 1 \}.$$
 (6.2.2)

That the right hand side of (6.2.2) is dominated by the left is Schwarz's inequality: $|\langle f, \mathbf{P}f \rangle_{\pi}| \leq ||f||_{2,\pi} ||\mathbf{P}f||_{2,\pi}$. To prove the opposite inequality, let $f \in L_0^2(\pi)$ with $||f||_{2,\pi} = 1$ be given, assume that $||\mathbf{P}f||_{2,\pi} > 0$, and set $g = \frac{\mathbf{P}f}{||\mathbf{P}f||_{2,\pi}}$. Then $g \in L_0^2(\pi)$ and $||g||_{2,\pi} = 1$. Hence, if γ denotes the supremum on the right hand side of (6.2.2), then, from the symmetry of \mathbf{P} ,

$$4\|\mathbf{P}f\|_{2,\pi} = 4\langle g, \mathbf{P}f \rangle_{\pi} = \langle (f+g), \mathbf{P}(f+g) \rangle_{\pi} - \langle (f-g), \mathbf{P}(f-g) \rangle_{\pi}$$
$$\leq \gamma \left(\|f+g\|_{2,\pi}^2 + \|f-g\|_{2,\pi}^2 \right) = 2\gamma \left(\|f\|_{2,\pi}^2 + \|g\|_{2,\pi}^2 \right) = 4\gamma.$$

The advantage won for us by (6.2.2) is that, in conjunction with (6.2.1), it shows that

$$eta=eta_+\wedgeeta_-$$

where
$$\beta_{\pm} \equiv \inf\{\langle f, (\mathbf{I} \mp \mathbf{P}) f \rangle_{\pi} : f \in L^{2}_{0}(\pi) \& ||f||_{2,\pi} = 1\}.$$
 (6.2.3)

Notice that, because $(\mathbf{I} - \mathbf{P})f = (\mathbf{I} - \mathbf{P})(f - \langle f \rangle_{\pi})$,

$$\beta_{+} = \inf\left\{\left\langle f, \left(\mathbf{I} - \mathbf{P}\right)f\right\rangle_{\pi} : f \in L^{2}(\pi) \& \operatorname{Var}_{\pi}(f) \equiv \left\| f - \langle f \rangle_{\pi} \right\|_{\pi}^{2} = 1\right\}.$$
 (6.2.4)

At the same time, because

$$\langle (f+c), (\mathbf{I}+\mathbf{P})(f+c) \rangle_{\pi} = \langle f, (\mathbf{I}+\mathbf{P})f \rangle_{\pi} + c^2 \text{ for } f \in L^2_0(\pi) \text{ and } c \in \mathbb{R},$$

it is clear that the infimum in the definition of β_{-} is the same whether we consider all $f \in L^{2}(\pi)$ or just those in $L_{0}^{2}(\pi)$. Hence, another expression for β_{-} is

$$\beta_{-} = \inf\{\langle f, (\mathbf{I} + \mathbf{P}) f \rangle_{\pi} : f \in L^{2}_{0}(\pi) \& ||f||_{2,\pi} = 1\}.$$
(6.2.5)

Two comments are in order here. First, when **P** is *non-negative definite*, abbreviated by $\mathbf{P} \ge 0$, in the sense that $\langle f, \mathbf{P}f \rangle_{\pi} \ge 0$ for all $f \in L^2(\pi)$, then $\beta_+ \le 1 \le \beta_-$, and so $\beta = \beta_+$. Hence,

$$\mathbf{P} \ge 0 \quad \Longrightarrow \quad \beta = \inf\{\{f, (\mathbf{I} - \mathbf{P})f\}_{\pi} : f \in L^2(\pi) \& \operatorname{Var}_{\pi}(f) = 1\}.$$
(6.2.6)

Second, by Theorem 6.1.14, we know that, in general, $\beta_{-} = 0$ unless the chain is aperiodic, and (cf. (6.1.11) and the discussion at the beginning of Sect. 6.1.4), when \mathbb{S} is finite, that $\beta > 0$ if and only if \mathbb{S} is aperiodic.

The expressions for β_+ and β_- in (6.2.4) and (6.2.5) lead to important calculational tools. Observe that, by (6.1.1),

$$\langle f, (\mathbf{I} - \mathbf{P}) f \rangle_{\pi} = \sum_{(i,j)} f(i)(\pi)_i (\mathbf{P})_{ij} (f(i) - f(j))$$

$$= \sum_{(i,j)} f(i)(\pi)_j (\mathbf{P})_{ji} (f(i) - f(j))$$

$$= \sum_{(i,j)} f(j)(\pi)_i (\mathbf{P})_{ij} (f(j) - f(i)).$$

Hence, when we add the second expression to the last, we find that

$$\langle f, (\mathbf{I} - \mathbf{P}) f \rangle_{\pi} = \mathcal{E}(f, f) \equiv \frac{1}{2} \sum_{i \neq j} (\pi)_i (\mathbf{P})_{ij} (f(j) - f(i))^2.$$
 (6.2.7)

Because the quadratic form $\mathcal{E}(f, f)$ is a discrete analog of the famous quadratic form $\frac{1}{2} \int |\nabla f|^2(x) dx$ introduced by Dirichlet, it is called a *Dirichlet form*. Extending this metaphor, one interprets β_+ as the *Poincaré constant*

$$\beta_{+} = \inf \{ \mathcal{E}(f, f) : f \in L^{2}(\pi) \& \operatorname{Var}_{\pi}(f) = 1 \}$$
(6.2.8)

in the Poincaré inequality

$$\beta_{+} \operatorname{Var}_{\pi}(f) \le \mathcal{E}(f, f), \quad f \in L^{2}(\pi).$$
(6.2.9)

To make an analogous application of (6.2.5), observe that

$$\sum_{(i,j)} (f(j) + f(i))^2 (\boldsymbol{\pi})_i (\mathbf{P})_{ij} = \langle \mathbf{P} f^2 \rangle_{\boldsymbol{\pi}} + 2 \langle f, \mathbf{P} f \rangle_{\boldsymbol{\pi}} + \|f\|_{2,\boldsymbol{\pi}}^2$$
$$= 2 \|f\|_{2,\boldsymbol{\pi}}^2 + 2 \langle f, \mathbf{P} f \rangle_{\boldsymbol{\pi}},$$

and therefore

$$\langle f, (\mathbf{I} + \mathbf{P}) f \rangle_{\pi} = \tilde{\mathcal{E}}(f, f) \equiv \frac{1}{2} \sum_{(i,j)} (\boldsymbol{\pi})_i (\mathbf{P})_{ij} (f(j) + f(i))^2.$$
 (6.2.10)

Hence

$$\beta_{-} = \inf \{ \tilde{\mathcal{E}}(f, f) : f \in L^{2}_{0}(\pi) \& \| f \|_{2,\pi} = 1 \}.$$
(6.2.11)

In order to give an immediate application of (6.2.9) and (6.2.11), suppose (cf. Exercise 2.4.3) that $(\mathbf{P})_{ij} \ge \epsilon_j$ for all (i, j), set $\epsilon = \sum_j \epsilon_j$, assume that $\epsilon > 0$, and define the probability vector $\boldsymbol{\mu}$ so that $(\boldsymbol{\mu})_i = \frac{\epsilon_i}{\epsilon}$. Then, by Schwarz's inequality for expectations with respect to $\boldsymbol{\mu}$ and the variational characterization of $\operatorname{Var}_{\boldsymbol{\pi}}(f)$ as the minimum value of $a \rightsquigarrow \langle (f-a)^2 \rangle_{\boldsymbol{\pi}}$,

$$2\mathcal{E}(f,f) = \sum_{(i,j)} (\boldsymbol{\pi})_i (\mathbf{P})_{ij} (f(j) - f(i))^2$$

$$\geq \epsilon \sum_{(i,j)} (f(j) - f(i))^2 (\boldsymbol{\pi})_i (\boldsymbol{\mu})_j \geq \epsilon \sum_i (\langle f \rangle_{\boldsymbol{\mu}} - f(i))^2 (\boldsymbol{\pi})_i \geq \epsilon \operatorname{Var}_{\boldsymbol{\pi}}(f),$$

and, similarly, $\tilde{\mathcal{E}}(f, f) \geq \frac{\epsilon}{2} \operatorname{Var}_{\pi}(f)$. Hence, by (6.2.9), (6.2.11), and (6.2.3), $\beta \geq \frac{\epsilon}{2}$. Of course, this result is a significantly weaker than the one we get by combining Exercise 2.4.3 with the reasoning in Sect. 6.1.5. Namely, by that exercise we know that $\|\delta_i \mathbf{P}^n - \pi\|_{v} \leq 2(1-\epsilon)^n$, and so the reasoning in Sect. 6.1.5 tells us that $\beta \geq \epsilon$, which is twice as good as the estimate we are getting here.

6.2.2 Estimating β_+

The origin of many applications of (6.2.9) and (6.2.11) to the estimation of β_+ and β_- is the simple observation that

$$\operatorname{Var}_{\pi}(f) = \frac{1}{2} \sum_{ij} (f(i) - f(j))^2 (\pi)_i (\pi)_j, \qquad (6.2.12)$$

which is easily checked by expanding $(f(i) - f(j))^2$ and seeing that the sum on the right equals $2\langle f^2 \rangle_{\pi} - 2\langle f \rangle_{\pi}^2$.

The importance of (6.2.12) is that it expresses $\operatorname{Var}_{\pi}(f)$ in terms of difference between the values of f at different points in \mathbb{S} , and clearly $\mathcal{E}(f, f)$ is also given in terms of such differences. However, the differences that appear in $\mathcal{E}(f, f)$ are only between the values of f at pairs of points (i, j) for which $(\mathbf{P})_{ij} > 0$, whereas the right hand of side of (6.2.12) entails sampling all pairs (i, j). Thus, in order to estimate $\operatorname{Var}_{\pi}(f)$ in terms of $\mathcal{E}(f, f)$, it is necessary to choose, for each (i, j) with $i \neq j$, a directed path $p(i, j) = (k_0, \dots, k_n) \in \mathbb{S}^{n+1}$ with $k_0 = i$ and $k_n = j$, which is *allowable* in the sense $(\mathbf{P})_{k_{m-1}k_m} > 0$ for each $1 \leq m \leq n$, and to write

$$(f(i) - f(j))^2 = \left(\sum_{e \in p(i,j)} \Delta_e f\right)^2,$$

where the summation in *e* is taken over the oriented segments (k_{m-1}, k_m) in the path p(i, j), and, for $e = (k, \ell)$, $\Delta_e f \equiv f(\ell) - f(k)$. At this point there are various ways in which one can proceed. For example, given any $\{\alpha(e) : e \in p(i, j)\} \subseteq (0, \infty)$, Schwarz's inequality (cf. Exercise 1.3.1) shows that the quantity on the right hand side of (6.2.12) is dominated by $\frac{1}{2}$ times

$$\left(\sum_{e \in p(i,j)} \frac{\alpha(e)}{\rho(e)}\right) \left(\sum_{e \in p(i,j)} \frac{\rho(e)}{\alpha(e)} (\Delta_e f)^2\right)$$
$$\leq \max_{e \in p(i,j)} \frac{1}{\alpha(e)} \left(\sum_{e \in p(i,j)} \frac{\alpha(e)}{\rho(e)}\right) \sum_{e \in p(i,j)} (\Delta_e f)^2 \rho(e).$$

where $\rho(e) \equiv (\pi)_k(\mathbf{P})_{k\ell}$ when $e = (k, \ell)$. Thus, for any selection of paths $\mathcal{P} = \{p(i, j) : (i, j) \in \mathbb{S}^2 \setminus D\}$ (*D* here denotes the diagonal $\{(i, j) \in \mathbb{S}^2 : i = j\}$) and coefficients $\mathcal{A} = \{\alpha(e, p) : e \in p \in \mathcal{P}\} \subseteq (0, \infty)$,

$$\begin{aligned} \operatorname{Var}_{\pi}(f) &\leq \frac{1}{2} \sum_{p \in \mathcal{P}} w_{A}(p) \sum_{e \in p} (\Delta_{e} f)^{2} \rho(e) \\ &= \frac{1}{2} \sum_{e} (\Delta_{e} f)^{2} \rho(e) \left(\sum_{p \ni e} w_{A}(p) \right) \leq W(\mathcal{P}, \mathcal{A}) \mathcal{E}(f, f), \end{aligned}$$

where

$$w_{\mathcal{A}}(p) \equiv (\boldsymbol{\pi})_{i}(\boldsymbol{\pi})_{j}\left(\max_{e \in p} \frac{1}{\alpha(e, p)}\right) \sum_{e \in p} \frac{\alpha(e, p)}{\rho(e)} \quad \text{if } p = p(i, j),$$

and

$$W(\mathcal{P}, \mathcal{A}) = \sup_{e} \sum_{p \ni e} w_{\mathcal{A}}(p).$$

Hence, we have now shown that

$$\beta_+ \ge \frac{1}{W(\mathcal{P}, \mathcal{A})} \tag{6.2.13}$$

for every choice of allowable paths \mathcal{P} and coefficients \mathcal{A} .

The most effective applications of (6.2.13) depend on making a choice of \mathcal{P} and \mathcal{A} takes advantage of the particular situation under consideration. Given a selection

 \mathcal{P} of paths, one of the most frequently made choices of \mathcal{A} is to take $\alpha(e, p) \equiv 1$. In this case, (6.2.13) gives

$$\beta_{+} \geq \frac{1}{W(\mathcal{P})} \quad \text{where } W(\mathcal{P}) \equiv \sup_{e} \sum_{p \geq e} \sum_{e' \in p} \frac{(\pi(p))_{-}(\pi(p))_{+}}{\rho(e')}, \tag{6.2.14}$$

where $(\boldsymbol{\pi}(p))_{-} = (\boldsymbol{\pi})_i$ and $(\boldsymbol{\pi}(p))_{+} = (\boldsymbol{\pi})_j$ when *p* begins at *i* and ends at *j*.

Finally, it should be recognized that, in general, (6.2.13) gives no information. Indeed, although irreducibility guarantees that there always is at least one path connecting every pair of points, when S is infinite there is no guarantee that \mathcal{P} and \mathcal{A} can be chosen so that $W(\mathcal{P}, \mathcal{A}) < \infty$. Moreover, even when S is finite, and therefore $W(\mathcal{P}, \mathcal{A}) < \infty$ for every choice, only a judicious choice will make (6.2.13) yield a good estimate.

6.2.3 Estimating β_{-}

The estimation of β_{-} starting from (6.2.11) is a bit more contrived than the one of β_{+} starting from (6.2.9). For one thing, we already know (cf. Theorem 6.1.14) that $\beta_{-} = 0$ unless the chain is aperiodic. Thus, we will now require that the chain be aperiodic. As a consequence of aperiodicity and irreducibility, we know (cf. (3.1.14)) that, for each $i \in \mathbb{S}$, there always is a path $p(i) = (k_0, \dots, k_{2n+1})$ which is allowable (i.e. ($\mathbf{P})_{k_{m-1}k_m} > 0$ for each $1 \le m \le 2n + 1$), is closed (i.e., $k_0 = k_{2n+1}$), and starts at i (i.e., $k_0 = i$). Note our insistence that this path have an odd number of steps. The reason for our doing so is that when the number of steps is odd, an elementary exercise in telescoping sums shows that

$$2f(i) = \sum_{m=0}^{2n} (-1)^m (f(k_m) + f(k_{m+1})).$$

Thus, if $\tilde{\mathcal{P}}$ be a selection of such paths, one path p(i) for each $i \in \mathbb{S}$, and we make an associated choice of coefficients $\mathcal{A} = \{\alpha(e, p) : e \in p \in \tilde{\mathcal{P}}\} \subseteq (0, \infty)$, then, just as in the preceding section,

$$\begin{aligned} 4\|f\|_{2,\pi}^2 &= \sum_i \left(\sum_{e \in p(i)} (-1)^{m(e)} \tilde{\Delta}_e f\right)^2 (\pi)_i \\ &\leq \sum_{p \in \tilde{\mathcal{P}}} \pi(p) \left(\sum_{e \in p} \frac{\alpha(e,p)}{\rho(e)}\right) \left(\sum_{e \in p} (\tilde{\Delta}_e f)^2 \frac{\rho(e)}{\alpha(e,p)}\right), \end{aligned}$$

where, when $p = (k_0, ..., k_{2n+1})$, $\pi(p) = (\pi)_{k_0}$, and, for $0 \le m \le 2n$, m(e) = mand $\tilde{\Delta}_e f \equiv f(k_m) + f(k_{m+1})$ if $e = (k_m, k_{m+1})$. Hence, if

$$\tilde{w}(p) = \pi(p) \left(\max_{e \in p} \frac{1}{\alpha(e, p)} \right) \sum_{e \in p} \frac{\alpha(e, p)}{\rho(e)}$$

then

$$2\|f\|_{2,\pi}^2 \le \tilde{W}(\tilde{\mathcal{P}},\mathcal{A})\tilde{\mathcal{E}}(f,f) \quad \text{where } \tilde{W}(\tilde{\mathcal{P}},\mathcal{A}) \equiv \sup_e \sum_{p \ni e} \tilde{w}(p)$$

and, by (6.2.11), this proves that

$$\beta_{-} \geq rac{2}{ ilde{W}(ilde{\mathcal{P}}, \mathcal{A})}$$

for any choice of allowable paths $\tilde{\mathcal{P}}$ and coefficients \mathcal{A} satisfying the stated requirements. When we take $\alpha(e, p) \equiv 1$, this specializes to

$$\beta_{-} \ge \frac{2}{\tilde{W}(\tilde{\mathcal{P}})} \quad \text{where } \tilde{W}_{1}(\tilde{\mathcal{P}}) \equiv \sup_{e} \sum_{p \ni e} \sum_{e' \in p} \frac{\pi(p)}{\rho(e')}, \tag{6.2.15}$$

It should be emphasized that the preceding method for getting estimates on β_{-} is inherently flawed in that it appears incapable of recognizing spectral properties of **P** like non-negative definiteness. In particular, when **P** is non-negative definite, then $\beta_{-} \ge 1$, but it seems unlikely that one could get that conclusion out of the arguments being used here. Thus, (6.2.15) should used only as a last resort, when nothing else seems to work.

6.3 Reversible Markov Processes in Continuous Time

Here we will see what the preceding theory looks like in the continuous time context and will learn that it is both easier and more æsthetically pleasing there.

We will be working with the notation and theory developed in Chap. 5. In particular, \mathbf{Q} will be a *Q*-matrix, and we will assume that \mathbb{S} is irreducible with respect to \mathbf{Q} .

6.3.1 Criterion for Reversibility

Let **Q** be given, assume that the associated Markov process never explodes (cf. Sect. 5.3.1), and use $t \rightsquigarrow \mathbf{P}(t)$ to denote the transition probability function determined

by **Q** (cf. Theorem 5.3.6). Our purpose in this subsection is to show that if $\hat{\mu}$ is a probability vector for which the *detailed balance condition*

$$(\hat{\boldsymbol{\mu}})_i(\mathbf{Q})_{ij} = (\hat{\boldsymbol{\mu}})_j(\mathbf{Q})_{ji} \quad \text{for all } (i,j) \in \mathbb{S}^2, \tag{6.3.1}$$

holds relative to \mathbf{Q} , then the detailed balance condition

$$(\hat{\boldsymbol{\mu}})_i \left(\mathbf{P}(t) \right)_{ij} = (\hat{\boldsymbol{\mu}})_j \left(\mathbf{P}(t) \right)_{ji} \quad \text{for all } t > 0 \text{ and } (i, j) \in \mathbb{S}^2$$
(6.3.2)

also holds.

The proof that (6.3.1) implies (6.3.2) is trivial in the case when the canonical rates for **Q** are bounded. Indeed, all that we need to do in that case is first use a simple inductive argument to check that (6.3.1) implies $(\hat{\mu})_i (\mathbf{Q}^n)_{ij} = (\hat{\mu})_j (\mathbf{Q}^n)_{ji}$ for all $n \ge 0$ and then use the expression for $\mathbf{P}(t)$ given in (5.1.9). When the rates are unbounded, we will use the approximation procedure introduced in Sect. 5.3.1. Namely, refer to Sect. 5.3.1, and take $\mathbf{Q}^{(N)}$ corresponding to the choice rates $\Re^{(N)}$ described there. Equivalently, take $(\mathbf{Q}^{(N)})_{ij}$ to be $(\mathbf{Q})_{ij}$ if $i \in F_N$ and 0 when $i \notin F_N$. Using induction, one finds first that $((\mathbf{Q}^{(N)})^n)_{ij} = 0$ for all $n \ge 1$ and $i \notin F_N$ and second that

$$(\hat{\boldsymbol{\mu}})_i ((\mathbf{Q}^{(N)})^n)_{ij} = (\hat{\boldsymbol{\mu}})_j ((\mathbf{Q}^{(N)})^n)_{ji} \text{ for all } n \ge 0 \text{ and } (i, j) \in F_N^2.$$

Hence, if { $\mathbf{P}^{(N)}(t) : t > 0$ } is the semigroup determined by $\mathbf{Q}^{(N)}$, then, since the rates for $\mathbf{Q}^{(N)}$ are bounded, (5.1.9) shows that

$$\left(\mathbf{P}^{(N)}(t) \right)_{ij} = \delta_{i,j} \quad \text{if } i \notin F_N (\hat{\boldsymbol{\mu}})_i \left(\mathbf{P}^{(N)}(t) \right)_{ij} = (\hat{\boldsymbol{\mu}})_j \left(\mathbf{P}^{(N)}(t) \right)_{ji} \quad \text{if } (i,j) \in F_N^2.$$

$$(6.3.3)$$

Therefore, because, by (5.3.4), $(\mathbf{P}^{(N)}(t))_{ij} \longrightarrow (\mathbf{P}(t))_{ij}$, we are done.

As a consequence of the preceding, we now know that (6.3.1) implies that $\hat{\mu}$ is stationary for **P**(*t*). Hence, because we are assuming that **Q** is irreducible, the results in Sect. 5.4.2 and Sect. 5.4.3 allow us to identify $\hat{\mu}$ as the probability vector $\hat{\pi} \equiv \hat{\pi}^{\mathbb{S}}$ introduced in Theorem 5.4.6 and discussed in Sect. 5.4.3. To summarize, *if* $\hat{\mu}$ *is a probability vector for which* (6.3.1) *holds, then* $\hat{\mu} = \hat{\pi}$.

6.3.2 Convergence in $L^2(\hat{\pi})$ for Bounded Rates

In view of the results just obtained, from now on we will be assuming that $\hat{\pi}$ is a probability vector for which (6.3.1) holds when $\hat{\mu} = \hat{\pi}$. In particular, this means that

$$(\hat{\boldsymbol{\pi}})_i \left(\mathbf{P}(t) \right)_{ij} = (\hat{\boldsymbol{\pi}})_j \left(\mathbf{P}(t) \right)_{ji} \quad \text{for all } t > 0 \text{ and } (i, j) \in \mathbb{S}^2.$$
(6.3.4)

Knowing (6.3.4), one is tempted to use the ideas in Sect. 6.2 to get an estimate on the rate, as measured by convergence in $L^2(\hat{\pi})$, at which $\mathbf{P}(t)f$ tends to $\langle f \rangle_{\hat{\pi}}$. To be precise, first note that, for each h > 0,

$$\langle f, \mathbf{P}(h) f \rangle_{\hat{\pi}} = \left\| \mathbf{P}\left(\frac{h}{2}\right) f \right\|_{2,\hat{\pi}}^2 \ge 0,$$

and therefore (cf. (6.1.13), (6.2.6), and (6.2.7)) that

$$\begin{split} \beta(h) &\equiv 1 - \sup\{\left| \left\langle f, \mathbf{P}(h) f \right\rangle_{\hat{\pi}} \right| : f \in L_0^2(\hat{\pi}) \text{ with } \|f\|_{2,\hat{\pi}} = 1 \} \\ &= \inf\{\left\langle f, \mathbf{I} - \mathbf{P}(h) f \right\rangle_{\hat{\pi}} : \operatorname{Var}_{\hat{\pi}}(f) = 1 \} = \inf\{\mathcal{E}_h(f, f) : \operatorname{Var}_{\hat{\pi}}(f) = 1 \}, \end{split}$$

where

$$\mathcal{E}_{h}(f,f) \equiv \frac{1}{2} \sum_{i \neq j} (\hat{\boldsymbol{\pi}})_{i} \big(\mathbf{P}(h) \big)_{ij} \big(f(j) - f(i) \big)^{2}$$

is the Dirichlet form for $\mathbf{P}(h)$ on $L^2(\hat{\pi})$. Hence (cf. (6.1.12)), for any t > 0 and $n \in \mathbb{Z}^+$,

$$\left\|\mathbf{P}(t)f - \langle f \rangle_{\hat{\pi}}\right\|_{L^{2}(\hat{\pi})} \leq \left(1 - \beta\left(\frac{t}{n}\right)\right)^{n} \left\|f - \langle f \rangle_{\hat{\pi}}\right\|_{2,\hat{\pi}}.$$

To take the next step, we add the assumption that the canonical rates are bounded. Then, because, by (5.2.14),

$$\mathbf{P}(t) = \mathbf{I} + \int_0^t \mathbf{Q} \mathbf{P}(\tau) d\tau = \mathbf{I} + t\mathbf{Q} + \int_0^t (t - \tau) \mathbf{Q}^2 \mathbf{P}(\tau) d\tau,$$
$$\|\mathbf{P}(t) - \mathbf{I} - t\mathbf{Q}\|_{u,v} \le \frac{\|\mathbf{Q}\|_{u,v}^2 t^2}{2}.$$

From this it follows that, uniformly in *f* satisfying $\operatorname{Var}_{\hat{\pi}}(f) = 1$,

$$\lim_{h \searrow 0} \frac{\mathcal{E}_h(f, f)}{h} = \mathcal{E}^{\mathbf{Q}}(f, f)$$

where

$$\mathcal{E}^{\mathbf{Q}}(f,f) \equiv \frac{1}{2} \sum_{i \neq j} (\hat{\pi})_i (\mathbf{Q})_{ij} (f(j) - f(i))^2;$$
(6.3.5)

and from this it follows that the limit $\lim_{h \to 0} h^{-1}\beta(h)$ exists and is equal to

$$\lambda \equiv \inf \left\{ \mathcal{E}^{\mathbf{Q}}(f, f) : f \in L^2(\boldsymbol{\pi}) \& \operatorname{Var}_{\hat{\boldsymbol{\pi}}}(f) = 1 \right\}.$$
(6.3.6)

Thus, at least when the rates are bounded, we know that

$$\left\|\mathbf{P}(t)f - \langle f \rangle_{\hat{\pi}}\right\|_{2,\hat{\pi}} \le e^{-\lambda t} \left\|f - \langle f \rangle_{\hat{\pi}}\right\|_{2,\hat{\pi}}.$$
(6.3.7)

6.3.3 $L^2(\hat{\pi})$ -Convergence Rate in General

When the rates are unbounded, the preceding line of reasoning is too naïve. In order to treat the unbounded case, one needs to make some additional observations, all of which have their origins in the following lemma.³

Lemma 6.3.8 Given $f \in L^2(\hat{\pi})$, the function $t \in [0, \infty) \mapsto \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2$ is continuous, non-increasing, non-negative, and convex. In particular, $t \in (0, \infty) \mapsto \frac{\langle f, (\mathbf{I}-\mathbf{P}(t))f \rangle_{\hat{\pi}}}{t}$ is non-increasing, and therefore

$$\lim_{h \searrow 0} \frac{\|f\|_{2,\hat{\pi}}^2 - \|\mathbf{P}(h)f\|_{2,\hat{\pi}}^2}{h} \quad exists \text{ in } [0,\infty].$$

Moreover (cf. (6.3.5)),

$$\lim_{h \searrow 0} \frac{\|f\|_{2,\hat{\pi}}^2 - \|\mathbf{P}(h)f\|_{2,\hat{\pi}}^2}{h} \ge 2\mathcal{E}^{\mathbf{Q}}(f,f).$$

Proof Let $f \in L^2(\hat{\pi})$ be given, and (cf. the notation in Sect. 5.3.1) define $f_N = \mathbf{1}_{F_N} f$. Then, by Lebesgue's dominated convergence theorem,

$$||f - f_N||_{2,\hat{\pi}} \longrightarrow 0 \text{ as } N \to \infty.$$

Because $\|\mathbf{P}(t)g\|_{2,\hat{\pi}} \le \|g\|_{2,\hat{\pi}}$ for all t > 0 and $g \in L^2(\hat{\pi})$, we know that

$$\left| \left\| \mathbf{P}(t) f \right\|_{2,\hat{\pi}} - \left\| \mathbf{P}(t) f_N \right\|_{2,\hat{\pi}} \right| \le \left\| \mathbf{P}(t) (f - f_N) \right\|_{2,\hat{\pi}} \le \| f - f_N \|_{2,\hat{\pi}} \longrightarrow 0$$

uniformly in $t \in (0, \infty)$ as $N \to \infty$. Hence, by part (a) of Exercise 6.6.1 below, in order to prove the initial statement, it suffices to do so when f vanishes off of F_M for some M. Now let f be a function which vanishes off of F_M , and set $\psi(t) = \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2$. At the same time, set $\psi_N(t) = \|\mathbf{P}^{(N)}(t)f\|_{2,\hat{\pi}}^2$ for $N \ge M$. Then, because by (5.3.4), $\psi_N \longrightarrow \psi$ uniformly on finite intervals, another application of part (a) in Exercise 6.6.1 allows us to restrict our attention to the ψ_N 's. That is, we will have proved that ψ is a continuous, non-increasing, non-negative, convex function as soon as we show that each ψ_N is. The non-negativity requires no comment. To prove the other properties, we apply (5.2.14) to see that

$$\dot{\psi}_N(t) = \left\langle \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f, \mathbf{P}^{(N)}(t) f \right\rangle_{\hat{\pi}} + \left\langle \mathbf{P}^{(N)}(t) f, \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f \right\rangle_{\hat{\pi}}.$$

Next, by the first line of (6.3.3), we know that, because $N \ge M$, $\mathbf{P}^{(N)}(t) f$ vanishes off of F_N , and so, because $(\hat{\boldsymbol{\pi}})_i \mathbf{Q}_{ij}^{(N)} = (\hat{\boldsymbol{\pi}})_j \mathbf{Q}_{ji}^{(N)}$ for $(i, j) \in F_N^2$, the preceding becomes

$$\dot{\psi}_N(t) = 2 \langle \mathbf{P}^{(N)}(t) f, \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f \rangle_{\hat{\pi}}.$$

³If one knows spectral theory, especially Stone's theorem, the rather cumbersome argument that follows can be avoided.

Similarly, we see that

$$\begin{split} \ddot{\psi}_{N}(t) &= 2 \langle \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f, \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f \rangle_{\hat{\pi}} \\ &+ 2 \langle \mathbf{P}^{(N)}(t) f, \left(\mathbf{Q}^{(N)} \right)^{2} \mathbf{P}^{(N)}(t) f \rangle_{\hat{\pi}} \\ &= 4 \langle \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f, \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f \rangle_{\hat{\pi}} = 4 \| \mathbf{Q}^{(N)} \mathbf{P}^{(N)}(t) f \|_{2,\hat{\pi}}^{2} \ge 0. \end{split}$$

Clearly, the second of these proves the convexity of ψ_N . In order to see that the first implies that ψ_N is non-increasing, we will show that

$$\langle g, -\mathbf{Q}^{(N)}g \rangle_{\hat{\pi}} = \frac{1}{2} \sum_{\substack{(i,j) \in F_N^2 \\ i \neq j}} (\hat{\pi})_i (\mathbf{Q})_{ij} (g(j) - g(i))^2 + \sum_{i \in F_N} (\hat{\pi})_i V^{(N)}(i)g(i)^2$$

if
$$g = 0$$
 off F_N and $V^{(N)}(i) \equiv \sum_{j \notin F_N} Q_{ij}$ for $i \in F_N$. (*)

To check (*), first observe that

$$\begin{split} \left\langle g, -\mathbf{Q}^{(N)}g \right\rangle_{\hat{\pi}} &= -\sum_{\substack{(i,j) \in F_N^2 \\ i \neq j}} (\hat{\pi})_i(\mathbf{Q})_{ij}g(i)g(j) \\ &= -\sum_{\substack{(i,j) \in F_N^2 \\ i \neq j}} (\hat{\pi})_i(\mathbf{Q})_{ij}g(i) \big(g(j) - g(i)\big) + \sum_{i \in F_N} (\hat{\pi})_i V^{(N)}(i)g(i)^2. \end{split}$$

Next, use $(\hat{\boldsymbol{\pi}})_i(\mathbf{Q})_{ij} = (\hat{\boldsymbol{\pi}})_j(\mathbf{Q})_{ji}$ for $(i, j) \in F_N^2$ to see that

$$-\sum_{\substack{(i,j)\in F_N\\i\neq j}} (\hat{\pi})_i(\mathbf{Q})_{ij}g(i) (g(j) - g(i)) = \sum_{\substack{(i,j)\in F_N\\i\neq j}} (\hat{\pi})_i(\mathbf{Q})_{ij}g(j) (g(j) - g(i)),$$

and thereby arrive at (*). Finally, apply (*) with $g = \mathbf{P}^{(N)}(t) f$ to conclude that $\dot{\psi}_N \leq 0$.

Turning to the second and third assertions, let f be any element of $L^2(\hat{\pi})$. Now that we know that the corresponding ψ is a continuous, non-increasing, non-negative, convex function, it is easy (cf. part (d) in Exercise 6.6.1) to check that $t \rightsquigarrow \frac{\psi(0) - \psi(t)}{t}$ is non-increasing and therefore that $\lim_{h \searrow 0} \frac{\psi(0) - \psi(h)}{h}$ exists in $[0, \infty]$. Next, remember (6.2.7), apply it when $\mathbf{P} = \mathbf{P}(2h)$, and conclude that

$$\psi(0) - \psi(h) = \left\langle f, \left(\mathbf{I} - \mathbf{P}(2h)\right) f \right\rangle_{\hat{\pi}} = \frac{1}{2} \sum_{\substack{(i,j)\\i \neq j}} (\hat{\pi})_i \left(\mathbf{P}(2h)\right)_{ij} \left(f(j) - f(i)\right)^2.$$

Hence, since, by (5.3.5),

$$\lim_{h \searrow 0} \frac{(\mathbf{P}(2h))_{ij}}{h} = 2(\mathbf{Q})_{ij} \quad \text{for } i \neq j,$$

the required inequality follows after an application of Fatou's lemma.

Lemma 6.3.9 If 0 < s < t, then for any $f \in L^2(\hat{\pi})$

$$\frac{\|f\|_{2,\hat{\pi}}^2}{s} \ge \frac{\|\mathbf{P}(s)f\|_{2,\hat{\pi}}^2 - \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2}{t-s} \ge 2\mathcal{E}^{\mathbf{Q}}\big(\mathbf{P}(t)f, \mathbf{P}(t)f\big).$$

Proof Set $\psi(t) = \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2$. We know that ψ is a continuous, non-increasing, non-negative, convex function. Hence, by part (a) of Exercise 5.5.2,

$$\frac{\psi(0)}{s} \ge \frac{\psi(0) - \psi(s)}{s} \ge \frac{\psi(s) - \psi(t)}{t - s} \ge \frac{\psi(t) - \psi(t + h)}{h}$$

for any h > 0. Moreover, because,

$$\frac{\psi(t) - \psi(t+h)}{h} = \frac{\|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2 - \|\mathbf{P}(h)\mathbf{P}(t)f\|_{2,\hat{\pi}}^2}{h}$$

the last part of Lemma 6.3.8 applied with $\mathbf{P}(t)f$ replacing f yields the second asserted estimate.

With the preceding at hand, we can now complete our program. Namely, by writing

$$\|f\|_{2,\hat{\pi}}^2 - \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2 = \sum_{m=0}^{n-1} \left(\left\|\mathbf{P}\left(\frac{mt}{n}\right)\right\|_{2,\hat{\pi}}^2 - \left\|\mathbf{P}\left(\frac{(m+1)t}{n}\right)\right\|_{2,\hat{\pi}}^2 \right),$$

we can use the result in Lemma 6.3.9 to obtain the estimate

$$\|f\|_{2,\hat{\pi}}^2 - \left\|\mathbf{P}(t)f\right\|_{2,\hat{\pi}}^2 \ge \frac{2t}{n} \sum_{m=1}^n \mathcal{E}^{\mathbf{Q}}\left(\mathbf{P}\left(\frac{mt}{n}\right)f, \mathbf{P}\left(\frac{mt}{n}\right)f\right).$$

Hence, if λ is defined as in (6.3.6), then, for any $f \in L_0^2(\hat{\pi})$,

$$\|f\|_{2,\hat{\pi}}^2 - \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2 \ge \frac{2\lambda t}{n} \sum_{m=1}^n \left\|\mathbf{P}\left(\frac{mt}{n}\right)f\right\|_{2,\hat{\pi}}^2,$$

which, when $n \to \infty$, leads to

$$\|f\|_{2,\hat{\pi}}^{2} - \|\mathbf{P}(t)f\|_{2,\hat{\pi}}^{2} \ge 2\lambda \int_{0}^{t} \|\mathbf{P}(\tau)f\|_{2,\hat{\pi}}^{2} d\tau.$$

Finally, by Gronwall's inequality (cf. Exercise 4.2.2), the preceding yields the estimate $\|\mathbf{P}(t)f\|_{2,\hat{\pi}}^2 \leq e^{-2\lambda t} \|f\|_{2,\hat{\pi}}^2$. After replacing a general $f \in L^2(\hat{\pi})$ by $f - \langle f \rangle_{\hat{\pi}}$, we have now proved that (6.3.7) holds even when the canonical rates are unbounded.

6.3.4 Estimating λ

Proceeding in the exactly the same way that we did in Sect. 6.2.2, we can estimate the λ in (6.3.6) in the same way as we estimated β_+ there. Namely, we make a selection \mathcal{P} consisting of paths p(i, j), one for each from pair $(i, j) \in \mathbb{S} \setminus D$, with the properties that if $p(i, j) = (k_0, \dots, k_n)$, then $k_0 = i, k_n = j$, and p(i, j) is *allowable* in the sense that $(\mathbf{Q})_{k_{m-1}k_m} > 0$ for each $1 \le m \le n$. Then, just as in Sect. 6.2.2, we can say that

$$\lambda \ge \frac{1}{W(\mathcal{P})} \quad \text{where } W(\mathcal{P}) \equiv \sup_{e} \sum_{p \ni e} \sum_{e' \in p} \frac{(\hat{\pi}(p))_{-}(\hat{\pi}(p))_{+}}{\rho(e')}, \tag{6.3.10}$$

where the supremum is over oriented edges $e = (k, \ell)$ with $(\mathbf{Q})_{k\ell} > 0$, the first sum is over $p \in \mathcal{P}$ in which the edge *e* appears, the second sum is over edges *e'* which appear in the path *p*, $(\hat{\pi}(p))_{-} = (\hat{\pi})_i$ if the path *p* starts at *i*, $(\hat{\pi}(p))_{+} = (\hat{\pi})_j$ if the path *p* ends at *j*, and $\rho(e') = (\hat{\pi})_k(\mathbf{Q})_{k\ell}$ if $e' = (k, \ell)$. See [1] for further applications of these results.

6.4 Gibbs States and Glauber Dynamics

Loosely speaking, the physical principle underlying statistical mechanics can be summarized in the statement that, when a system is in equilibrium, *states with lower energy are more likely than those with higher energy*. In fact, J.W. Gibbs sharpened this statement by saying that the probability of a state *i* will be proportional to $e^{-\frac{H(i)}{kT}}$, where *k* is the Boltzmann constant, *T* is temperature, and H(i) is the energy of the system when it is in state *i*. For this reason, a distribution which assigns probabilities in this Gibbsian manner is called a *Gibbs state*.

Since a Gibbs state is to be a model of equilibrium, it is only reasonable to ask what is the dynamics for which it is the equilibrium. From our point of view, this means that we should seek a Markov process for which the Gibbs state is the stationary distribution. Further, because dynamics in physics should be reversible, we should be looking for Markov processes which are reversible with respect to the Gibbs state, and, because such processes were introduced in this context by R. Glauber, we will call a Markov process which is reversible with respect to a Gibbs state a *Glauber dynamics* for that Gibbs state.

In this section, we will give a rather simplistic treatment of Gibbs states and their associated Glauber dynamics.

6.4.1 Formulation

Throughout this section, we will be working in the following setting. As usual, \mathbb{S} is either a finite or countably infinite space. On \mathbb{S} there is given some "natural" background assignment $\mathbf{v} \in (0, \infty)^{\mathbb{S}}$ of (not necessarily summable) weights, which should be thought of as a row vector. In many applications, \mathbf{v} is uniform: it assigns each *i* weight 1, but in other situations it is convenient to not have to assume that it is uniform. Next, there is a function $H : \mathbb{S} \longrightarrow [0, \infty)$ (alias, the energy function) with the property that

$$Z(\beta) \equiv \sum_{i \in \mathbb{S}} e^{-\beta H(i)}(\mathbf{v})_i < \infty \quad \text{for each } \beta \in (0, \infty).$$
(6.4.1)

In the physics metaphor, $\beta = \frac{1}{kT}$ is, apart from Boltzmann's constant *k*, the reciprocal temperature, and physicists would call $\beta \rightsquigarrow Z(\beta)$ the *partition function*. Finally, for each $\beta \in (0, \infty)$, the *Gibbs state* $\gamma(\beta)$ is the probability vector given by

$$\left(\boldsymbol{\gamma}(\beta)\right)_{i} = \frac{1}{Z(\beta)} e^{-\beta H(i)}(\boldsymbol{v})_{i} \quad \text{for } i \in \mathbb{S}.$$
 (6.4.2)

From a physical standpoint, everything of interest is encoded in the partition function. For example, it is elementary to compute both the average and variance of the energy by taking logarithmic derivatives:

$$\langle H \rangle_{\boldsymbol{\gamma}(\beta)} = -\frac{d}{d\beta} \log Z(\beta) \quad \text{and} \quad \operatorname{Var}_{\boldsymbol{\gamma}(\beta)}(H) = \frac{d^2}{d\beta^2} \log Z(\beta).$$
 (6.4.3)

The final ingredient is the description of the Glauber dynamics. For this purpose, we start with a matrix \mathbf{A} all of whose entries are non-negative and whose diagonal entries are 0. Further, we assume that \mathbf{A} is irreducible in the sense that

$$\sup_{n \ge 0} (\mathbf{A}^n)_{ij} > 0 \quad \text{for all } (i, j) \in \mathbb{S}^2$$
(6.4.4)

and that it is reversible in the sense that

$$(\mathbf{v})_i(\mathbf{A})_{ij} = (\mathbf{v})_j(\mathbf{A})_{ji}$$
 for all $(i, j) \in \mathbb{S}^2$. (6.4.5)

Finally, we insist that

$$\sum_{j \in \mathbb{S}} e^{-\beta H(j)}(\mathbf{A})_{ij} < \infty \quad \text{for each } i \in \mathbb{S} \text{ and } \beta > 0.$$
 (6.4.6)

At this point there are many ways in which to construct a Glauber dynamics. However, for our purposes, the one which will serve us best is the one whose Q-matrix is given by

$$\left(\mathbf{Q}(\beta) \right)_{ij} = e^{-\beta (H(j) - H(i))^+} (\mathbf{A})_{ij} \quad \text{when } j \neq i$$

$$\left(\mathbf{Q}(\beta) \right)_{ii} = -\sum_{j \neq i} \left(\mathbf{Q}(\beta) \right)_{ij},$$

$$(6.4.7)$$

where $a^+ \equiv a \lor 0$ is the non-negative part of the number $a \in \mathbb{R}$. Because,

$$\left(\boldsymbol{\gamma}(\beta)\right)_{i} \left(\mathbf{Q}(\beta)\right)_{ij} = Z(\beta)^{-1} e^{-\beta H(i) \vee H(j)} (\boldsymbol{\nu})_{i} (\mathbf{A})_{ij} \quad \text{for } i \neq j,$$
(6.4.8)

 $\gamma(\beta)$ clearly is reversible for $\mathbf{Q}(\beta)$. There are many other possibilities, and the optimal choice is often dictated by special features of the situation under consideration. However, whatever choice is made, it should be made in such a way that, for each $\beta > 0$, $\mathbf{Q}(\beta)$ determines a Markov process which never explodes.

6.4.2 The Dirichlet Form

In this subsection we will modify the ideas developed in Sect. 6.2.3 to get a lower bound on

$$\lambda_{\beta} \equiv \inf \left\{ \mathcal{E}_{\beta}(f, f) : \operatorname{Var}_{\beta}(f) = 1 \right\}$$

where $\mathcal{E}_{\beta}(f, f) \equiv \frac{1}{2} \sum_{j \neq i} (\boldsymbol{\gamma}(\beta))_{i} (\mathbf{Q}(\beta))_{ij} (f(j) - f(i))^{2}$ (6.4.9)

and $\operatorname{Var}_{\beta}(f)$ is shorthand for $\operatorname{Var}_{\gamma(\beta)}(f)$, the variance of f with respect to $\gamma(\beta)$. For this purpose, we introduce the notation

$$\operatorname{Elev}(p) = \max_{0 \le m \le n} H(i_m)$$
 and $e(p) = \operatorname{Elev}(p) - H(i_0) - H(i_n)$

for an allowable path $p = (i_0, ..., i_n)$. Assuming that $\mathbf{Q}(\beta)$ is given by (6.4.7), one sees that, when $p = (i_0, ..., i_n)$

$$w_{\beta}(p) \equiv \sum_{m=1}^{n} \frac{(\boldsymbol{\gamma}(\beta))_{i_0}(\boldsymbol{\gamma}(\beta))_{i_n}}{(\boldsymbol{\gamma}(\beta))_{i_{m-1}}(\mathbf{Q}(\beta))_{i_{m-1}i_m}} \le Z(\beta)^{-1} e^{\beta e(p)} w(p),$$

where $w(p) \equiv \sum_{m=1}^{n} \frac{(\boldsymbol{\nu})_{i_0}(\boldsymbol{\nu})_{i_n}}{(\boldsymbol{\nu})_{i_{m-1}} \mathbf{A}_{i_{m-1}i_m}}.$

Hence, for any choice of paths \mathcal{P} , we know that (cf. (6.3.10))

$$W_{\beta}(\mathcal{P}) \equiv \sup_{e} \sum_{p \ni e} w_{\beta}(p) \le Z(\beta)^{-1} e^{\beta E(\mathcal{P})} W(\mathcal{P}),$$

where $W(\mathcal{P}) \equiv \sup_{e} \sum_{p \ni e} w(p)$ and $E(\mathcal{P}) \equiv \sup_{p \in \mathcal{P}} e(p),$

and therefore that

$$\lambda_{\beta} \ge \frac{Z(\beta)e^{-\beta E(\mathcal{P})}}{W(\mathcal{P})}.$$
(6.4.10)

On the one hand, it is clear that (6.4.10) gives information only when $W(\mathcal{P}) < \infty$. At the same time, it shows that, at least if ones interest is in large β 's, then it is important to choose \mathcal{P} so that $E(\mathcal{P})$ is as small as possible. When S is finite, reconciling these two creates no problem. Indeed, the finiteness of S guarantees that $W(\mathcal{P})$ will be finite for every choice of allowable paths. In addition, finiteness allows one to find for each (i, j) a path p(i, j) which minimizes Elev(p) among allowable paths from *i* to *j*, and clearly any \mathcal{P} consisting of such paths will minimize $E(\mathcal{P})$. Of course, it is sensible to choose such a \mathcal{P} so as to minimize $W(\mathcal{P})$ as well. In any case, whenever S is finite and \mathcal{P} consists of paths p(i, j) which minimize Elev(p) among paths p between i and j, $E(\mathcal{P})$ has a nice interpretation. Namely, think of S as being sites on a map and of H as giving the altitude of the sites. That is, in this metaphor, H(i)is the distance of *i* "above sea level." Without loss in generality, we will assume that at least one site k_0 is at sea level: $H(k_0) = 0.4$ When such an k_0 exists, the metaphorical interpretation of $E(\mathcal{P})$ is as the least upper bound on the altitude a hiker must gain, no matter where he starts or what allowable path he chooses to follow, in order to reach sea level. To see this, first observe that if p and p' are a pair of allowable paths and if the end point of p is the initial point of p', then the path q is allowable and $\text{Elev}(q) = \text{Elev}(p) \lor \text{Elev}(p')$ when q is obtained by concatenating p and p': if $p = (i_0, \ldots, i_n)$ and $p' = (i'_0, \ldots, i'_n)$, then $q = (i_0, \ldots, i_n, i'_1, \ldots, i'_n)$. Hence, for any (i, j), $e(p(i, j)) \le e(p(i, k_0)) \lor e(p(j, k_0))$, from which it should be clear that $E(\mathcal{P}) = \max_i e(p(i, k_0))$. Finally, since, for each i, $e(p(i, k_0)) = H(\ell) - H(i)$, where ℓ is a highest point along the path $p(i, k_0)$, the explanation is complete. When \mathbb{S} is infinite, the same interpretation is valid in various circumstances. For example, it applies when H "tends to infinity at infinity" in the sense that $\{i : H(i) \le M\}$ is finite for each $M < \infty$.

When S is finite, we can show that, at least for large β , (6.4.10) is quite good. To be precise, we have the following result.

Theorem 6.4.11 Assume that S is finite and that $Q(\beta)$ is given by (6.4.7). Set $\mathfrak{m} = \min_{i \in S} H(i)$ and $S_0 = \{i : H(i) = \mathfrak{m}\}$, and let \mathfrak{e} be the minimum value $E(\mathcal{P})$

⁴If that is not already so, we can make it so by choosing k_0 to be a point at which *H* takes its minimum value and replacing *H* by $H - H(k_0)$. Such a replacement leaves both $\gamma(\beta)$ and $\mathbf{Q}(\beta)$ as well as the quantity on the right hand side of (6.4.10) unchanged.

takes as \mathcal{P} runs over all selections of allowable paths. Then $\mathfrak{e} \geq -\mathfrak{m}$, and $\mathfrak{e} = -\mathfrak{m}$ if and only if for each $(i, j) \in \mathbb{S} \times \mathbb{S}_0$ there is an allowable path p from i to j with $\operatorname{Elev}(p) = H(i)$. (See also Exercise 6.6.13 below.) More generally, whatever the value of \mathfrak{e} , there exist constants $0 < c_- \leq c_+ < \infty$, which are independent of H, such that

$$c_{-}e^{-\beta(\mathfrak{e}+\mathfrak{m})} \leq \lambda_{\beta} \leq c_{+}e^{-\beta(\mathfrak{e}+\mathfrak{m})} \quad \text{for all } \beta \geq 0.$$

Proof Because neither $\gamma(\beta)$ nor $\mathbf{Q}(\beta)$ is changed if *H* is replaced by $H - \mathfrak{m}$ whereas \mathfrak{e} changes to $\mathfrak{e} + \mathfrak{m}$, we may and will assume that $\mathfrak{m} = 0$.

Choose a collection $\mathcal{P} = \{p(i, j) : (i, j) \in \mathbb{S}^2\}$ of allowable paths so that, for each (i, j), e(p(i, j)) minimizes e(p) over allowable paths from *i* to *j*. Next choose and fix a $k_0 \in \mathbb{S}_0$. By the reasoning given above,

$$\mathbf{e} = \max_{i \in \mathbb{S}} e(p(i, k_0)). \tag{(*)}$$

In particular, since $e(p(i, k_0)) = \text{Elev}(p(i, k_0)) - H(i) \ge 0$ for all $i \in \mathbb{S}$, this proves that $\mathfrak{e} \ge 0$ and makes it clear that $\mathfrak{e} = 0$ if and only if $H(i) = \text{Elev}(p(i, k_0))$ for all $i \in \mathbb{S}$.

Turning to the lower bound for λ_{β} , observe that, because $\mathfrak{m} = 0$, $Z(\beta) \ge (\mathfrak{v})_{k_0} > 0$ and therefore, by (6.4.10), that we can take $c_- = \frac{(\mathfrak{v})_{k_0}}{W(\mathcal{P})}$.

Finally, to prove the upper bound, choose $\ell_0 \in \mathbb{S} \setminus \{k_0\}$ so that $e(p_0) = \mathfrak{e}$ when $p_0 \equiv p(\ell_0, k_0)$, let Γ be the set of $i \in \mathbb{S}$ with the property that either $i = k_0$ or $\text{Elev}(p(i)) < \text{Elev}(p_0)$ for the path $p(i) \equiv p(i, k_0) \in \mathcal{P}$ from i to k_0 , and set $f = \mathbf{1}_{\Gamma}$. Then, because $k_0 \in \Gamma$ and $\ell_0 \notin \Gamma$, (cf. (6.2.12)),

$$\operatorname{Var}_{\beta}(f) = \left(\sum_{i \in \Gamma} (\boldsymbol{\gamma}(\beta))_{i}\right) \left(\sum_{j \notin \Gamma} (\boldsymbol{\gamma}(\beta))_{j}\right)$$
$$\geq \left(\boldsymbol{\gamma}(\beta)\right)_{k_{0}} (\boldsymbol{\gamma}(\beta))_{\ell_{0}} = \frac{(\boldsymbol{\nu})_{k_{0}}(\boldsymbol{\nu})_{\ell_{0}}}{Z(\beta)^{2}} e^{-\beta H(\ell_{0})}.$$

At the same time

$$\mathcal{E}_{\beta}(f,f) = \sum_{(i,j)\in\Gamma\times\Gamma\mathcal{C}} (\boldsymbol{\gamma}(\beta))_i (\mathbf{Q}(\beta))_{ij} = \frac{1}{Z(\beta)} \sum_{(i,j)\in\Gamma\times\Gamma\mathcal{C}} (\boldsymbol{\nu})_i (\mathbf{A})_{ij} e^{-\beta H(i)\vee H(j)}.$$

If $i \in \Gamma \setminus \{k_0\}$, $j \notin \Gamma$, and $(\mathbf{A})_{ij} > 0$, then $H(j) \ge \operatorname{Elev}(p_0)$. To see this, consider the path q obtained by going in one step from j to i and then following p(i) from i to k_0 . Clearly, q is an allowable path from j to k_0 , and therefore $\operatorname{Elev}(q) \ge \operatorname{Elev}(p(j)) \ge \operatorname{Elev}(p_0)$. But this means that

$$\operatorname{Elev}(p_0) \leq \operatorname{Elev}(p(j)) \leq \operatorname{Elev}(q) = \operatorname{Elev}(p(i)) \lor H(j),$$

which, together with $\text{Elev}(p(i)) < \text{Elev}(p_0)$, forces the conclusion that $H(j) \ge$ $\text{Elev}(p_0)$. Even easier is the observation that $H(j) \ge \text{Elev}(p_0)$ if $j \notin \Gamma$ and $(\mathbf{A})_{k_0 j} > 0$, since in that case the path (j, k_0) is allowable and

$$H(j) = \operatorname{Elev}((j, k_0)) \ge \operatorname{Elev}(p(j)) \ge \operatorname{Elev}(p_0).$$

Hence, after plugging this into the preceding expression for $\mathcal{E}_{\beta}(f, f)$, we get

$$\mathcal{E}_{\beta}(f,f) \leq \frac{e^{-\beta \operatorname{Elev}(p_0)}}{Z(\beta)} \sum_{(i,j) \in \Gamma \times \Gamma^{\complement}} (\mathbf{v})_i(\mathbf{A})_{ij},$$

which, because $\mathfrak{e} = \operatorname{Elev}(p_0) - H(\ell_0)$, means that

$$\lambda_{\beta} \leq \frac{\mathcal{E}_{\beta}(f,f)}{\operatorname{Var}_{\beta}(f)} \leq \frac{Z(\beta)}{(\boldsymbol{\nu})_{k_{0}}(\boldsymbol{\nu})_{\boldsymbol{\ell}_{0}}} \left(\sum_{(i,j)\in\Gamma\times\Gamma\mathcal{C}}(\boldsymbol{\nu})_{i}(\mathbf{A})_{ij}\right) e^{-\beta\boldsymbol{\mathfrak{e}}}.$$

Finally, since $Z(\beta) \le \|\mathbf{v}\|_{v}$, the upper bound follows.

6.5 Simulated Annealing

This concluding section deals with an application of the ideas in the preceding section. Namely, given a function $H : \mathbb{S} \longrightarrow [0, \infty)$, we want to describe a procedure, variously known as the *simulated annealing* or the *Metropolis algorithm*, for locating a place where H achieves its minimum value.

In order to understand the intuition that underlies this procedure, let A be a matrix of the sort discussed in Sect. 6.4, assume that 0 is the minimum value of H, set $\mathbb{S}_0 =$ $\{i: H(i) = 0\}$, and think about dynamic procedures which would lead you from any initial point to \mathbb{S}_0 via paths which are allowable according to **A** (i.e., $\mathbf{A}_{k\ell} > 0$ if k and ℓ are successive points along the path). One procedure is based on the steepest decent strategy. That is, if one is at k, one moves to any one of the points ℓ for which $A_{k\ell} > 0$ and $H(\ell)$ is minimal if $H(\ell) \le H(k)$ for at least one such point, and one stays put if $H(\ell) > H(k)$ for every ℓ with $(\mathbf{A})_{k\ell} > 0$. This procedure works beautifully as long as you avoid, in the topographic metaphor suggested earlier, getting trapped in some "mountain valley." The point is that the steepest decent procedure is the most efficient strategy for getting to some local minimum of H. However, if that minimum is not global, then, in general, you will get stuck! Thus, if you are going to avoid this fate, occasionally you will have to go "up hill" even when you may have the option to go "down hill." However, unless you have a detailed a priori knowledge of the whole terrain, there is no way to know when you should decide to do so. For this reason, it may be best to abandon rationality and let the decision be made randomly. Of course, after a while, you should hope that you will have worked your way out of the mountain valleys and that a steepest decent strategy should become increasingly reasonable.

6.5.1 The Algorithm

In order to eliminate as many technicalities as possible, we will assume throughout that S is finite and has at least 2 elements. Next, let $H : \mathbb{S} \longrightarrow [0, \infty)$ be the function for which we want to locate a place where it achieves its minimum, and, without loss in generality, we will assume that 0 is its minimum. Now take \mathbf{v} so that $(\mathbf{v})_i = 1$ for all $i \in \mathbb{S}$, and choose a matrix **A** so that $(\mathbf{A})_{ii} = 0$ for all $i \in \mathbb{S}$, $(\mathbf{A})_{ij} = (\mathbf{A})_{ji} \ge 0$ if $j \ne i$, and **A** is irreducible (cf. (6.4.4)) on S. In practice, the selection of **A** should be made so that the evaluation of H(j) - H(i) when $(\mathbf{A})_{ij} > 0$ is as "easy as possible." For example, if S has some sort of natural neighborhood structure with respect to which S is connected and the computation of H(j) - H(i) when j is a neighbor of i requires very little time, then it is reasonable to take **A** so that $(\mathbf{A})_{ij} = 0$ unless $j \ne i$ is a neighbor of i.

Now define $\gamma(\beta)$ as in (6.4.2) and $\mathbf{Q}(\beta)$ as in (6.4.7). Clearly, $\gamma(0)$ is just the normalized, uniform distribution on \mathbb{S} : $(\gamma(0))_i = L^{-1}$, where $L \equiv \#\mathbb{S} \ge 2$ is the number of elements in \mathbb{S} . On the one hand, as β gets larger, $\gamma(\beta)$ becomes more concentrated on \mathbb{S}_0 . More precisely, since $Z(\beta) \ge \#\mathbb{S}_0 \ge 1$

$$\langle \mathbf{1}_{\mathbb{S}_0} \mathbf{C} \rangle_{\boldsymbol{\gamma}(\beta)} \le L e^{-\beta\delta}, \quad \text{where } \delta \equiv \min \{ H(j) : j \notin \mathbb{S}_0 \}.$$
 (6.5.1)

On the other hand, as β gets larger, Theorem 6.4.11 says that, at least when e > 0, $\lambda(\beta)$ will be getting smaller. Thus, we are confronted by a conflict.

In view of the introductory discussion, this conflict between the virtues of taking β large, which is tantamount to adopting an approximately steepest decent strategy, versus those of taking β small, which is tantamount to keeping things fluid and thereby diminishing the danger of getting trapped, should be expected. Moreover, a resolution is suggested at the end of that discussion. Namely, in order to maximize the advantages of each, one should start with $\beta = 0$ and allow β to increase with time.⁵ That is, we will make β an increasing, continuous function $t \rightsquigarrow \beta(t)$ with $\beta(0) = 0$. In the interest of unencumbering our formulae, we will adopt the notation

$$Z(t) = Z(\beta(t)), \qquad \mathbf{\gamma}_t = \mathbf{\gamma}(\beta(t)), \qquad \langle \cdot \rangle_t = \langle \cdot \rangle_{\mathbf{\gamma}_t}, \qquad \| \cdot \|_{2,t} = \| \cdot \|_{2,\mathbf{\gamma}_t},$$
$$\operatorname{Var}_t = \operatorname{Var}_{\mathbf{\gamma}_t}, \qquad \mathbf{Q}(t) = \mathbf{Q}(\beta(t)), \qquad \mathcal{E}_t = \mathcal{E}_{\beta(t)}, \quad \text{and} \quad \lambda_t = \lambda_{\beta(t)}.$$

Because, in the physical model, β is proportional to the reciprocal of temperature and β increases with time, $t \rightsquigarrow \beta(t)$ is called the *cooling schedule*.

⁵Actually, there is good reason to doubt that monotonically increasing β may not be the best strategy. Indeed, the name "simulated annealing" derives from the idea that what one wants to do is simulate the annealing process familiar to chemists, material scientists, skilled carpenters, and followers of Metropolis. Namely, what these people do is alternately heat and cool to achieve their goal, and there is reason to believe we should be following their example. However, I have chosen not to follow them on the unforgivable, but understandable, grounds that my analysis is capable of handling only the monotone case.

6.5.2 Construction of the Transition Probabilities

The *Q*-matrix here being time dependent means that the associated transition probabilities will be *time-inhomogeneous*. Thus, instead of $t \rightsquigarrow \mathbf{Q}(t)$ determining a one parameter family of transition probability matrices, for each $s \in [0, \infty)$ it will determine a map $t \rightsquigarrow \mathbf{P}(s, t)$ from $[s, \infty)$ into transition probability matrices by the *time-inhomogeneous Kolmogorov forward equation*

$$\frac{d}{dt}\mathbf{P}(s,t) = \mathbf{P}(s,t)\mathbf{Q}(t) \quad \text{on } (s,\infty) \text{ with } \mathbf{P}(s,s) = \mathbf{I}.$$
(6.5.2)

Although (6.5.2) is not exactly covered by our earlier analysis of Kolmogorov equations, it nearly is. To see this, we solve (6.5.2) via an approximation procedure in which $\mathbf{Q}(t)$ is replaced on the right hand side by

$$\mathbf{Q}^{(N)}(t) \equiv \mathbf{Q}(\lfloor t \rfloor_N)$$
 where $\lfloor t \rfloor_N = \frac{m}{N}$ for $t \in \left[\frac{m}{N}, \frac{(m+1)}{N}\right]$.

The solution $t \rightsquigarrow \mathbf{P}^{(N)}(s, t)$ to the resulting equation is then given by the prescription $\mathbf{P}^{(N)}(s, s) = \mathbf{I}$ and

$$\mathbf{P}^{(N)}(s,t) = \mathbf{P}^{(N)}(s,s \vee \lfloor t \rfloor_N) \sum_{m=0}^{\infty} \frac{(t-s \vee \lfloor t \rfloor_N)^m}{m!} \mathbf{Q}(s \vee \lfloor t \rfloor_N)^m \quad \text{for } t > s.$$

As this construction makes obvious, $\mathbf{P}^{(N)}(s, t)$ is a transition probability matrix for each $N \ge 1$ and $t \ge s$, and $(s, t) \rightsquigarrow \mathbf{P}^{(N)}(s, t)$ is continuous. Moreover,

$$\begin{aligned} \left\| \mathbf{P}^{(N)}(s,t) - \mathbf{P}^{(M)}(s,t) \right\|_{u,v} \\ &\leq \int_{s}^{t} \left\| \mathbf{Q}(\lfloor \tau \rfloor_{N}) - \mathbf{Q}(\lfloor \tau \rfloor_{M}) \right\|_{u,v} \left\| \mathbf{P}^{(N)}(s,\tau) \right\|_{u,v} d\tau \\ &+ \int_{s}^{t} \left\| \mathbf{Q}(\lfloor \tau \rfloor_{M}) \right\|_{u,v} \left\| \mathbf{P}^{(N)}(s,\tau) - \mathbf{P}^{(M)}(s,\tau) \right\|_{u,v} d\tau \end{aligned}$$

But

 $\left\|\mathbf{Q}(\tau)\right\|_{\mathbf{u},\mathbf{v}} \leq \|\mathbf{A}\|_{\mathbf{u},\mathbf{v}} \quad \text{and} \quad \left\|\mathbf{Q}(\tau') - \mathbf{Q}(\tau)\right\|_{\mathbf{u},\mathbf{v}} \leq \|\mathbf{A}\|_{\mathbf{u},\mathbf{v}}\|H\|_{\mathbf{u}}\left|\beta(\tau') - \beta(\tau)\right|,$

and so

$$\begin{aligned} \left\| \mathbf{P}^{(N)}(s,t) - \mathbf{P}^{(M)}(s,t) \right\|_{u,v} &\leq \|\mathbf{A}\|_{u,v} \|H\|_{u} \int_{s}^{t} \left| \beta \left(\lfloor \tau \rfloor_{N} \right) - \beta \left(\lfloor \tau \rfloor_{M} \right) \right| d\tau \\ &+ \|\mathbf{A}\|_{u,v} \int_{s}^{t} \left\| \mathbf{P}^{(N)}(\tau) - \mathbf{P}^{(M)}(\tau) \right\|_{u,v} d\tau. \end{aligned}$$

Hence, after an application of Gronwall's inequality, we find that

$$\sup_{0 \le s \le t \le T} \left\| \mathbf{P}^{(N)}(s,t) - \mathbf{P}^{(M)}(s,t) \right\|_{u,v}$$

$$\leq \|\mathbf{A}\|_{u,v} \|H\|_{u} e^{\|\mathbf{A}\|_{u,v}T} \int_{0}^{T} \left| \beta \left(\lfloor \tau \rfloor_{N} \right) - \beta \left(\lfloor \tau \rfloor_{M} \right) \right| d\tau.$$

Because $\tau \rightsquigarrow \beta(\tau)$ is continuous, this proves that the sequence { $\mathbf{P}^{(N)}(s, t) : N \ge 1$ } is Cauchy convergent in the sense that, for each T > 0,

$$\lim_{M \to \infty} \sup_{N \ge M} \sup_{0 \le s \le t \le T} \left\| \mathbf{P}^{(N)}(s,t) - \mathbf{P}^{(M)}(s,t) \right\|_{u,v} = 0.$$

As a consequence, we know that there exists a continuous $(s, t) \rightsquigarrow \mathbf{P}(s, t)$ to which the $\mathbf{P}^{(N)}(s, t)$'s converge with respect to $\|\cdot\|_{u,v}$ uniformly on finite intervals. In particular, for each $t \ge s$, $\mathbf{P}(s, t)$ is a transition probability matrix and $t \in [s, \infty) \mapsto$ $\mathbf{P}(s, t)$ is a continuous solution to

$$\mathbf{P}(s,t) = \mathbf{I} + \int_{s}^{t} \mathbf{P}(s,\tau) \mathbf{Q}(\tau) d\tau, \quad t \in [s,\infty),$$

which is the equivalent integrated form of (6.5.2). Furthermore, if $t \in [s, \infty) \mapsto \mu_t \in M_1(\mathbb{S})$ is continuously differentiable, then

$$\dot{\boldsymbol{\mu}}_t \equiv \frac{d}{dt} \boldsymbol{\mu}_t = \boldsymbol{\mu}_t \mathbf{Q}(t) \quad \text{for } t \in [s, \infty) \quad \Longleftrightarrow \quad \boldsymbol{\mu}_t = \boldsymbol{\mu}_s \mathbf{P}(s, t) \quad \text{for } t \in [s, \infty).$$
(6.5.3)

Since the "if" assertion is trivial, we turn to the "only if" statement. Thus, suppose that $t \in [s, \infty) \mapsto \mu_t \in M_1(\mathbb{S})$ satisfying $\dot{\mu}_t = \mu_t \mathbf{Q}(t)$ is given, and set $\boldsymbol{\omega}_t = \mu_t - \mu_s \mathbf{P}(s, t)$. Then

$$\boldsymbol{\omega}_t = \int_s^t \boldsymbol{\omega}_\tau \mathbf{Q}(\tau) \, d\tau,$$

and so,

$$\|\boldsymbol{\omega}_t\|_{\mathbf{v}} \leq \|\mathbf{A}\|_{\mathbf{u},\mathbf{v}} \int_s^t \|\boldsymbol{\omega}_{\tau}\|_{\mathbf{v}} d\tau.$$

Hence, after another application of Gronwall's inequality, we see that $\omega_t = 0$ for all $t \ge s$.

Of course, by applying this uniqueness result when $\mu_t = \delta_i \mathbf{P}(s, t)$ for each $i \in \mathbb{S}$, we learn that $(s, t) \rightsquigarrow \mathbf{P}(s, t)$ is the one and only solution to (6.5.2). In addition, it leads to the following time-inhomogeneous version of the Chapman–Kolmogorov equation:

$$\mathbf{P}(r,t) = \mathbf{P}(r,s)\mathbf{P}(s,t) \quad \text{for } 0 \le r \le s \le t.$$
(6.5.4)

Indeed, set $\mu_t = \delta_i \mathbf{P}(r, t)$ for $t \ge s$, note that $t \rightsquigarrow \mu_t$ satisfies (6.5.3) with $\mu_s = \delta_i \mathbf{P}(r, s)$, and conclude that $\mu_t = \delta_i \mathbf{P}(r, s) \mathbf{P}(s, t)$.

6.5.3 Description of the Markov Process

Given a probability vector μ , we now want to construct a Markov process $\{X(t) : t \ge 0\}$ that has μ as its initial distribution and $(s, t) \rightsquigarrow \mathbf{P}(s, t)$ as its transition mechanism, in the sense that

$$\mathbb{P}(X(0) = i) = (\boldsymbol{\mu})_i \& \mathbb{P}(X(t) = j \mid X(\sigma), \ \sigma \in [0, s]) = \mathbf{P}(s, t)_{X(s)j}.$$
(6.5.5)

The idea which we will use is basically the same as the one which we used in Sects. 5.2.1 & 2.1.1. However, life here is made more complicated by the fact that the time inhomogeneity forces us to have an uncountable number of random variables at hand: a pair for each $(t, i) \in [0, \infty) \times \mathbb{S}$. To handle this situation, take, without loss in generality, $\mathbb{S} = \{1, \ldots, L\}$ and, for $(t, i, j) \in [0, \infty) \times \mathbb{S}^2$, set $S(t, i, j) = \sum_{\ell=1}^{j} e^{-\beta(t)(H(\ell) - H(i))^+} (\mathbf{A})_{i\ell}$, take S(t, i, 0) = 0, and define

$$\Psi(t, i, u) = \begin{cases} j & \text{if } \frac{S(t, i, j-1)}{S(t, i, L)} \le u < \frac{S(t, i, j)}{S(t, i, L)} \\ i & \text{if } u > 1. \end{cases}$$

Also, determine $\mathcal{T}: [0,\infty) \times \mathbb{S} \times [0,\infty) \longrightarrow [0,\infty)$ by

$$\int_{s}^{s+\mathcal{T}(s,i,\xi)} S(\tau,i,L) \, d\tau = \xi$$

Next, take X_0 to be an S-valued random variable with distribution μ , let $\{E_n : n \ge 1\}$ be a sequence of unit exponential random variables that are independent of each other and of X_0 , and let $\{U_n : n \ge 1\}$ be a sequence of mutually independent random variables that are uniformly distributed on [0, 1) and are independent of $\sigma(\{X_0\} \cup \{E_n : n \ge 1\})$. Finally, set $J_0 = 0$ and $X(0) = X_0$, and, when $n \ge 1$, use induction to define

$$J_n - J_{n-1} = \mathcal{T}(J_{n-1}, X(J_{n-1}), E_n), \qquad X(J_n) = \Psi(J_n, X(J_{n-1}), U_n), \text{ and}$$
$$X(t) = X(J_{n-1}) \quad \text{for } J_{n-1} \le t < J_n.$$

Without substantial change, the reasoning given in Sect. 2.1.1 combined with that in Sect. 5.2.2 and Exercise 5.5.12 allow one to show that (6.5.5) holds.

6.5.4 Choosing a Cooling Schedule

In this section, we will give a rational basis on which to choose the cooling schedule $t \rightsquigarrow \beta(t)$. For this purpose, it is essential to keep in mind what it is that we are attempting to do. Namely, we are trying to have the Markov process $\{X(t) : t \ge 0\}$ seek out the set $\mathbb{S}_0 = \{j : H(j) = 0\}$ in the sense that, as $t \to \infty$, $\mathbb{P}(X(t) \notin \mathbb{S}_0)$ should tend to 0 as fast as possible, and the way we hope to accomplish this is by making the distribution of X(t) look as much like γ_t as possible. Thus, on the one hand, we need to give $\{X(t) : t \ge 0\}$ enough time to equilibrate, so that the distribution of X(t) will look a lot like γ_t . On the other hand, in spite of the fact that it may inhibit equilibration, unless we make $\beta(t)$ increase to infinity, there is no reason for our wanting to make the distribution of X(t) look like γ_t .

In order to understand how to deal with the concerns raised above, let μ be a fixed initial distribution, and let μ_t be the distribution at time $t \ge 0$ of the Markov process $\{X(t) : t \ge 0\}$ described in the Sect. 6.5.3 with initial distribution μ . Equivalently, $\mu_t = \mu \mathbf{P}(0, t)$, where $\{\mathbf{P}(s, t) : 0 \le s \le t < \infty\}$ is the family of transition probability matrices constructed in Sect. 6.5.2. Next, define $f_t : \mathbb{S} \longrightarrow [0, \infty)$ so that

$$f_t(i) = \frac{(\boldsymbol{\mu}_t)_i}{(\boldsymbol{\gamma}_t)_i} \quad \text{for } t \ge 0 \text{ and } i \in \mathbb{S}.$$

It should be obvious that the size of f_t provides a good measure of the extent to which μ_t resembles γ_t . For example, by Schwarz's inequality and (6.5.1),

$$\mathbb{P}(X(t) \notin \mathbb{S}_{0}) = \langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}} \rangle_{\mu_{t}} = \langle f_{t} \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}} \rangle_{t}$$

$$\leq \|f_{t}\|_{2,t} \sqrt{\langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}} \rangle_{t}} \leq L^{\frac{1}{2}} \|f_{t}\|_{2,t} e^{-\frac{\beta(t)\delta}{2}}, \qquad (6.5.6)$$

and so we will have made progress if we can keep $||f_t||_{2,t}$ under control.

With the preceding in mind, assume that $t \rightsquigarrow \beta(t)$ is continuously differentiable, note that this assumption makes $t \rightsquigarrow ||f_t||_{2,t}^2$ also continuously differentiable, and, in fact, that

$$\begin{aligned} \frac{d}{dt} \|f_t\|_{2,t}^2 &= \frac{d}{dt} \left(\frac{1}{Z(t)} \sum_{i \in \mathbb{S}} (f_t(i))^2 e^{-\beta(t)H(i)} \right) \\ &= 2\langle f_t, \dot{f}_t \rangle_t - \dot{\beta}(t) \langle H - \langle H \rangle_t, f_t^2 \rangle_t, \end{aligned}$$

since, by (6.4.3), $\dot{Z}(t) = -\dot{\beta}(t)Z(t)\langle H \rangle_t$. On the other hand, we can compute this same derivative another way. Namely, because $\langle \mathbf{P}(0,t)g \rangle_{\boldsymbol{\mu}} = \langle g \rangle_{\boldsymbol{\mu}_t} = \langle f_t, g \rangle_t$ for any function g,

$$\|f_t\|_{2,t}^2 = \langle f_t \rangle_{\boldsymbol{\mu}_t} = \langle \mathbf{P}(0,t) f_t \rangle_{\boldsymbol{\mu}},$$

and so we can use (6.5.2) to see that

$$\frac{d}{dt} \|f_t\|_{2,t}^2 = \left\langle \mathbf{P}(0,t)\mathbf{Q}(t)f_t \right\rangle_{\boldsymbol{\mu}} + \left\langle \mathbf{P}(0,t)\dot{f}_t \right\rangle_{\boldsymbol{\mu}} = -\mathcal{E}_t(f_t,f_t) + \langle f_t,\dot{f}_t \rangle_t.$$

Thus, after combining these to eliminate the term containing f_t , we arrive at

$$\begin{aligned} \frac{d}{dt} \|f_t\|_{2,t}^2 &= -2\mathcal{E}_t(f_t, f_t) + \dot{\beta}(t) \langle H - \langle H \rangle_t, f_t^2 \rangle_t \\ &\leq -(2\lambda_t - \|H\|_{\mathbf{u}} \dot{\beta}(t)) \|f_t\|_{2,t}^2 + 2\lambda_t, \end{aligned}$$

where, in passing to the second line, we have used the fact that $\langle f_t \rangle_t = 1$ and therefore that $\operatorname{Var}_t(f) = \|f_t\|_{2,t}^2 - 1$. Putting all this together, we now know that

$$\|H\|_{\mathbf{u}}\dot{\beta}(t) \leq \lambda_t \implies \frac{d}{dt} \|f_t\|_{2,t}^2 \leq -\lambda_t \|f_t\|_{2,t}^2 + 2\lambda_t.$$

The preceding differential inequality for $||f||_{2,t}^2$ is easy to integrate. Namely, it says that

$$\frac{d}{dt} \left(e^{\Lambda(t)} \| f_t \|_{2,t}^2 \right) \le 2\lambda_t e^{\Lambda(t)} \quad \text{where } \Lambda(t) = \int_0^t \lambda_\tau \, d\tau.$$

Hence,

$$\|f_t\|_{2,t}^2 \le e^{-\Lambda(t)} \|f_0\|_{2,0}^2 + 2(1 - e^{-\Lambda(t)}) \le \|f_0\|_{2,0}^2 \lor 2.$$

Moreover, since $(\gamma_0)_i = L^{-1}$, where $L = \# \mathbb{S} \ge 2$, $\|f_0\|_{2,0}^2 \le \|f_0\|_u \langle f \rangle_0 \le L$, and so

$$\|H\|_{\mathbf{u}}\dot{\beta}(t) \le \lambda_t \implies \|f_t\|_{2,t} \le L^{\frac{1}{2}}.$$
(6.5.7)

The final step is to figure out how to choose $t \rightsquigarrow \beta(t)$ so that it satisfies the condition in (6.5.7); and, of course, we are only interested in the case when $\mathbb{S}_0 \neq \mathbb{S}$ or, equivalently, $||H||_u > 0$. From Theorem 6.4.11 we know that $\lambda_t \ge c_- e^{-\beta(t)\mathfrak{e}}$. Hence, we can take

$$\beta(t) = \begin{cases} \frac{1}{\mathfrak{e}} \log(1 + \frac{c_{-}\mathfrak{e}t}{\|H\|_{\mathfrak{u}}}) & \text{when } \mathfrak{e} > 0\\ \frac{c_{-t}}{\|H\|_{\mathfrak{u}}} & \text{when } \mathfrak{e} = 0. \end{cases}$$
(6.5.8)

After putting this together with (6.5.7) and (6.5.6), we have now proved that when $\beta(t)$ is given by (6.5.8), then

$$\mathbb{P}(X(t) \notin \mathbb{S}_0) \le L \begin{cases} (1 + \frac{c - \mathfrak{e}t}{\|H\|_u})^{-\alpha} & \text{with } \alpha = \frac{\delta}{2\mathfrak{e}} \text{ when } \mathfrak{e} > 0\\ e^{-\frac{\delta c - t}{2\|H\|_u}} & \text{when } \mathfrak{e} = 0. \end{cases}$$
(6.5.9)

Remark The result in (6.5.9) when e = 0 deserves some further comment. In particular, it should be observed that e = 0 does not guarantee success for a steepest decent strategy. Indeed, e = 0 only means that each *i* can be connected to S₀ by an allowable path along which *H* is non-increasing (cf. Exercise 6.6.13), it does not rule out the possibility that, when using steepest decent, one will choose a bad path and get stuck. Thus, even in this situation, one needs enough randomness to hunt around until one finds a good path.

6.5.5 Small Improvements

An observation, which is really only of interest when e > 0, is that one carry out the same sort of analysis to control the $||f_t||_{q,t} \equiv (\langle |f_t|^q \rangle_t)^{\frac{1}{q}}$ for each $q \in [2, \infty)$. The reason for wanting to do so is that it allows one to improve the rate of convergence in (6.5.6) and therefore in (6.5.9). Namely, by (6.6.3),

$$\langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}\rangle_{\boldsymbol{\mu}_{t}} = \langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}\rangle_{t} \frac{\langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}f_{t}\rangle_{t}}{\langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}\rangle_{t}} \leq \langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}\rangle_{t} \left(\frac{\langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}f_{t}^{q}\rangle_{t}}{\langle \mathbf{1}_{\mathbb{S}_{0}\mathbb{C}}\rangle_{t}}\right)^{\frac{1}{q}} \leq L^{\frac{1}{q}} \|f_{t}\|_{q,t} e^{-\frac{\beta(t)\delta}{q'}},$$

where $q' \equiv \frac{q}{q-1}$ is the Hölder conjugate of q. Hence, if $\theta \in [\frac{1}{2}, 1)$ and $q = \frac{1}{1-\theta}$, then one can choose a cooling schedule which makes $\mathbb{P}(X(t) \notin \mathbb{S}_0)$ go to 0 at least as fast as $t^{-\frac{\theta\delta}{\varepsilon}}$.

To carry this out, one begins by computing $\frac{d}{dt} ||f_t||_{q,t}^q$ twice, once for each of the following expressions:

$$\|f_t\|_{q,t}^q = \langle f_t^q \rangle_t \quad \text{and} \quad \|f_t\|_{q,t}^q = \langle \mathbf{P}(0,t)f_t^{q-1} \rangle_{\boldsymbol{\mu}}.$$

One then eliminates \dot{f}_t from the resulting equations and thereby arrives at

$$\frac{d}{dt} \|f_t\|_{q,t}^q = -q\mathcal{E}_t(f_t, f_t^{q-1}) - \frac{\dot{\beta}(t)}{q'-1} \langle f_t^q, H - \langle H \rangle_t \rangle_t$$

where $q' = \frac{q}{q-1}$ and

$$\mathcal{E}_{t}(\varphi,\psi) \equiv \frac{1}{2} \sum_{i \neq j} (\boldsymbol{\gamma}_{t})_{i} \big(\mathbf{Q}(t) \big)_{ij} \big(\varphi(j) - \varphi(i) \big) \big(\psi(j) - \psi(i) \big) = - \big\langle \varphi, \mathbf{Q}(t) \psi \big\rangle_{t}.$$

At this point one has to show that

$$\mathcal{E}_t(f_t, f_t^{q-1}) \geq \frac{4(q-1)}{q^2} \mathcal{E}_t(f_t^{\frac{q}{2}}, f_t^{\frac{q}{2}}),$$

and a little thought makes it clear that this inequality comes down to checking that, for any pair $(a, b) \in [0, \infty)^2$,

$$\left(b^{\frac{q}{2}}-a^{\frac{q}{2}}\right)^2 \le \frac{q^2}{4(q-1)}(b-a)\left(b^{q-1}-a^{q-1}\right),$$

which, when looked at correctly, follows from the fundamental theorem of calculus plus Schwarz's inequality. Hence, in conjunction with the preceding and (6.3.6), we find that

$$\frac{d}{dt} \|f_t\|_{q,t}^q \le -\frac{1}{q'} \Big(4\lambda_t - q\dot{\beta}(t) \|H\|_{\mathsf{u}} \Big) \|f_t\|_{q,t}^q + \frac{4\lambda_t}{q'} \Big\langle f_t^{\frac{q}{2}} \Big\rangle_t^2.$$

In order to proceed further, we must learn how to control $\langle f_t^{\frac{q}{2}} \rangle_t^2$ in terms of $||f_t||_{q,t}^q$. In the case when q = 2, $\langle f_t^{\frac{q}{2}} \rangle_t^2$ caused no problem because we knew it was equal to 1. When q > 2, we no longer have so much control over it. Nonetheless, by first writing $\langle f_t^{\frac{q}{2}} \rangle_t^2 = \langle f_t^{\frac{q}{2}-1} \rangle_{\mu_t}^2$ and then using part (c) of Exercise 6.6.2 to see that

$$\left\langle f_{t}^{\frac{q}{2}-1}\right\rangle_{\boldsymbol{\mu}_{t}}^{2} \leq \left\langle f_{t}^{q-1}\right\rangle_{\boldsymbol{\mu}_{t}}^{\frac{q-2}{q-1}} = \left\langle f_{t}^{q}\right\rangle_{t}^{\frac{q-2}{q-1}},$$

we arrive at $\langle f_t^{\frac{q}{2}} \rangle_t^2 \leq \langle f_t^{q} \rangle_t^{\frac{q-2}{q-1}}$. Armed with this estimate, we obtain the differential inequality

$$\frac{d}{dt}\|f_t\|_{q,t}^q \le -\frac{1}{q'} \left(4\lambda_t - q\|H\|_{\mathbf{u}}\dot{\beta}(t)\right)\|f_t\|_{q,t}^q + \frac{4\lambda_t}{q'} \left(\|f_t\|_{q,t}^q\right)^{1-\frac{q'}{q}}.$$

Finally, by taking $\beta(t) = \frac{1}{e} \log(1 + \frac{3c - et}{q ||H||_u})$, the preceding inequality can be replaced by

$$\frac{d}{dt} \|f_t\|_{q,t}^q \le -\frac{\lambda_t}{q'} \|f_t\|_{q,t}^q + \frac{4\lambda_t}{q'} \left(\|f_t\|_{q,t}^q\right)^{1-\frac{q'}{q}},$$

which, after integration, can be made to yield

$$||f_t||_{q,t} \le 2^{\frac{1}{q'}} \lor ||f_0||_{q,0},$$

at which point the rest of the argument is the same as the one when q = 2.

Remark Actually, it is possible to do even better if one is prepared to combine the preceding line of reasoning, which is basically a consequence of Poincaré's inequality, with analytic ideas which come under the general heading of Sobolev inequalities. The interested reader might want to consult [3], which is the source from which the contents of this whole section are derived.

6.6 Exercises

Exercise 6.6.1 A function $\psi : [0, \infty) \longrightarrow \mathbb{R}$ is said to be *convex* if the graph of ψ lies below the secant connecting any pair of points on its graph. That is, if it satisfies

$$\psi((1-\theta)s+\theta t) \le (1-\theta)\psi(s) + \theta\psi(t) \quad \text{for all } 0 \le s < t \text{ and } \theta \in [0,1]. \quad (*)$$

This exercise deals with various properties of convex functions, all of which turn on the property that the slope of a convex function is non-decreasing.

(a) If {ψ_n}[∞]₁ ∪ {ψ} are functions on [0, ∞) and ψ_n(t) → ψ(t) for each t ∈ [0, ∞), show that ψ is non-increasing if each of the ψ_n's is and that ψ is convex if each of the ψ_n's is.

(b) If \u03c6 : [0,∞) → R is continuous and twice continuously differentiable on (0,∞), show that \u03c6 is convex on [0,∞) if and only if \u03c6 ≥ 0 on (0,∞).

Hint: The "only if" part is an easy consequence of

$$\ddot{\psi}(t) = \lim_{h \searrow 0} \frac{\psi(t+h) + \psi(t-h) - 2\psi(t)}{h^2} \quad \text{for } t \in (0,\infty).$$

To prove the "if" statement, let 0 < s < t be given, and for $\epsilon > 0$ set

$$\varphi_{\epsilon}(\theta) = \psi \left((1-\theta)s + \theta t \right) - (1-\theta)\psi(s) - \theta \psi(s) - \epsilon \theta (1-\theta), \quad \theta \in [0,1].$$

Note that $\varphi_{\epsilon}(0) = 0 = \varphi_{\epsilon}(1)$ and that $\ddot{\varphi}_{\epsilon} > 0$ on (0, 1). Hence, by the second derivative test, φ_{ϵ} cannot achieve a maximum value in (0, 1). Now let $\epsilon \searrow 0$.

(c) If ψ is convex on $[0, \infty)$, show that, for each $s \in [0, \infty)$,

$$t \in (s, \infty) \longmapsto \frac{\psi(s) - \psi(t)}{t - s}$$
 is non-increasing.

(d) If ψ is convex on $[0, \infty)$ and $0 \le s < t \le u < w$, show that

$$\frac{\psi(s) - \psi(t)}{t - s} \ge \frac{\psi(u) - \psi(w)}{w - u}$$

Hint: Reduce to the case when u = t.

Exercise 6.6.2 Given a probability vector $\boldsymbol{\mu} \in [0, 1]^{\mathbb{S}}$, there are many ways to prove that $\langle f \rangle_{\boldsymbol{\mu}}^2 \leq \langle f^2 \rangle_{\boldsymbol{\mu}}$ for any $f \in L^2(\boldsymbol{\mu})$. For example, one can get this inequality as an application of Schwarz's inequality $|\langle f, g \rangle_{\boldsymbol{\mu}}| \leq ||f||_{2,\boldsymbol{\mu}} ||g||_{2,\boldsymbol{\mu}}$ by taking $g = \mathbf{1}$. Alternatively, one can use $0 \leq \operatorname{Var}_{\boldsymbol{\mu}}(f) = \langle f^2 \rangle_{\boldsymbol{\mu}} - \langle f \rangle_{\boldsymbol{\mu}}^2$. However, neither of these approaches reveals the essential role that convexity plays here. Namely, the purpose of this exercise is to show that for any non-decreasing, continuous, convex function $\boldsymbol{\psi} : [0, \infty) \longrightarrow [0, \infty)$ and any $f : \mathbb{S} \longrightarrow [0, \infty)$,

$$\psi(\langle f \rangle_{\boldsymbol{\mu}}) \le \langle \psi \circ f \rangle_{\boldsymbol{\mu}}, \tag{6.6.3}$$

where the meaning of the left hand side when $\langle f \rangle_{\mu} = \infty$ is given by taking $\psi(\infty) \equiv \lim_{t \neq \infty} \psi(t)$. The inequality (6.6.3) is an example of more general statement known as *Jensen's inequality* (cf. Theorem 2.4.15 in [8]).

(a) Use induction on $n \ge 2$ to show that

$$\psi\left(\sum_{k=1}^{n} \theta_k x_k\right) \le \sum_{k=1}^{n} \theta_k \psi(x_k)$$

for all $(\theta_1, \dots, \theta_n) \in [0, 1]^n$ with $\sum_{m=1}^{n} \theta_k = 1$ and $(x_1, \dots, x_n) \in [0, \infty)^n$.

(b) Let $\{F_N\}_1^\infty$ be a non-decreasing exhaustion of \mathbb{S} by finite sets satisfying $\mu(F_N) \equiv \sum_{i \in F_N} (\boldsymbol{\mu})_i > 0$, apply part (a) to see that

$$\psi\left(\frac{\sum_{i\in F_N} f(i)(\boldsymbol{\mu})_i}{\mu(F_N)}\right) \leq \frac{\sum_{i\in F_N} \psi(f(i))(\boldsymbol{\mu})_i}{\mu(F_N)} \leq \frac{\langle\psi\circ f\rangle_{\boldsymbol{\mu}}}{\mu(F_N)}$$

for each N, and get the asserted result after letting $N \to \infty$.

- (c) As an application of (6.6.3), show that, for any $0 and <math>f : \mathbb{S} \longrightarrow [0, \infty), \langle f^p \rangle_{\mu}^{\frac{1}{p}} \le \langle f^q \rangle_{\mu}^{\frac{1}{q}}.$
- (d) Let $p \in (0, \infty)$ and $\{a_{\ell} : 1 \le \ell \le m\} \subseteq [0, \infty)$. Show that

$$\left(\sum_{\ell=1}^m a_\ell\right)^p \le m^{(p-1)^+} \sum_{\ell=1}^m a_\ell^p.$$

Hint: When $p \ge 1$, use part (c) with p = 1 and q = p for the measure that assigns mass $\frac{1}{m}$ to the points $1, \ldots, m$ and the function $f(\ell) = a_{\ell}$. When $p \in (0, 1)$, first reduce to the case when n = 2, and in that case show that it suffices to show that $(1 + x)^p \le 1 + x^p$ for $x \ge 0$.

Exercise 6.6.4 This exercise deals with the material in Sect. 6.1.2 and demonstrates that, with the exception of (6.1.10), more or less everything in that section extends to general irreducible, positive recurrent **P**'s, whether or not they are reversible. Again let $\pi = \pi^{\mathbb{S}}$ be the unique **P**-stationary probability vector. In addition, for the exercise which follows, it will be important to consider the space $L^2(\pi; \mathbb{C})$ consisting of those $f: \mathbb{S} \longrightarrow \mathbb{C}$ for which $|f| \in L^2(\pi)$.

- (a) Define \mathbf{P}^{\top} as in (6.1.2), show that $1 \ge (\mathbf{P}\mathbf{P}^{\top})_{ii} = (\pi)_i \sum_{j \in \mathbb{S}} \frac{(\mathbf{P})_{ij}^2}{(\pi)_j}$, and conclude that the series in the definition $\mathbf{P}f(i) \equiv \sum_{j \in \mathbb{S}} f(j)(\mathbf{P})_{ij}$ is absolutely convergent for each $i \in \mathbb{S}$ and $f \in L^2(\pi; \mathbb{C})$.
- (b) Show that $\|\mathbf{P}f\|_{2,\pi} \le \|f\|_{2,\pi}$ for all $f \in L^2(\pi; \mathbb{C})$, and conclude that (6.1.8) and (6.1.9) extend to the present setting for all $f \in L^2(\pi; \mathbb{C})$.

Exercise 6.6.5 Continuing with the program initiated in Exercise 6.6.4, we will now see that reversibility plays only a minor role in Sect. 6.1.4. Thus, let **P** be any irreducible transition probability on S which is positive recurrent, let *d* be its period, and set $\theta_d = e^{\sqrt{-1}2\pi d^{-1}}$.

(a) Show that for each $0 \le m < d$ there is a function $f_m : \mathbb{S} \longrightarrow \mathbb{C}$ with the properties that $|f_m| \equiv 1$ and $\mathbf{P} f_m = \theta_d^m f_m$.

Hint: Choose a cyclic decomposition $(\mathbb{S}_0, \ldots, \mathbb{S}_{d-1})$ as in Sect. 4.1.8, and use (4.1.19).

(b) Given $\alpha \in \mathbb{R} \setminus \{0\}$, set $\theta_{\alpha} = e^{\sqrt{-1}2\pi\alpha^{-1}}$, and show that there exists an $f \in L^2(\pi; \mathbb{C}) \setminus \{0\}$ satisfying $\mathbf{P}f = \theta_{\alpha}f$ if and only if $d = m\alpha$ for some $m \in \mathbb{Z} \setminus \{0\}$.

Hint: By part (a), it suffices to show that no f exists unless $d = m\alpha$ for some nonzero $m \in \mathbb{Z}$. Thus, suppose that f exists for some α which is not a rational number of the form $\frac{d}{m}$, choose $i \in \mathbb{S}$ so that $f(i) \neq 0$, and get a contradiction with the fact that $\lim_{n\to\infty} \mathbf{P}^{nd} f(i)$ exists.

- (c) Suppose that f ∈ L²(S; C) is a non-trivial, bounded solution to Pf = θ_d^m f for some m ∈ Z, and let (S₀, ..., S_{d-1}) be a cyclic decomposition of S. Show that, for each 0 ≤ r < d, f ↾ S_r ≡ θ_d^{-rm}c₀, where c₀ ∈ C \ {0}. In particular, up to a multiplicative constant, for each integer 0 ≤ m < d there is exactly one non-trivial, bounded f ∈ L²(S; C) satisfying Pf = θ_d^m f.
- (d) Let *H* denote the linear subspace of f ∈ L²(π; C) satisfying Pf = θf for some θ ∈ C with |θ| = 1. By combining parts (b) and (c), show that *d* is the dimension of *H* as a vector space over C.
- (e) Assume that $(\mathbf{P})_{ij} > 0 \implies (\mathbf{P})_{ji} > 0$ for all $(i, j) \in \mathbb{S}^2$. Show that $d \le 2$ and that d = 2 if and only if there is a non-trivial, bounded $f : \mathbb{S} \longrightarrow \mathbb{R}$ satisfying $\mathbf{P}f = -f$. Thus, this is the only property of reversible transition probability matrices of which we made essential use in Theorem 6.1.14.

Exercise 6.6.6 This exercise provides another way to think about the relationship between non-negative definiteness and aperiodicity. Namely, let **P** be a not necessarily irreducible transition probability on \mathbb{S} , and assume that $\boldsymbol{\mu}$ is a probability vector for which the detailed balance condition $(\boldsymbol{\mu})_i(\mathbf{P})_{ij} = (\boldsymbol{\mu})_j(\mathbf{P})_{ji}$, $(i, j) \in \mathbb{S}^2$ holds. Further, assume that **P** is non-negative definite in $L^2(\boldsymbol{\mu})$: $\langle f, \mathbf{P}f \rangle_{\boldsymbol{\mu}} \ge 0$ for all bounded $f : \mathbb{S} \longrightarrow \mathbb{R}$. Show that $(\boldsymbol{\mu})_i > 0 \implies (\mathbf{P})_{ii} > 0$ and therefore that *i* is aperiodic if $(\boldsymbol{\mu})_i > 0$. What follows are steps which lead to this conclusion.

- (a) Define the matrix **A** so that $(\mathbf{A})_{ij} = \langle \mathbf{1}_{\{i\}}, \mathbf{P1}_{\{j\}} \rangle_{\mu}$, and show that **A** is symmetric (i.e., $(\mathbf{A})_{ij} = (\mathbf{A})_{ji}$) and non-negative definite in the sense that $\sum_{ij} (\mathbf{A})_{ij} (\mathbf{x})_i (\mathbf{x})_j \geq 0$ for any $\mathbf{x} \in \mathbb{R}^{\mathbb{S}}$ with only a finite number of non-vanishing entries.
- (b) Given $i \neq j$, consider the plane $\{\alpha \mathbf{1}_{\{i\}} + \beta \mathbf{1}_{\{j\}} : \alpha, \beta \in \mathbb{R}\}$ in $L^2(\pi)$, and, using the argument with which we derived (6.1.5), show that $(\mathbf{A})_{ij}^2 \leq (\mathbf{A})_{ii}(\mathbf{A})_{jj}$. In particular, if $(\mathbf{A})_{ii} = 0$, then $(\mathbf{A})_{ij} = 0$ for all $j \in \mathbb{S}$.
- (c) Complete the proof by noting that $\sum_{j \in \mathbb{S}} (\mathbf{A})_{ij} = (\boldsymbol{\mu})_i$.
- (d) After examining the argument, show that we did not need **P** to be non-negative definite but only that, for a given $i \in S$, each of the 2 × 2 submatrices

$$\begin{pmatrix} \langle \mathbf{1}_{\{i\}}, \mathbf{P}\mathbf{1}_{\{i\}}\rangle_{\boldsymbol{\mu}} & \langle \mathbf{1}_{\{i\}}, \mathbf{P}\mathbf{1}_{\{j\}}\rangle_{\boldsymbol{\mu}} \\ \langle \mathbf{1}_{\{j\}}, \mathbf{P}\mathbf{1}_{\{i\}}\rangle_{\boldsymbol{\mu}} & \langle \mathbf{1}_{\{j\}}, \mathbf{P}\mathbf{1}_{\{j\}}\rangle_{\boldsymbol{\mu}} \end{pmatrix}$$

be.

Exercise 6.6.7 Let **P** be an irreducible, positive recurrent transition probability with stationary probability π . Refer to (6.1.2), and show that the first construction in (6.1.3) is again irreducible whereas the second one need not be. Also, show that

the period of the first construction is never greater than that of **P** and that the second construction is always non-negative definite. Thus, if $\mathbf{P}^{\top}\mathbf{P}$ is irreducible, it is necessarily aperiodic.

Exercise 6.6.8 Assume that **P** is irreducible and reversible with respect to π . Because the estimate in (6.2.5) is so weak, it is better to avoid it whenever possible. One way to avoid it is to note that **P**² is necessarily non-negative definite and therefore that the associated β equals β_+ .

- (a) Show \mathbf{P}^2 is irreducible if and only if \mathbf{P} is aperiodic.
- (b) Assume that **P** is aperiodic, define

$$\mathcal{E}^{(2)}(f,f) = \frac{1}{2} \sum_{i,j \in \mathbb{S}} (\boldsymbol{\pi})_i (\mathbf{P}^2)_{ij} (f(j) - f(i))^2,$$

and set

$$\beta^{(2)} = \inf \{ \mathcal{E}(f, f) : f \in L^2_0(\pi) \text{ and } \| f \|_{2,\pi} = 1 \}$$

Show that

$$\left\|\mathbf{P}^{n}f-\langle f\rangle\right)_{\boldsymbol{\pi}}\right\|_{2,\boldsymbol{\pi}}\leq\left(1-\beta^{(2)}\right)^{\lfloor\frac{n}{2}\rfloor}$$

Exercise 6.6.9 As in Sect. 3.3, let $\Gamma = (V, \mathcal{E})$ be a connected, finite graph in which no vertex has an edge to itself and there are no double edges, and take $\mathbf{P}_{vw} = \frac{1}{d_v}$, where d_v is the degree of v if $\{v, w\} \in \mathcal{E}$ and $\mathbf{P}_{v,w} = 0$ if $\{v, w\} \notin \mathcal{E}$. As was observed earlier, the associated Markov chain is irreducible.

- (a) Show that the chain is aperiodic if and only if the graph is not bipartite. That is, if and only if there is no non-empty V' ⊊ V with the property that every edge connects a point in V' to one in V \ V'.
- (b) Determine the probability vector $\boldsymbol{\pi}$ by $(\boldsymbol{\pi})_v = \frac{d_v}{2\#\mathcal{E}}$, and show that $\boldsymbol{\pi}$ is a reversible probability vector for **P**.
- (c) Choose a set $\mathcal{P} = \{p(v, w) : v \neq w\}$ of allowable paths, as in Sect. 6.2.2, and show that

$$\beta_{+} \ge \frac{2\#\mathcal{E}}{D^{2}L(\mathcal{P})B(\mathcal{P})},\tag{6.6.10}$$

where $D \equiv \max_{v \in V} d_v$, $L(\mathcal{P})$ is the maximal length (i.e., the maximal number of edges in) of the paths in \mathcal{P} , and $B(\mathcal{P})$, the *bottleneck coefficient*, is the maximal number of paths $p \in \mathcal{P}$ which cross over an edge $e \in E$. Obviously, if one chooses \mathcal{P} to consist of geodesics (i.e., paths of minimal length connecting their end points), then $L(\mathcal{P})$ is just the diameter of Γ , and, as such, is as small as possible. On the other hand, because it may force there to be bad bottlenecks (i.e., many paths traversing a given edge), choosing geodesics may not be optimal. (d) Assuming that the graph is not bipartite, choose $\mathcal{P} = \{p(v) : v \in V\}$ to be a set of allowable closed paths of odd length, and show that

$$\beta_{-} \ge \frac{2}{DL(\mathcal{P})B(\mathcal{P})}.$$
(6.6.11)

Exercise 6.6.12 Here is an example to which the considerations in Exercises 6.6.9 and 6.6.8 apply and give pretty good results. Let $N \ge 2$, take *V* to be the subset of the complex plane \mathbb{C} consisting of the *N*th roots of unity, and let \mathcal{E} be the collection of pairs of adjacent roots of unity. That is, pairs of the form $\{e^{\frac{\sqrt{-12\pi}(m-1)}{N}}, e^{\frac{\sqrt{-12\pi}m}{N}}\}$. Finally, take **P** and π accordingly, as in Exercise 6.6.9.

- (a) Show that \mathbf{P} is always irreducible but that it is aperiodic if and only if N is odd.
- (b) Assuming that N is odd, use (c) in the Exercise 6.6.9 to show that $\beta_+ \ge \frac{2}{N}$ and $\beta_- \ge \frac{2}{N^2}$.
- (c) Again assume that N is odd, and use Exercise 6.6.8 to show that

$$\left\|\mathbf{P}^{n}f - \langle f \rangle_{\boldsymbol{\pi}}\right\|_{2,\boldsymbol{\pi}} \leq \left(1 - \frac{4}{N}\right)^{\lfloor \frac{n}{2} \rfloor} \left\|f - \langle f \rangle_{\boldsymbol{\pi}}\right\|_{2,\boldsymbol{\pi}}$$

When N is large, this is a far better estimate than the one that comes from (b).

Exercise 6.6.13 Define \mathfrak{e} and \mathfrak{m} as in Theorem 6.4.11, and show that $\mathfrak{e} = -\mathfrak{m}$ if and only if for each $(i, j) \in \mathbb{S} \times \mathbb{S}_0$ there is an allowable path (i_0, \ldots, i_n) starting at i and ending at j along which H is non-increasing. That is, $i_0 = i$, $i_n = j$, and, for each $1 \le m \le n$, $\mathbf{A}_{i_{m-1}i_m} > 0$ and $H(i_m) \le H(i_{m-1})$.

Exercise 6.6.14 In this exercise I will give a very cursory introduction to a class of reversible Markov processes which provide somewhat naïve mathematical models of certain physical systems. In the literature, these are often called, for reasons which will be clear shortly, *spin-flip systems*, and they are among the earliest examples of Glauber dynamics. Here the state space $\mathbb{S} = \{-1, 1\}^N$ is to be thought of as the configuration space for a system of N particles, each of which has "spin" +1 or -1. Because it is more conventional, I will use $\boldsymbol{\omega} = (\omega_1, \ldots, \omega_N)$ or $\boldsymbol{\eta} = (\eta_1, \ldots, \eta_N)$ to denote generic elements of \mathbb{S} . Given $\boldsymbol{\omega} \in \mathbb{S}$ and $1 \le k \le N$, $\hat{\boldsymbol{\omega}}^k$ will be the configuration obtained from $\boldsymbol{\omega}$ by "flipping" its *k*th spin. That is, the $\hat{\boldsymbol{\omega}}^k = (\omega_1, \ldots, \omega_{k-1}, -\omega_k, \omega_{k+1}, \ldots, \omega_N)$. Next, given \mathbb{R} -valued functions *f* and *g* on \mathbb{S} , define

$$\Gamma(f,g)(\boldsymbol{\omega}) = \sum_{k=1}^{N} \left(f\left(\hat{\boldsymbol{\omega}}^{k}\right) - f(\boldsymbol{\omega}) \right) \left(g\left(\hat{\boldsymbol{\omega}}^{k}\right) - g(\boldsymbol{\omega}) \right),$$

which is a discrete analog of the dot product of the gradient of f with the gradient of g. Finally, given a probability vector $\boldsymbol{\mu}$ with $(\boldsymbol{\mu})_{\boldsymbol{\omega}} > 0$ for all $\boldsymbol{\omega} \in \mathbb{S}$, define

$$\mathcal{E}^{\mu}(f,g) = \frac{1}{2} \langle \Gamma(f,g) \rangle_{\mu}, \qquad \rho_{k}^{\mu}(\omega) = \frac{(\mu)_{\omega} + (\mu)_{\hat{\omega}^{k}}}{2(\mu)_{\omega}} \quad \text{and}$$
$$\left(\mathbf{Q}^{\mu}\right)_{\omega\eta} = \begin{cases} \rho_{k}^{\mu}(\omega) & \text{if } \eta = \hat{\omega}^{k} \\ 0 & \text{if } \eta \notin \{\omega, \hat{\omega}^{1}, \dots, \hat{\omega}^{N}\} \\ -\sum_{k=1}^{N} \rho_{k}^{\mu}(\omega) & \text{if } \eta = \omega. \end{cases}$$

(a) Check that \mathbf{Q}^{μ} is an irreducible *Q*-matrix on \mathbb{S} and that the detailed balance condition $(\mu)_{\omega}\mathbf{Q}^{\mu}_{\omega\eta} = (\mu)_{\eta}\mathbf{Q}^{\mu}_{\eta\omega}$ holds. In addition, show that

$$-\langle g, \mathbf{Q}^{\boldsymbol{\mu}} f \rangle_{\boldsymbol{\mu}} = \mathcal{E}^{\boldsymbol{\mu}}(f, g).$$

(b) Let λ be the uniform probability vector on S. That is, $(\lambda)_{\omega} = 2^{-N}$ for each $\omega \in S$. Show that

$$\operatorname{Var}_{\mu}(f) \leq M_{\mu} \operatorname{Var}_{\lambda}(f), \quad \text{where } M_{\mu} \equiv 2^{N} \max_{\omega \in \mathbb{S}} \mu_{\omega}$$

and

$$\mathcal{E}^{\lambda}(f, f) \leq \frac{1}{m_{\mu}} \mathcal{E}^{\mu}(f, f) \text{ where } m_{\mu} = 2^{N} \min_{\omega \in \mathbb{S}} \mu_{\omega}$$

- (c) For each $S \subseteq \{1, ..., N\}$, define $\chi_S : \mathbb{S} \longrightarrow \{-1, 1\}$ so that $\chi_S(\omega) = \prod_{k \in S} \omega_k$. In particular, $\chi_{\emptyset} = \mathbf{1}$. Show that $\{\chi_S : S \subseteq \{1, ..., N\}\}$ is an orthonormal basis in $L^2(\lambda)$ and that $\mathbf{Q}^{\lambda}\chi_S = -2(\#S)\chi_S$, and conclude from these that $\operatorname{Var}_{\lambda}(f) \leq \frac{1}{2}\mathcal{E}^{\lambda}(f, f)$.
- (d) By combining (b) with (c), show that $\beta_{\mu} \operatorname{Var}_{\mu} \leq \mathcal{E}^{\mu}(f, f)$ where $\beta_{\mu} \equiv \frac{2m_{\mu}}{M_{\mu}}$. In particular, if $t \rightsquigarrow \mathbf{P}_{t}^{\mu}$ is the transition probability function determined by \mathbf{Q}^{μ} , conclude that

$$\left\|\mathbf{P}_{t}^{\boldsymbol{\mu}}f-\langle f\rangle_{\boldsymbol{\mu}}\right\|_{2,\boldsymbol{\mu}}\leq e^{-t\beta_{\boldsymbol{\mu}}}\left\|f-\langle f\rangle_{\boldsymbol{\mu}}\right\|_{2,\boldsymbol{\mu}}.$$

Exercise 6.6.15 Refer to parts (b) and (c) in Exercise 6.6.14. It is somewhat surprising that the spectral gap for the uniform probability λ is 2, independent of *N*. In particular, this means that if { $\mathbf{P}^{(N)}(t) : t \ge 0$ } is the semigroup determined by the *Q*-matrix

$$(\mathbf{Q}^{(N)})_{\omega\eta} = \begin{cases} 1 & \text{if } \eta = \hat{\omega}^k \\ -N & \text{if } \eta = \omega \\ 0 & \text{otherwise} \end{cases}$$

on $\{-1, 1\}^N$, then $\|\mathbf{P}^{(N)}(t)f - \langle f \rangle_{\boldsymbol{\lambda}^{(N)}}\|_{2,\boldsymbol{\lambda}^{(N)}} \le e^{-2t} \|f\|_{2,\boldsymbol{\lambda}^{(N)}}$ for $t \ge 0$ and $f \in L^2(\boldsymbol{\lambda}^{(N)})$, where $\boldsymbol{\lambda}^{(N)}$ is the uniform probability measure on $\{-1, 1\}^N$. In fact, the situation here provides convincing evidence of that the theory developed in this chapter works in situations where Doeblin's theory is doomed to failure. Indeed, the purpose of this exercise is to prove that, for any t > 0 and $\boldsymbol{\omega} \in \{-1, 1\}^N$,

 $\lim_{N\to\infty} \|\boldsymbol{\mu}^{(N)}(t,\boldsymbol{\omega}) - \boldsymbol{\lambda}^{(N)}\|_{v} = 2$ when $(\boldsymbol{\mu}^{(N)}(t,\boldsymbol{\omega}))_{\eta} = (\mathbf{P}^{(N)}(t))_{\boldsymbol{\omega}\eta}$. That is, although the L^2 -estimate is independent of N, the variation estimate deteriorates as N gets larger.

- (a) Begin by showing that $\|\mu^{(N)}(t, \omega) \lambda^{(N)}\|_{v}$ is independent of ω .
- (b) Show that for any two probability vectors $\boldsymbol{\nu}$ and $\boldsymbol{\nu}'$ on a countable space \mathbb{S} , $2 \ge \|\boldsymbol{\nu} \boldsymbol{\nu}'\|_{v} \ge 2|\boldsymbol{\nu}(A) \boldsymbol{\nu}'(A)|$ for all $A \subseteq \mathbb{S}$.
- (c) For $1 \le k \le N$, define the random variable X_k on $\{-1, 1\}^N$ so that $X_k(\eta) = \eta_k$ if $\eta = (\eta_1, \ldots, \eta_N)$. Show that, under $\lambda^{(N)}$, the X_k 's are mutually independent, $\{-1, 1\}$ -valued Bernoulli random variables with expectation 0. Next, let t > 0 be given, and set $\mu^{(N)} = \mu^{(N)}(t, \omega)$, where ω is the element of $\{-1, 1\}$ whose coordinates are all -1. Show that, under $\mu^{(N)}$, the X_k 's are mutually independent, $\{-1, 1\}$ -valued Bernoulli random variables with expectation value $-e^{-2t}$.
- (d) Continuing in the setting of (c), let $A^{(N)}$ be the set of η for which $\frac{1}{N} \sum_{k=1}^{N} X_{k}(\eta) \leq -\frac{1}{2}e^{-2t}$, and show that

$$\boldsymbol{\mu}^{(N)}(A^{(N)}) - \boldsymbol{\lambda}^{(N)}(A^{(N)}) \ge 1 - \frac{8e^{4t}}{N}.$$

In fact, by using the sort of estimate developed at the end of Sect. 1.2.4, especially (1.2.16), one can sharpen this and get

$$\mu^{(N)}(A^{(N)}) - \lambda^{(N)}(A^{(N)}) \ge 1 - 2\exp\left(-\frac{Ne^{-4t}}{8}\right).$$

Hint: Use the usual Chebyshev estimate with which the weak law is proved.

(e) By combining the preceding, conclude that $\|\boldsymbol{\mu}^{(N)}(t,\boldsymbol{\omega}) - \boldsymbol{\lambda}^{(N)}\|_{v} \ge 2(1 - \frac{8e^{4t}}{N})$ for all $t > 0, N \in \mathbb{Z}^{+}$, and $\boldsymbol{\omega} \in \{-1, 1\}^{N}$.

Chapter 7 A Minimal Introduction to Measure Theory

On Easter 1933, A.N. Kolmogorov published *Foundations of Probability*, a book which laid the foundations on which most of probability theory has rested ever since. Because Kolmogorov's model is given in terms of Lebesgue's theory of measures and integration, its full appreciation requires a thorough understanding of that theory. Thus, although it is far too sketchy to provide anything approaching a thorough understanding, this chapter is an attempt to provide an introduction to Lebesgue's ideas and Kolmogorov's application of them in his model of probability theory.

7.1 A Description of Lebesgue's Measure Theory

In this section I will introduce the terminology used in Lebesgue's theory. However, I will systematically avoid giving rigorous proofs of any hard results. There are many places in which these proofs can be found, one of them being [8].

7.1.1 Measure Spaces

The essential components in measure theory are a set Ω , the *space*, a collection \mathcal{F} of subsets of Ω , the collection of *measurable*, *subsets*, and a function μ from \mathcal{F} into $[0, \infty]$, called the *measure*. Being a space on which a measure might exist, the pair (Ω, \mathcal{F}) is called a *measurable space*, and when a measurable space (Ω, \mathcal{F}) comes equipped with a measure μ , the triple $(\Omega, \mathcal{F}, \mu)$ is called a *measure space*.

In order to avoid stupid trivialities, we will always assume that the space Ω is non-empty. Also, we will assume that the collection \mathcal{F} of measurable sets forms a σ -algebra over Ω :

$$\Omega \in \mathcal{F}, \qquad A \in \mathcal{F} \implies A\mathbb{C} \equiv \Omega \setminus A \in \mathcal{F}, \quad \text{and}$$
$$\{A_n\}_1^\infty \subseteq \mathcal{F} \implies \bigcup_1^\infty A_n \in \mathcal{F}.$$

D.W. Stroock, *An Introduction to Markov Processes*, Graduate Texts in Mathematics 230, 179 DOI 10.1007/978-3-642-40523-5_7, © Springer-Verlag Berlin Heidelberg 2014

It is important to emphasize that, as distinguished from point-set topology, only finite or countable set theoretic operations are permitted in measure theory.

Using elementary set theoretic manipulations, it is easy to check that

$$A, B \in \mathcal{F} \implies A \cap B \in \mathcal{F} \text{ and } B \setminus A \in \mathcal{F}$$

 $\{A_n\}_1^\infty \subseteq \mathcal{F} \implies \bigcap_1^\infty A_n \in \mathcal{F}.$

Finally, the measure μ is a function which assigns¹ 0 to \emptyset and is *countably additive* in the sense that

$$\{A_n\}_1^\infty \subseteq \mathcal{F} \quad \text{and} \quad A_m \cap A_n = \emptyset \quad \text{when } m \neq n$$
$$\implies \quad \mu\left(\bigcup_{1}^\infty A_n\right) = \sum_{1}^\infty \mu(A_n). \tag{7.1.1}$$

In particular, for $A, B \in \mathcal{F}$,

$$A \subseteq B \implies \mu(B) = \mu(A) + \mu(B \setminus A) \ge \mu(A)$$

$$\mu(A \cap B) < \infty \implies \mu(A \cup B) = \mu(A) + \mu(B) - \mu(A \cap B).$$

(7.1.2)

The first line comes from writing *B* as the union of the disjoint sets *A* and $B \setminus A$, and the second line comes from writing $A \cup B$ as the union of the disjoint sets *A* and $B \setminus (A \cap B)$ and then applying the first line to *B* and $A \cap B$. The finiteness condition is needed when one moves the term $\mu(A \cap B)$ to the left hand side in $\mu(B) = \mu(A \cap B) + \mu(B \setminus (A \cap B))$. That is, one wants to avoid having to subtract ∞ from ∞ .

When the set Ω is finite or countable, one has no problem constructing measure spaces. In this case one can take $\mathcal{F} = \{A : A \subseteq \Omega\}$, the set of all subsets of Ω , make any assignment of $\omega \in \Omega \longmapsto \mu(\{\omega\}) \in [0, \infty]$, at which point countable additivity demands that we take

$$\mu(A) = \sum_{\omega \in A} \mu(\{\omega\}) \quad \text{for } A \subseteq \Omega,$$

where summation over the empty set is taken to be 0. However, when Ω is uncountable, it is far from obvious that interesting measures can be constructed on a non-trivial collection of measurable sets. Indeed, it is reasonable to think that Lebesgue's most significant achievement was his construction of a measure space in which $\Omega = \mathbb{R}$, \mathcal{F} is a σ -algebra of which every interval (open, closed, or semiclosed) is an element, and μ assigns each interval its length (i.e., $\mu(I) = b - a$ if I is an interval whose right and left end points are b and a).

¹In view of additivity, it is clear that either $\mu(\emptyset) = 0$ or $\mu(A) = \infty$ for all $A \in \mathcal{F}$. Indeed, by additivity, $\mu(\emptyset) = \mu(\emptyset \cup \emptyset) = 2\mu(\emptyset)$, and therefore $\mu(\emptyset)$ is either 0 or ∞ . Moreover, if $\mu(\emptyset) = \infty$, then $\mu(A) = \mu(A \cup \emptyset) = \mu(A) + \mu(\emptyset) = \infty$ for all $A \in \mathcal{F}$.

Although the terminology is misleading, a measure space $(\Omega, \mathcal{F}, \mu)$ is said to be *finite* if $\mu(\Omega) < \infty$. That is, the "finiteness" here is not determined by the size of Ω but instead by how large Ω looks to μ . Even if a measure space is not finite, it may be decomposable into a countable number of finite pieces, in which case it is said to be σ -*finite*. Equivalently, $(\Omega, \mathcal{F}, \mu)$ is σ -finite if there exists $\{\Omega_n\}_1^\infty \subseteq \mathcal{F}$ such that² $\Omega_n \nearrow \Omega$ and $\mu(\Omega_n) < \infty$ for each $n \ge 1$. Thus, for example, both Lebesgue's measure on \mathbb{R} and counting measure on \mathbb{Z} are σ -finite but not finite.

7.1.2 Some Consequences of Countable Additivity

Countable additivity is the *sine qua non* in this subject. In particular, it leads to the following continuity properties of measures:

$$\{A_n\}_1^{\infty} \subseteq \mathcal{F} \quad \text{and} \quad A_n \nearrow A \implies \mu(A_n) \nearrow \mu(A)$$

$$\{A_n\}_1^{\infty} \subseteq \mathcal{F}, \qquad \mu(A_1) < \infty, \quad \text{and} \quad A_n \searrow A \implies \mu(A_n) \searrow \mu(A).$$

(7.1.3)

Although I do not intend to prove many of the results discussed in this chapter, the proofs of those in (7.1.3) are too basic and easy to omit. To prove the first line, simply take $B_1 = A_1$ and $B_{n+1} = A_{n+1} \setminus A_n$. Then $B_m \cap B_n = \emptyset$ when $m \neq n$, $\bigcup_{i=1}^{n} B_m = A_n$ for all $n \ge 1$, and $\bigcup_{i=1}^{\infty} B_m = A$. Hence, by (7.1.1),

$$\mu(A_n) = \sum_{1}^{n} \mu(B_m) \nearrow \sum_{1}^{\infty} \mu(B_m) = \mu\left(\bigcup_{1}^{\infty} B_m\right) = \mu(A).$$

To prove the second line, begin by noting that $\mu(A_1) = \mu(A_n) + \mu(A_1 \setminus A_n)$ and $\mu(A_1) = \mu(A) + \mu(A_1 \setminus A)$. Hence, since $A_1 \setminus A_n \nearrow A_1 \setminus A$ and $\mu(A_1) < \infty$, we have $\mu(A_1) - \mu(A_n) \nearrow \mu(A_1) - \mu(A)$ and therefore that $\mu(A_n) \searrow \mu(A)$. Just as in the proof of second line in (7.1.2), we need the finiteness condition here to avoid being forced to subtract ∞ from ∞ .

Another important consequence of countable additivity is *countable subadditiv-ity*:

$$\{A_n\}_1^\infty \subseteq \mathcal{F} \implies \mu\left(\bigcup_{1}^\infty A_n\right) \le \sum_{1}^\infty \mu(A_n).$$
 (7.1.4)

Like the preceding, this is easy. Set $B_1 = A_1$ and $B_{n+1} = A_{n+1} \setminus \bigcup_{i=1}^{n} A_m$. Then $\mu(B_n) \le \mu(A_n)$ and

$$\mu\left(\bigcup_{1}^{\infty}A_{n}\right) = \mu\left(\bigcup_{1}^{\infty}B_{n}\right) = \sum_{1}^{\infty}\mu(B_{n}) \leq \sum_{1}^{\infty}\mu(A_{n}).$$

²I write $A_n \nearrow A$ when $A_n \subseteq A_{n+1}$ for all $n \ge 1$ and $A = \bigcup_{1=1}^{\infty} A_n$. Similarly, $A_n \searrow A$ means that $A_n \supseteq A_{n+1}$ for all $n \ge 1$ and $A = \bigcap_{1=1}^{\infty} A_n$. Obviously, $A_n \nearrow A$ if and only if $A_n \complement \searrow A \complement$.

A particularly important consequence of (7.1.4) is the fact that

$$\mu\left(\bigcup_{1}^{\infty} A_n\right) = 0 \quad \text{if } \mu(A_n) = 0 \text{ for each } n \ge 1.$$
 (7.1.5)

That is, the countable union of sets each of which has measure 0 is again a set having measure 0. Here one begins to see the reason for restricting oneself to countable operations in measure theory. Namely, it is certainly not true that the uncountable union of sets having measure 0 will necessarily have measure 0. For example, in the case of Lebesgue's measure on \mathbb{R} , the second line of (7.1.3) implies

$$\mu(\lbrace x\rbrace) = \lim_{\delta \searrow 0} \mu((x - \delta, x + \delta)) = \lim_{\delta \searrow 0} 2\delta = 0$$

for each point $x \in \mathbb{R}$, and yet $(0, 1) = \bigcup_{x \in (0, 1)} \{x\}$ has measure 1.

7.1.3 Generating σ -Algebras

Very often one wants to make sure that a certain collection of subsets will be among the measurable subsets, and for this reason it is important to know the following constructions. First, suppose that C is a collection of subsets of Ω . Then there is a smallest σ -algebra $\sigma(C)$, called the σ -algebra generated by C, over Ω which contains C. To see that $\sigma(C)$ always exists, consider the collection of all the σ -algebras over Ω which contain C. This collection is non-empty because $\{A : A \subseteq \Omega\}$ is an element. In addition, as is easily verified, the intersection of any collection of σ -algebras is again a σ -algebra. Hence, $\sigma(C)$ is the intersection of all the σ -algebras which contain C. When Ω is a topological space and C is the collection of all open subsets of Ω , then $\sigma(C)$ is called the *Borel* σ -algebra and is denoted by \mathcal{B}_{Ω} .

One of the most important reasons for knowing how a σ -algebra is generated is that one can often check properties of measures on $\sigma(\mathcal{C})$ by making sure that the property holds on \mathcal{C} . A particularly useful example of such a result is the one in the following uniqueness theorem.

Theorem 7.1.6 Suppose that (Ω, \mathcal{F}) is a measurable space and that $\mathcal{C} \subseteq \mathcal{F}$ includes Ω and is closed under intersection (i.e., $A \cap B \in \mathcal{C}$ whenever $A, B \in \mathcal{C}$). If μ and ν are a pair of finite measures on (Ω, \mathcal{F}) and if $\mu(A) = \nu(A)$ for each $A \in \mathcal{C}$, then $\mu(A) = \nu(A)$ for all $A \in \sigma(\mathcal{C})$.

Proof We will say that $S \subseteq \mathcal{F}$ is good if

(i) $A, B \in S$ and $A \subseteq B \implies B \setminus A \in S$.

- (ii) $A, B \in S$ and $A \cap B = \emptyset \implies A \cup B \in S$.
- (iii) $\{A_n\}_1^\infty \subseteq S$ and $A_n \nearrow A \implies A \in S$.

Notice that if S is good, $\Omega \in S$, and $A, B \in S \implies A \cap B \in S$, then S is a σ -algebra. Indeed, because of (i) and (iii), all that one has to do is check that $A \cup B \in S$ whenever $A, B \in S$. But, because $A \cup B = (A \setminus (A \cap B)) \cup B$, this is clear from (i), (ii), and the fact that S is closed under intersections. In addition, observe that if S is good, then, for any $D \subseteq F$,

$$\mathcal{S}' \equiv \{A \in \mathcal{S} : A \cap B \in \mathcal{S} \text{ for all } B \in \mathcal{D}\}$$

is again good.

Now set $\mathcal{B} = \{A \in \mathcal{F} : \mu(A) = \nu(A)\}$. From the properties of finite measures, in particular, (7.1.2) and (7.1.3), it is easy to check that \mathcal{B} is good. Moreover, by assumption, $\mathcal{C} \subseteq \mathcal{B}$. Thus, if

$$\mathcal{B}' = \{ A \in \mathcal{B} : A \cap C \in \mathcal{B} \text{ for all } C \in \mathcal{C} \},\$$

then, by the preceding observation, \mathcal{B}' is again good. In addition, because \mathcal{C} is closed under intersection, $\mathcal{C} \subseteq \mathcal{B}'$. Similarly,

$$\mathcal{B}'' \equiv \{ A \in \mathcal{B}' : A \cap B \in \mathcal{B}' \text{ for all } B \in \mathcal{B}' \}$$

is also good, and, by the definition of $\mathcal{B}', \mathcal{C} \subseteq \mathcal{B}''$. Finally, if $A, A' \in \mathcal{B}''$ and $B \in \mathcal{B}'$, then $(A \cap A') \cap B = A \cap (A' \cap B) \in \mathcal{B}'$, and so $A \cap A' \in \mathcal{B}''$. Hence, \mathcal{B}'' is a σ -algebra which contains $\mathcal{C}, \mathcal{B}'' \subseteq \mathcal{B}$, and therefore μ equals ν on $\sigma(\mathcal{C})$.

7.1.4 Measurable Functions

Given a pair $(\Omega_1, \mathcal{F}_1)$ and $(\Omega_2, \mathcal{F}_2)$ of measurable spaces, we will say that the map $F : \Omega_1 \longrightarrow \Omega_2$ is *measurable* if the inverse image of sets in \mathcal{F}_2 are elements of $\mathcal{F}_1: F^{-1}(\Gamma) \in \mathcal{F}_1$ for every $\Gamma \in \mathcal{F}_2$.³ Notice that if $\mathcal{F}_2 = \sigma(\mathcal{C})$, F is measurable if $F^{-1}(\mathcal{C}) \in \mathcal{F}_1$ for each $\mathcal{C} \in \mathcal{C}$. In particular, if Ω_1 and Ω_2 are topological spaces and $\mathcal{F}_i = \mathcal{B}_{\Omega_i}$, then every continuous map from Ω_1 to Ω_2 is measurable.

It is important to know that when $\Omega_2 = \mathbb{R}$ and $\mathcal{F}_2 = \mathcal{B}_{\mathbb{R}}$, measurability is preserved under sequential limit operations. To be precise, if $\{f_n\}_1^\infty$ is a sequence of \mathbb{R} -valued measurable functions from (Ω, \mathcal{F}) to $(\mathbb{R}, \mathcal{B}_{\mathbb{R}})$, then

$$\omega \rightsquigarrow \sup_{n} f_{n}(\omega), \qquad \omega \rightsquigarrow \inf_{n} f_{n}(\omega),$$
 are measurable. (7.1.7)
$$\omega \rightsquigarrow \overline{\lim_{n \to \infty}} f_{n}(\omega), \quad \text{and} \quad \omega \rightsquigarrow \underline{\lim_{n \to \infty}} f_{n}(\omega)$$

For example, the first of these can be proved by the following line of reasoning. Begin with the observation that $\mathcal{B}_{\mathbb{R}} = \sigma(\mathcal{C})$ when $\mathcal{C} = \{(a, \infty) : a \in \mathbb{R}\}$. Hence

³The reader should notice the striking similarity between this definition and the one for continuity in terms of inverse images of open sets.

 $f: \Omega \longrightarrow \mathbb{R}$ will be measurable if and only if⁴ {f > a} for each $a \in \mathbb{R}$, and so the first function in (7.1.7) is measurable because

$$\left\{\sup_{n} f_{n} > a\right\} = \bigcup_{n} \{f_{n} > a\}$$
 for every $a \in \mathbb{R}$.

As a consequence of the second line in (7.1.7), we know that the set Δ of points ω at which the limit $\lim_{n\to\infty} f_n(\omega)$ exists is measurable, and $\omega \rightsquigarrow \lim_{n\to\infty} f_n(\omega)$ is measurable if $\Delta = \Omega$.

Measurable functions give rise to important instances of the construction in the preceding subsection. Suppose that the space Ω_1 and the measurable space $(\Omega_2, \mathcal{F}_2)$ are given, and let \mathfrak{F} be some collection of maps from Ω_1 into Ω_2 . Then the σ -algebra $\sigma(\mathfrak{F})$ over Ω_1 generated by \mathfrak{F} is the smallest σ -algebra over Ω_1 with respect to which every element of \mathfrak{F} is measurable. Equivalently, $\sigma(\mathfrak{F})$ is equal to the set $\{F^{-1}(\Gamma) : F \in \mathfrak{F} \& \Gamma \in \mathcal{F}_2\}$. Finally, suppose that $\mathcal{F}_2 = \sigma(\mathcal{C}_2)$, where $\mathcal{C}_2 \subseteq \mathcal{F}_2$ contains Ω_2 and is closed under intersection (i.e., $A \cap B \in \mathcal{C}_2$ if $A, B \in \mathcal{C}_2$), and let \mathcal{C}_1 be the collection of sets of the form $F_1^{-1}(A_1) \cap \cdots \cap F_n^{-1}(A_n)$ for $n \in \mathbb{Z}^+$, $\{F_1, \ldots, F_n\} \subseteq \mathfrak{F}$, and $\{A_1, \ldots, A_n\} \subseteq \mathcal{C}_2$. Then $\sigma(\mathfrak{F}) = \sigma(\mathcal{C}_1)$, $\Omega_1 \in \mathcal{C}_1$, and \mathcal{C}_1 is closed under intersection. In particular, by Theorem 7.1.6, *if* μ and ν are a pair of finite measures on $(\Omega_1, \mathcal{F}_1)$ and

$$\mu\left(\left\{\omega_1: F_1(\omega_1) \in A_1, \dots, F_n(\omega_1) \in A_n\right\}\right)$$
$$= \nu\left(\left\{\omega_1: F_1(\omega_1) \in A_1, \dots, F_n(\omega_1) \in A_n\right\}\right)$$

for all $n \in \mathbb{Z}^+$, $\{F_1, \ldots, F_n\} \subseteq \mathfrak{F}$, and $\{A_1, \ldots, A_n\} \subseteq \mathcal{C}_2$, then μ equals ν on $\sigma(\mathfrak{F})$.

7.1.5 Lebesgue Integration

Given a measure space $(\Omega, \mathcal{F}, \mu)$, Lebesgue's theory of integration begins by giving a prescription for defining the integral with respect to μ of all non-negative, measurable functions (i.e., all measurable maps f from (Ω, \mathcal{F}) into⁵ ($[0, \infty]$; $\mathcal{B}_{[0,\infty]}$). His theory says that when $\mathbf{1}_A$ is the indicator function of a set $A \in \mathcal{F}$ and $a \in [0, \infty)$, then the integral of the function⁶ $a\mathbf{1}_A$ should be equal to $a\mu(A)$. He then insists that the integral should be additive in the sense that the integral of $f_1 + f_2$ should be sum of the integrals of f_1 and f_2 . In particular, this means that if f is a nonnegative, measurable function which is *simple*, in the sense that it takes on only a

⁴When it causes no ambiguity, I use $\{F \in \Gamma\}$ to stand for $\{\omega : F(\omega) \in \Gamma\}$.

⁵In this context, we are thinking of $[0, \infty]$ as the compact metric space obtained by mapping [0, 1] onto $[0, \infty]$ via the map $t \in [0, 1] \longmapsto \tan(\frac{\pi}{2}t)$.

⁶In measure theory, the convention which works best is to take $0 \infty = 0$.

A TI

finite number of values, then the integral $\int f d\mu$ of f must be

$$\sum_{x\in[0,\infty)}x\mu(f^{-1}(\{x\})),$$

where, because $f^{-1}({x}) = \emptyset$ for all but a finite number of *x*, the sum involves only a finite number of non-zero terms. Of course, before one can insist on this additivity of the integral, one is obliged to check that additivity is consistent. Specifically, it is necessary to show that $\sum_{1}^{n} a_m \mu(A_m) = \sum_{1}^{n'} a'_{m'} \mu(A'_{m'})$ when $\sum_{1}^{n} a_m \mathbf{1}_{A_m} = f =$ $\sum_{1}^{n'} a'_m \mathbf{1}_{A'_{m'}}$. However, once this consistency has been established, one knows how to define the integral of any non-negative, measurable, simple function in such a way that the integral is additive and gives the obvious answer for indicator functions. In particular, additivity implies *monotonicity*: $\int f dv \leq \int g d\mu$ when $f \leq g$.

To complete Lebesgue's program for non-negative functions, one has to first observe that if $f: E \longrightarrow [0, \infty]$ is measurable, then there exists a sequence $\{\varphi_n\}_1^\infty$ of non-negative, measurable, simple functions with the property that $\varphi_n(\omega) \nearrow f(\omega)$ for each $\omega \in \Omega$. For example, one can take

$$\varphi_n = \sum_{m=1}^{4^n} \frac{m}{2^n} \mathbf{1}_{A_{m,n}} \quad \text{where } A_{m,n} = \left\{ \omega : m2^{-n} \le f(\omega) < (m+1)2^{-n} \right\}$$

Given $\{\varphi_n\}_1^\infty$, what Lebesgue says is that $\int f d\mu = \lim_{n\to\infty} \int \varphi_n d\mu$. Indeed, by monotonicity, $\int \varphi_n d\mu$ is non-decreasing in *n*, and therefore the indicated limit necessarily exists. On the other hand, just as there was earlier, there is a consistency problem which one must resolve before adopting Lebesgue's definition. This time the problem comes from the fact that there are myriad ways in which to construct the approximating simple functions φ_n , and one must make sure that the limit does not depend on which approximation scheme one chooses. That is, it is necessary to check that if $\{\varphi_n\}_1^\infty$ and $\{\psi_n\}_1^\infty$ are two non-decreasing sequences of non-negative, simple, measurable functions such that $\lim_{n\to\infty} \varphi_n(\omega) = \lim_{n\to\infty} \psi_n(\omega)$ for each $\omega \in \Omega$, then $\lim_{n\to\infty} \int \varphi_n d\mu = \lim_{n\to\infty} \int \psi_n d\mu$; and it is at this step that the full power of countable additivity must be brought to bear.

Having defined $\int f d\mu$ for all non-negative, measurable f's, one must check that the resulting integral is homogeneous and additive: $\int af d\mu = a \int f d\mu$ for $a \in [0, \infty]$ and $\int (f_1 + f_2) d\mu = \int f_1 d\mu + \int f_2 d\mu$. However, both these properties are easily seen to be inherited from the case of simple f's. Thus, the only remaining challenge in Lebesgue's construction is to get away from the restriction to nonnegative functions and extend the integral to measurable functions which can take both signs. On the other hand, if one wants the resulting theory to be linear, then there is no doubt about how this extension must be made. Namely, given a measurable function $f : E \longrightarrow [-\infty, \infty]$, it is not hard to show that $f^+ \equiv f \lor 0$ and $f^- \equiv -(f \land 0) = (-f)^+$ are non-negative measurable functions. Hence, because $f = f^+ - f^-$, linearity demands that $\int f d\mu = \int f^+ d\mu - \int f^- d\mu$; and this time there are two problems which have to be confronted. In the first place, at the very least, it is necessary to restrict ones attention to functions f for which at most one of the numbers $\int f^+ d\mu$ or $\int f^- d\mu$ is infinite, otherwise one ends up having to deal with $\infty - \infty$. Secondly, one must, once again, check consistency. That is, if $f = f_1 - f_2$, where f_1 and f_2 are non-negative and measurable, then one has to show that $\int f_1 d\mu - \int f_2 d\mu = \int f^+ d\mu - \int f^- d\mu$.

In most applications, the measurable functions which one integrates are either non-negative or have the property that $\int |f| d\mu < \infty$, in which case f is said to be an *integrable function*. Because $|a_1f_1 + a_2f_2| \le |a_1||f_1| + |a_2||f_2|$, the set of integrable functions forms a vector space over \mathbb{R} on which integration with respect to μ acts as a linear function. Finally, if f is a measurable function which is either non-negative or integrable and $A \in \mathcal{F}$, then the product $\mathbf{1}_A f$ is again a measurable function which is either non-negative or integrable and so

$$\int_{A} f \, d\mu \equiv \int \mathbf{1}_{A} f \, d\mu \tag{7.1.8}$$

is well defined.

7.1.6 Stability Properties of Lebesgue Integration

The power of Lebesgue's theory derives from the stability of the integral it defines, and its stability is made manifest in the following three famous theorems. Throughout, $(\Omega, \mathcal{F}, \mu)$ is a measure space to which all references about measurability and integration refer. Also, the functions here are assumed to take their values in $(-\infty, \infty]$.

Theorem 7.1.9 (Monotone Convergence) Suppose that $\{f_n\}_1^\infty$ is a non-decreasing sequence of measurable functions, and, for each $\omega \in \Omega$, set $f(\omega) = \lim_{n \to \infty} f_n(\omega)$. If there exists a fixed integrable function g which is dominated by each f_n , then $\int f_n d\mu \nearrow \int f d\mu$. If, instead, $\{f_n\}_1^\infty$ is non-increasing and if there exists a fixed integrable function g which dominates each f_n , then $\int f_n d\mu \searrow \int f d\mu$.

Theorem 7.1.10 (Fatou's Lemmas) Given any sequence $\{f_n\}_1^\infty$ of measurable functions, all of which dominate some fixed integrable function g,

$$\underline{\lim}_{n\to\infty}\int f_n\,d\mu\geq\int\underline{\lim}_{n\to\infty}f_n\,d\mu.$$

If, instead, there is some fixed integrable function g which dominates all of the f_n 's, then

$$\overline{\lim_{n\to\infty}}\int f_n\,d\mu\leq\int\overline{\lim_{n\to\infty}}\,f_n\,d\mu.$$

Theorem 7.1.11 (Lebesgue's Dominated Convergence) Suppose that $\{f_n\}_1^\infty$ is a sequence of measurable functions, and assume that there is a fixed integrable function g with the property that $\mu(\{\omega : |f_n(\omega)| > g(\omega)\}) = 0$ for each $n \ge 1$. Further, assume that there exists a measurable function f to which $\{f_n\}_1^\infty$ converges in either one of the following two senses:

(a)
$$\mu\left(\left\{\omega:f(\omega)\neq\lim_{n\to\infty}f_n(\omega)\right\}\right)=0$$

or

(b)
$$\lim_{n \to \infty} \mu(\{\omega : |f_n(\omega) - f(\omega)| \ge \epsilon\}) = 0 \text{ for all } \epsilon > 0.$$

Then

$$\left|\int f_n \, d\mu - \int f \, d\mu\right| \leq \int |f_n - f| \, d\mu \longrightarrow 0 \quad \text{as } n \to \infty.$$

Before proceeding, a word should be said about the conditions in Lebesgue's Dominated Convergence Theorem. Because integrals cannot "see" what is happening on a set of measure 0 (i.e., an $A \in \mathcal{F}$ for which $\mu(A) = 0$) it is natural that conditions which guarantee certain behavior of integrals should be conditions which need only hold on the complement of a set of measure 0. Thus, condition (a) says that $f_n(\omega) \longrightarrow f(\omega)$ for all ω outside of a set of measure 0. In the standard jargon, conditions that hold off of a set of measure 0 are said to hold *almost everywhere*. In this terminology, the first hypothesis is that $|f_n| \le g$ almost everywhere and (a) is saying that $\{f_n\}_{1}^{\infty}$ converges to f almost everywhere, and these statements would be abbreviated in the literature by something like $|f_n| \leq g$ a.e. and $f_n \longrightarrow f$ a.e.. The condition in (b) is related to, but significantly different from, the one in (a). In particular, it does not guarantee that $\{f_n(\omega)\}_1^\infty$ converges for any $\omega \in \Omega$. For example, take μ to be the measure of Lebesgue on \mathbb{R} described above, and take $f_{m+2^n}(\omega) =$ $\mathbf{1}_{[0,2^{-n})}(\omega - m2^{-n})$ for $n \ge 0$ and $0 \le m < 2^n$. Then $\mu(\{\omega : f_{m+2^n}(\omega) \ne 0\}) = 2^{-n}$, and so $\lim_{n\to\infty} \mu(\{\omega : |f_n(\omega)| \ge \epsilon\}) = 0$ for all $\epsilon > 0$. On the other hand, for each $\omega \in [0, 1), \overline{\lim}_{n \to \infty} f_n(\omega) = 1$ but $\underline{\lim}_{n \to \infty} f_n(\omega) = 0$. Thus, (b) most definitely does not imply (a). Conversely, although (a) implies (b) when $\mu(\Omega) < \infty$, (a) does not imply (b) when $\mu(\Omega) = \infty$. To wit, again take μ to be Lebesgue's measure, and consider $f_n = \mathbf{1}_{\mathbb{R} \setminus [-n,n]}$.

In connection with the preceding discussion, there is a basic estimate, known as *Markov's inequality*, that plays a central role in all measure theoretic analysis. Namely, because, for any $\lambda > 0$, $\lambda \mathbf{1}_{[\lambda,\infty]} \circ f \leq f \mathbf{1}_{[\lambda,\infty]} \circ f \leq |f|$,

$$\mu\left(\left\{\omega:f(\omega)\geq\lambda\right\}\right)\leq\frac{1}{\lambda}\int_{\left\{\omega:f(\omega)\geq\lambda\right\}}f\,d\mu\leq\frac{1}{\lambda}\int|f|\,d\mu.$$
(7.1.12)

In particular, this leads to the conclusion that

$$\lim_{n \to \infty} \int |f_n - f| \, d\mu = 0 \quad \Longrightarrow \quad \mu \left(\left\{ \omega : \left| f_n(\omega) - f(\omega) \right| \ge \epsilon \right\} \right) = 0$$

for all $\epsilon > 0$. In particular, the condition (b) is necessary for the conclusion in Lebesgue's theorem. In addition, (7.1.12) proves that

$$\int |f| d\mu = 0 \quad \Longrightarrow \quad \mu(\{\omega : |f(\omega)| \ge \epsilon\}) = 0 \quad \text{for all } \epsilon > 0,$$

and therefore, by (7.1.3),

$$\int |f| d\mu = 0 \implies \mu(\{\omega : f(\omega) \neq 0\}) = 0.$$
 (7.1.13)

Finally, the role of the Lebesgue dominant g is made clear by considering $\{f_n\}_1^\infty$ when either $f_n = n \mathbf{1}_{[0,\frac{1}{2}]}$ or $f_n = \mathbf{1}_{[n-1,n)}$ and μ is Lebesgue's measure.

7.1.7 Lebesgue Integration on Countable Spaces

In this subsection we will see what Lebesgue's theory looks like in the relatively trivial case when Ω is countable, \mathcal{F} is the collection of all subsets of Ω , and the measure μ is σ -finite. As we pointed out in Sect. 7.1.1, specifying a measure on (Ω, \mathcal{F}) is tantamount to assigning a non-negative number to each element of Ω , and, because we want our measures to be σ -finite, no element is to be assigned ∞ .

Because the elements of Ω can be counted, there is no reason to not count them. Thus, in the case when Ω is finite, there is no loss in generality if we write $\Omega = \{1, ..., N\}$, where $N = \#\Omega$ is the number of elements in Ω , and, similarly, when Ω is countably infinite, we might as well, at least for abstract purposes, think of its being the set \mathbb{Z}^+ of positive integers. In fact, in order to avoid having to worry about the finite and countably infinite cases separately, we will embed the finite case into the infinite one by simply noting that the theory for $\{1, ..., N\}$ is exactly the same as the theory for \mathbb{Z}^+ restricted to measures μ for which $\mu(\{\omega\}) = 0$ when $\omega > N$. Finally, in order to make the notation here conform with the notation in the rest of the book, I will use \mathbb{S} in place of Ω , *i*, *j*, or *k* to denote generic elements of \mathbb{S} , and will identify a measure μ with the row vector $\boldsymbol{\mu} \in [0, \infty)^{\mathbb{S}}$ given by $(\boldsymbol{\mu})_i = \mu(\{i\})$.

The first thing to observe is that

$$\int f \, d\mu = \sum_{i \in \mathbb{S}} f(i)(\boldsymbol{\mu})_i \tag{7.1.14}$$

whenever either $\int f^+ d\mu$ or $\int f^- d\mu < \infty$ is finite. Indeed, if $\varphi \ge 0$ is simple, a_1, \ldots, a_L are the distinct values taken by f, and $A_\ell = \{i : \varphi(i) = a_\ell\}$, then

$$\int \varphi \, d\mu = \sum_{\ell=1}^{L} a_{\ell} \mu(A_{\ell}) = \sum_{\ell=1}^{L} a_{\ell} \sum_{i \in A_{\ell}} (\boldsymbol{\mu})_{i}$$
$$= \sum_{\ell=1}^{L} \sum_{i \in A_{\ell}} \varphi(i)(\boldsymbol{\mu})_{i} = \sum_{i \in \mathbb{S}} \varphi(i)(\boldsymbol{\mu})_{i}.$$

Second, given any $f \ge 0$ on \mathbb{S} , set $\varphi_n(i) = f(i)$ if $1 \le i \le n$ and $\varphi_n(i) = 0$ if i > n. Then,

$$\int f \, d\mu = \lim_{n \to \infty} \int \varphi_n \, d\mu = \lim_{n \to \infty} \sum_{1 \le i \le n} f(i)(\boldsymbol{\mu})_i = \sum_{i \in \mathbb{S}} f(i)(\boldsymbol{\mu})_i.$$

Finally, if either $\int f^+ d\mu$ or $\int f^- d\mu$ is finite, then it is clear that

$$\int f \, d\mu = \int f^+ \, d\mu - \int f^- \, d\mu$$
$$= \sum_{\{i:f(i) \ge 0\}} f(i)(\mu)_i - \sum_{\{i:f(i) \le 0\}} f(i)(\mu)_i = \sum_{i \in \mathbb{S}} f(i)(\mu)_i.$$

We next want to see what the "big three" look like in this context.

The Monotone Convergence Theorem First observe that it suffices to treat the case when $0 \le f_n \nearrow f$. Indeed, one can reduce each of these statements to that case by replacing f_n with $f_n - g$ or $g - f_n$. When $0 \le f_n \nearrow f$, it is obvious that

$$0 \leq \sum_{i \in \mathbb{S}} f_n(i)(\boldsymbol{\mu})_i \leq \sum_{i \in \mathbb{S}} f_{n+1}(i)(\boldsymbol{\mu})_i \leq \sum_{i \in \mathbb{S}} f(i)(\boldsymbol{\mu})_i.$$

Thus, all that remains is to note that

$$\lim_{n \to \infty} \sum_{i \in \mathbb{S}} f_n(i)(\boldsymbol{\mu})_i \ge \lim_{n \to \infty} \sum_{i=1}^L f_n(i)(\boldsymbol{\mu})_i = \sum_{i=1}^L f(i)(\boldsymbol{\mu})_i$$

for each $L \in \mathbb{Z}^+$, and therefore that the desired result follows after one lets $L \nearrow \infty$.

Fatou's Lemma Again one can reduce to the case when $f_n \ge 0$ and the limit being taken is the limit inferior. But in this case,

$$\lim_{n \to \infty} \sum_{i \in \mathbb{S}} f_n(i)(\boldsymbol{\mu})_i = \lim_{m \to \infty} \inf_{n \ge m} \sum_{i \in \mathbb{S}} f_n(i)(\boldsymbol{\mu})_i$$
$$\geq \lim_{m \to \infty} \sum_{i \in \mathbb{S}} \inf_{n \ge m} f_n(i)(\boldsymbol{\mu})_i = \sum_{i \in \mathbb{S}} \lim_{n \to \infty} f_n(i)(\boldsymbol{\mu})_i,$$

where the last equality follows from the monotone convergence theorem applied to $0 \leq \inf_{n \geq m} f_n \nearrow \underline{\lim}_{n \to \infty} f_n$.

Lebesgue's Dominated Convergence Theorem First note that, if we eliminate those $i \in S$ for which $(\mu)_i = 0$, none of the conclusions change. Thus, we will

assume that $(\boldsymbol{\mu})_i > 0$ for all $i \in \mathbb{S}$ from now on. Next observe that, under this assumption, the hypotheses become $\sup_n |f_n(i)| \le g(i)$ and $f_n(i) \longrightarrow f(i)$ for each $i \in \mathbb{S}$. In particular, $|f| \le g$. Further, by considering $f_n - f$ instead of $\{f_n\}_1^\infty$ and replacing g by 2g, we may and will assume that $f \equiv 0$. Now let $\epsilon > 0$ be given, and choose L so that $\sum_{i>L} g(i)(\boldsymbol{\mu})_i < \epsilon$. Then

$$\left|\sum_{i\in\mathbb{S}}f_n(i)(\boldsymbol{\mu})_i\right| \leq \sum_{i\in\mathbb{S}}\left|f_n(i)\right|(\boldsymbol{\mu})_i \leq \sum_{i=1}^L \left|f_n(i)\right|(\boldsymbol{\mu})_i + \epsilon \longrightarrow \epsilon \quad \text{as } n \to \infty.$$

7.1.8 Fubini's Theorem

Fubini's theorem deals with products of measure spaces. Given measurable spaces $(\Omega_1, \mathcal{F}_1)$ and $(\Omega_2, \mathcal{F}_2)$, the product of \mathcal{F}_1 and \mathcal{F}_2 is the σ -algebra $\mathcal{F}_1 \times \mathcal{F}_2$ over $\Omega_1 \times \Omega_2$ which is generated by the set $\{A_1 \times A_2 : A_1 \in \mathcal{F}_1 \& A_2 \in \mathcal{F}_2\}$ of *measurable rectangles*. An important technical fact about this construction is that, if f is a measurable function on $(\Omega_1 \times \Omega_2, \mathcal{F}_1 \times \mathcal{F}_2)$, then, for each $\omega_1 \in \Omega_1$ and $\omega_2 \in \Omega_2, \omega_2 \rightsquigarrow f(\omega_1, \omega_2)$ and $\omega_1 \rightsquigarrow f(\omega_1, \omega_2)$ are measurable functions on, respectively, $(\Omega_2, \mathcal{F}_2)$ and $(\Omega_1, \mathcal{F}_1)$.

In the following statement, it is important to emphasize exactly which variable is being integrated. For this reason, I will use the more detailed notation $\int_{\Omega} f(\omega) \mu(d\omega)$ instead of the more abbreviated $\int f d\mu$.

Theorem 7.1.15 (Fubin)⁷ Let $(\Omega_1, \mathcal{F}_1, \mu_1)$ and $(\Omega_2, \mathcal{F}_2, \mu_2)$ be a pair of σ finite measure spaces, and set $\Omega = \Omega_1 \times \Omega_2$ and $\mathcal{F} = \mathcal{F}_1 \times \mathcal{F}_2$. Then there is a unique measure $\mu = \mu_1 \times \mu_2$ on (Ω, \mathcal{F}) with the property that $\mu(A_1 \times A_2) =$ $\mu_1(A_1)\mu_2(A_2)$ for all $A_1 \in \mathcal{F}_1$ and $A_2 \in \mathcal{F}_2$. Moreover, if f is a non-negative, measurable function on (Ω, \mathcal{F}) , then both

$$\omega_1 \rightsquigarrow \int_{\Omega_2} f(\omega_1, \omega_2) \,\mu_2(d\omega_2) \quad and \quad \omega_2 \rightsquigarrow \int_{\Omega_1} f(\omega_1, \omega_2) \,\mu_1(d\omega_1)$$

are measurable functions, and

$$\int_{\Omega_1} \left(\int_{\Omega_2} f(\omega_1, \omega_2) \,\mu_2(d\omega_2) \right) \mu_1(d\omega_1)$$

=
$$\int_{\Omega} f \, d\mu = \int_{\Omega_2} \left(\int_{\Omega_1} f(\omega_1, \omega_2) \,\mu_1(d\omega_1) \right) \mu_2(d\omega_2).$$

⁷Although this theorem is usually attributed Fubini, it seems that Tonelli deserves, but seldom receives, a good deal of credit for it.

Finally, for any integrable function f on $(\Omega, \mathcal{F}, \mu)$,

$$\begin{split} A_1 &= \left\{ \omega_1 : \int_{\Omega_2} \left| f(\omega_1, \omega_2) \right| \mu_2(d\omega_2) < \infty \right\} \in \mathcal{F}_1, \\ A_2 &= \left\{ \omega_2 : \int_{\Omega_1} \left| f(\omega_1, \omega_2) \right| \mu_1(d\omega_1) < \infty \right\} \in \mathcal{F}_2, \\ \omega_1 \rightsquigarrow f_1(\omega_1) \equiv \mathbf{1}_{A_1}(\omega_1) \int_{\Omega_2} f(\omega_1, \omega_2) \mu_2(d\omega_2) \end{split}$$

and

$$\omega_2 \rightsquigarrow f_2(\omega_2) \equiv \mathbf{1}_{A_2}(\omega_2) \int_{\Omega_1} f(\omega_1, \omega_2) \,\mu_1(d\omega_1)$$

are integrable, and

$$\int_{\Omega_1} f_1(\omega_1) \,\mu_1(d\omega_1) = \int_{\Omega} f(\omega) \,\mu(d\omega) = \int_{\Omega_2} f_2(\omega_2) \,\mu_2(d\omega_2).$$

In the case when Ω_1 and Ω_2 are countable, all this becomes very easy. Namely, in the notation that we used in Sect. 7.1.6, the measure $\mu_1 \times \mu_2$ corresponds to row vector $\boldsymbol{\mu}_1 \times \boldsymbol{\mu}_2 \in [0, \infty)^{\mathbb{S}_1 \times \mathbb{S}_2}$ given by $(\boldsymbol{\mu}_1 \times \boldsymbol{\mu}_2)_{(i_1, i_2)} = (\boldsymbol{\mu})_{i_1}(\boldsymbol{\mu})_{i_2}$, and so, by (7.1.14), Fubini's Theorem reduces to the statement that

$$\sum_{i_1 \in \mathbb{S}_1} \left(\sum_{i_2 \in \mathbb{S}_2} a_{i_1 i_2} \right) = \sum_{(i_1, i_2) \in \mathbb{S}_1 \times \mathbb{S}_2} a_{i_1 i_2} = \sum_{i_2 \in \mathbb{S}_2} \left(\sum_{i_i \in \mathbb{S}_1} a_{i_1 i_2} \right)$$

when $\{a_{i_1i_2} : (i_1, i_2) \in \mathbb{S}_1 \times \mathbb{S}_2\} \subseteq (-\infty, \infty]$ satisfies either $a_{i_1i_2} \ge 0$ for all (i_1, i_2) or $\sum_{(i_1, i_2)} |a_{i_1i_2}| < \infty$. In proving this, we may and will assume that $\mathbb{S}_1 = \mathbb{Z}^+ = \mathbb{S}_2$ throughout and will start with the case when $a_{i_1i_2} \ge 0$. Given any pair $(n_1, n_2) \in (\mathbb{Z}^+)^2$,

$$\sum_{\substack{(i_1,i_2)\in\mathbb{S}_1\times\mathbb{S}_2\\i_1\le n_1 \& i_2\le n_2}} a_{i_1i_2} \ge \sum_{\substack{(i_1,i_2)\in\mathbb{S}_1\times\mathbb{S}_2\\i_1\le n_1 \& i_2\le n_2}} a_{i_1i_2} = \sum_{i_2=1}^{n_2} \left(\sum_{i_1=1}^{n_1} a_{i_1i_2}\right).$$

Hence, by first letting $n_1 \rightarrow \infty$ and then letting $n_2 \rightarrow \infty$, we arrive at

$$\sum_{(i_1,i_2)\in\mathbb{S}_1\times\mathbb{S}_2}a_{i_1i_2}\geq \sum_{i_2\in\mathbb{S}_2}\left(\sum_{i_1\in\mathbb{S}_1}a_{i_1i_2}\right).$$

Similarly, for any $n \in \mathbb{Z}^+$,

$$\sum_{\substack{(i_1,i_2)\mathbb{S}_1\times\mathbb{S}_2\\i_1\vee i_2\leq n}} a_{i_1i_2} = \sum_{i_2=1}^n \left(\sum_{i_1=1}^n a_{i_1i_2}\right) \leq \sum_{i_2\in\mathbb{S}_1} \left(\sum_{i_2\in\mathbb{S}_2} a_{i_1i_2}\right),$$

and so, after letting $n \to \infty$, we get the opposite inequality. Next, when $\sum_{(i_1,i_2)} |a_{i_1i_2}| < \infty$,

$$\sum_{i_1 \in \mathbb{S}_1} |a_{i_1 i_2}| < \infty \quad \text{for all } i_2 \in \mathbb{S}_2,$$

and so

$$\sum_{(i_1,i_2)\in\mathbb{S}_1\times\mathbb{S}_2} a_{i_1i_2} = \sum_{(i_1,i_2)\in\mathbb{S}_1\times\mathbb{S}_2} a_{i_1i_2}^+ - \sum_{(i_1,i_2)\in\mathbb{S}_1\times\mathbb{S}_2} a_{i_1i_2}^-$$
$$= \sum_{i_2\in\mathbb{S}_2} \left(\sum_{i_1\in\mathbb{S}_1} a_{i_1i_2}^+\right) - \sum_{i_2\in\mathbb{S}_2} \left(\sum_{i_1\in\mathbb{S}_1} a_{i_1i_2}^-\right)$$
$$= \lim_{n\to\infty} \sum_{\substack{i_2\in\mathbb{S}_2\\i_2\le n}} \left(\sum_{i_1\in\mathbb{S}_1} a_{i_1i_2}\right) = \sum_{i_2\in\mathbb{S}_2} \left(\sum_{i_1\in\mathbb{S}_1} a_{i_1i_2}^-\right).$$

Finally, after reversing the roles 1 and 2, we get the relation with the order of summation reversed.

7.2 Modeling Probability

To understand how these considerations relate to probability theory, think about the problem of modeling the tosses of a fair coin. When the game ends after the *n*th toss, a Kolmogorov model is provided by taking $\Omega = \{0, 1\}^n$, \mathcal{F} the set of all subsets of Ω , and setting $\mu(\{\omega\}) = 2^{-n}$ for each $\omega \in \Omega$. More generally, any measure space in which Ω has total measure 1 can be thought of as a model of probability, for which reason such a measure space is called a *probability space*, the measure μ is called a *probability measure* and it is often denoted by \mathbb{P} . In this connection, when dealing with probability spaces, ones intuition is aided by extending the metaphor to other objects. For example, one calls Ω the sample space, its elements are called sample points, the elements of \mathcal{F} are called *events*, and the number that \mathbb{P} assigns an event is called the *probability* of that event. In addition, a measurable map is called a random variable, it tends to be denoted by X instead of F, and, when it is \mathbb{R} -valued, its integral is called its expected value. Moreover, the latter convention is reflected in the use of $\mathbb{E}[X]$, or, when more precision is required, $\mathbb{E}^{\mathbb{P}}[X]$ to denote $\int X d\mathbb{P}$. Also, $\mathbb{E}[X, A]$ or $\mathbb{E}^{\mathbb{P}}[X, A]$ is used to denote $\int_{A} X d\mathbb{P}$, the expected value of X on the event A. Finally, the *distribution* of a random variable whose values lie in the measurable space (E, \mathcal{B}) is the probability measure $\mu = X_* \mathbb{P}$ on (E, \mathcal{B}) given by $\mu(B) = \mathbb{P}(X \in B)$. In particular, when X is \mathbb{R} -valued, its distribution function F_X is defined so that

$$F_X(x) = (X_* \mathbb{P}) \big((-\infty, x] \big) = \mathbb{P}(X \le x).$$

Obviously, F_X is non-increasing. Moreover, as an application of (7.1.3), one can show that F_X is continuous from the right, in the sense that $F_X(x) = \lim_{y \searrow x} F_X(y)$ for each $x \in \mathbb{R}$, $\lim_{x \searrow -\infty} F_X(x) = 0$, and $\lim_{x \nearrow \infty} F_X(x) = 1$. At the same time, one sees that $\mathbb{P}(X < x) = F_X(x-) \equiv \lim_{y \nearrow x} F_X(y)$, and, as a consequence, that $\mathbb{P}(X = x) = F_X(x) - F_X(x-)$ is the jump in F_X at x.

7.2.1 Modeling Infinitely Many Tosses of a Fair Coin

Just as in the general theory of measure spaces, the construction of probability spaces presents no analytic (as opposed to combinatorial) problems as long as the sample space is finite or countable. The only change from the general theory is that the assignment of probabilities to the sample points must satisfy the condition that $\sum_{\omega \in \Omega} \mathbb{P}(\{\omega\}) = 1$. The technical analytic problems arise when Ω is uncountable. For example, suppose that, instead of stopping the game after *n* tosses, one thinks about a coin tossing game of indefinite duration. Clearly, the sample space will now be $\Omega = \{0, 1\}^{\mathbb{Z}^+}$. In addition, when $A \subseteq \Omega$ depends only on tosses 1 through *n*, then *A* should be a measurable event and the probability assigned to *A* should be the same as the probability that would have been assigned had the game stopped after the *n*th toss. That is, if $\Gamma \subseteq \{0, 1\}^n$ and $A = \{\omega \in \Omega : (\omega(1), \ldots, \omega(n)) \in \Gamma\}$, then $\mathbb{P}(A)$ should equal $2^{-n} \# \Gamma$.

Continuing with the example of an infinite coin tossing game, one sees (cf. (7.1.3)) that, for any fixed $\eta \in \Omega$, $\mathbb{P}(\{\eta\})$ is equal to

$$\lim_{n\to\infty} \mathbb{P}\big(\big\{\omega: \big(\omega(1),\ldots,\omega(n)\big) = \big(\eta(1),\ldots,\eta(n)\big)\big\}\big) = \lim_{n\to\infty} 2^{-n} = 0.$$

Hence, in this case, nothing is learned from the way in which probability is assigned to points: every sample points has probability 0. In fact, it is far from obvious that there exists a probability measure on Ω with the all these properties. Nonetheless, as is proved in Sect. 2.2 of [8], one does. To be precise, what is proved there is the following theorem.

Theorem 7.2.1 Let $\Omega = \{0, 1\}^{\mathbb{Z}^+}$ and let \mathcal{A} be the set of all subsets A of Ω of the form

$$A = \left\{ \omega \in \Omega : \omega(m) = \epsilon_m \text{ for } 1 \le m \le n \right\}, \tag{7.2.2}$$

where $n \ge 1$ and $(\epsilon_1, \ldots, \epsilon_n) \in \{0, 1\}^n$. If $\mathcal{B} = \sigma(\mathcal{A})$, then there is a unique probability measure \mathbb{P} on (Ω, \mathcal{B}) such that $\mathbb{P}(\mathcal{A}) = 2^{-n}$ if $\mathcal{A} \in \mathcal{A}$ is the one in (7.2.2).

⁸It is convenient here to identify Ω with the set a mappings ω from \mathbb{Z}^+ into $\{0, 1\}$. Thus, we will use $\omega(n)$ to denote the "*n*th coordinate" of ω .

7.3 Independent Random Variables

In Kolmogorov's model, independence is best described in terms of σ -algebras. If $(\Omega, \mathcal{F}, \mathbb{P})$ is a probability space and \mathcal{F}_1 and \mathcal{F}_2 are σ -subalgebras (i.e., are σ algebras which are subsets) of \mathcal{F} , then we say that \mathcal{F}_1 and \mathcal{F}_2 are *independent* if $\mathbb{P}(\Gamma_1 \cap \Gamma_2) = \mathbb{P}(\Gamma_1)\mathbb{P}(\Gamma_2)$ for all $\Gamma_1 \in \mathcal{F}_1$ and $\Gamma_2 \in \mathcal{F}_2$. It should be comforting to recognize that, when $A_1, A_2 \in \mathcal{F}$ and, for $i \in \{1, 2\}, \mathcal{F}_i = \sigma(\{A_i\}) = \{\emptyset, A_i, A_i \complement, \Omega\}$, then, as is easily checked, \mathcal{F}_1 is independent of \mathcal{F}_2 precisely when, in the terminology of elementary probability theory, " A_1 is independent of A_2 ": $\mathbb{P}(A_1 \cap A_2) = \mathbb{P}(A_1)\mathbb{P}(A_2)$.

The notion of independence gets inherited by random variables. Namely, the members of a collection $\{X_{\alpha} : \alpha \in \mathcal{I}\}$ of random variables on $(\Omega, \mathcal{F}, \mathbb{P})$ are said to be *mutually independent* if, for each pair of disjoint subsets \mathcal{J}_1 and \mathcal{J}_2 of \mathcal{I} , the σ -algebras $\sigma(\{X_{\alpha} : \alpha \in \mathcal{J}_1\})$ and $\sigma(\{X_{\alpha} : \alpha \in \mathcal{J}_2\})$ are independent. One can use Theorem 7.1.6 to show that this definition is equivalent to saying that if X_{α} takes its values in the measurable space $(E_{\alpha}, \mathcal{B}_{\alpha})$, then, for every finite subset $\{\alpha_m\}_1^n$ of distinct elements of \mathcal{I} and choice of $B_{\alpha_m} \in \mathcal{B}_{\alpha_m}$, $1 \le m \le n$,

$$\mathbb{P}(X_{\alpha_m} \in B_{\alpha_m} \text{ for } 1 \le m \le n) = \prod_{1}^{n} \mathbb{P}(X_{\alpha_m} \in B_{\alpha_m}).$$

As a dividend of this definition, it is essentially obvious that if $\{X_{\alpha} : \alpha \in \mathcal{I}\}$ are mutually independent and if, for each $\alpha \in \mathcal{I}$, F_{α} is a measurable map on the range of X_{α} , then $\{F_{\alpha}(X_{\alpha}) : \alpha \in \mathcal{I}\}$ are again mutually independent. Finally, by starting with simple functions, one can show that if $\{X_m\}_1^n$ are mutually independent and, for each $1 \le m \le n$, f_m is a measurable \mathbb{R} -valued function on the range of X_m , then

$$\mathbb{E}\big[f_1(X_1)\cdots f_n(X_n)\big] = \prod_1^n \mathbb{E}\big[f_m(X_m)\big]$$

whenever the f_m 's are all bounded or are all non-negative.

7.3.1 Existence of Lots of Independent Random Variables

As a consequence of Theorem 7.2.1, we know that there exist an infinite family of mutually independent random variables. Namely, if $(\Omega, \mathcal{B}, \mathbb{P})$ is the probability space discussed in that theorem and $X_m(\omega) = \omega(m)$ is the *m*th coordinate of ω , then, for any choice of $n \ge 1$ and $(\epsilon_1, \ldots, \epsilon_n) \in \{0, 1\}^n$,

$$\mathbb{P}(X_m = \epsilon_m, \ 1 \le m \le n) = 2^{-n} = \prod_{1}^n \mathbb{P}(X_m = \epsilon_m).$$

Thus, the random variables $\{X_m\}_{1}^{\infty}$ are mutually independent. Mutually independent random variables, like these, which take on only two values are called *Bernoulli* random variables.

As we are about to see, Bernoulli random variables can be used as building blocks to construct many other families of mutually independent random variables. The key to such constructions is contained in the following lemma.

Lemma 7.3.1 Given any family $\{B_m\}_1^\infty$ of mutually independent, $\{0, 1\}$ -valued Bernoulli random variables satisfying $\mathbb{P}(B_m = 0) = \frac{1}{2} = \mathbb{P}(B_m = 1)$ for all $m \in \mathbb{Z}^+$, set $U = \sum_{1}^{\infty} 2^{-m} B_m$. Then U is uniformly distributed on [0, 1). That is, $\mathbb{P}(U \le u)$ is 0 if u < 0, u if $u \in [0, 1)$, and 1 if $u \ge 1$.

Proof Given $N \ge 1$ and $0 \le n < 2^N$, we want to show that

$$\mathbb{P}(n2^{-N} < U \le (n+1)2^{-N}) = 2^{-N}.$$
(*)

To this end, note that $n2^{-N} < U \le (n+1)2^{-N}$ if and only if $\sum_{1}^{N} 2^{-m} B_m = (n+1)2^{-N}$ and $B_m = 0$ for all m > N, or $\sum_{1}^{N} 2^{-m} B_m = n2^{-N}$ and $B_m = 1$ for some m > N. Hence, since $\mathbb{P}(B_m = 0$ for all m > N) = 0, the left hand side of (*) is equal to the probability that $\sum_{1}^{N} 2^{-m} B_m = n2^{-N}$. However, elementary considerations show that, for any $0 \le n < 2^N$ there is exactly one choice of $(\epsilon_1, \ldots, \epsilon_N) \in \{0, 1\}^N$ for which $\sum_{1}^{N} 2^{-m} \epsilon_m = n2^{-N}$. Hence,

$$\mathbb{P}\left(\sum_{1}^{N} 2^{-m} B_m = n 2^{-N}\right) = \mathbb{P}(B_m = \epsilon_m \text{ for } 1 \le m \le N) = 2^{-N}.$$

Having proved (*), the rest is easy. Since

$$\mathbb{P}(U=0) = \mathbb{P}(B_m = 0 \text{ for all } m \ge 1) = 0,$$

(*) tells us that, for any $1 \le k \le 2^N$

$$\mathbb{P}(U \le k2^{-N}) = \sum_{m=0}^{k-1} \mathbb{P}(m2^{-N} < U \le (m+1)2^{-N}) = k2^{-N}.$$

Hence, because, by (7.1.3), $u \rightsquigarrow \mathbb{P}(U \le u)$ is continuous from the right, it is now clear that $F_U(u) = u$ for all $u \in [0, 1)$. Finally, since $\mathbb{P}(U \in [0, 1]) = 1$, this completes the proof.

Now let \mathcal{I} a non-empty, finite or countably infinite set. Then $\mathcal{I} \times \mathbb{Z}^+$ is countable, and so we can construct a 1-to-1 map $(\alpha, n) \rightsquigarrow N(\alpha, n)$ from $\mathcal{I} \times \mathbb{Z}^+$ onto \mathbb{Z}^+ . Next, for each $\alpha \in \mathcal{I}$, define $\omega \in \Omega = \{0, 1\}^{\mathbb{Z}^+} \longmapsto X_{\alpha}(\omega) \in \Omega$ so that the *n*th coordinate of $X_{\alpha}(\omega)$ is $\omega(N(\alpha, n))$. Then, as random variables on the probability space $(\Omega, \mathcal{B}, \mathbb{P})$ in Theorem 7.2.1, $\{X_{\alpha} : \alpha \in \mathcal{I}\}$ are mutually independent and each has distribution \mathbb{P} . Hence, if $\Phi : \Omega \longrightarrow [0, 1]$ is the continuous map given by

$$\Phi(\eta) \equiv \sum_{m=1}^{\infty} 2^{-m} \eta(m) \quad \text{for } \eta \in \Omega$$

and if $U_{\alpha} \equiv \Phi(X_{\alpha})$, then the random variables $\{U_{\alpha} : \alpha \in \mathcal{I}\}$ are mutually independent and, by Lemma 7.3.1, each is uniformly distributed on [0, 1].

The final step in this program is to combine the preceding construction with the well-known fact that any \mathbb{R} -valued random variable can be represented in terms of a uniform random variable. More precisely, a map $F : \mathbb{R} \longrightarrow [0, 1]$ is called a *distribution function* if *F* is non-decreasing, continuous from the right, and tends to 0 at $-\infty$ and 1 at $+\infty$. Given such an *F*, define

$$F^{-1}(u) = \inf \{ x \in \mathbb{R} : F(x) \ge u \} \text{ for } u \in [0, 1].$$

Notice that, by right continuity, $F(x) \ge u \iff F^{-1}(u) \le x$. Hence, if U is uniformly distributed on [0, 1], then $F^{-1}(U)$ is a random variable whose distribution function is F.

Theorem 7.3.2 Let $(\Omega, \mathcal{B}, \mathbb{P})$ be the probability space in Theorem 7.2.1. Given any finite or countably infinite index set \mathcal{I} and a collection $\{F_{\alpha} : \alpha \in \mathcal{I}\}$ of distribution functions, there exist mutually independent random variables $\{X_{\alpha} : \alpha \in \mathcal{I}\}$ on $(\Omega, \mathcal{B}, \mathbb{P})$ with the property that, for each $\alpha \in \mathcal{I}$, F_{α} is the distribution of X_{α} .

7.4 Conditional Probabilities and Expectations

Just as they are to independence, σ -algebras are central to Kolmogorov's definition of conditioning. To be precise, given a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, a sub- σ -algebra Σ , and a random variable X which is non-negative or integrable, Kolmogorov says that the random variable X_{Σ} is a *conditional expectation* of X given Σ if X_{Σ} is a non-negative or integrable random variable which is measurable with respect to Σ (i.e., $\sigma(\{X_{\Sigma}\}) \subseteq \Sigma$) and satisfies

$$\mathbb{E}[X_{\Sigma}, \Gamma] = \mathbb{E}[X, \Gamma] \quad \text{for all } \Gamma \in \Sigma.$$
(7.4.1)

When *X* is the indicator function of a set $B \in \mathcal{F}$, then the term conditional expectation is replaced to *conditional probability*.

To understand that this definition is an extension of the one given in elementary probability courses, begin by considering the case when Σ is the trivial σ -algebra { \emptyset, Ω }. Because only constant random variables are measurable with respect to { \emptyset, Ω }, it is clear that the one and only conditional expectation of X will be $\mathbb{E}[X]$. Next, suppose that $\Sigma = \sigma(\{A\}) = \{\emptyset, A, AC, \Omega\}$ for some $A \in \mathcal{F}$ with $\mathbb{P}(A) \in (0, 1)$. In this situation, it is an easy matter to check that, for any $B \in \mathcal{F}$,

$$\omega \leadsto \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)} \mathbf{1}_{A}(\omega) + \frac{\mathbb{P}(B \cap A\mathbb{C})}{\mathbb{P}(A\mathbb{C})} \mathbf{1}_{A\mathbb{C}}(\omega)$$

is a conditional probability of *B* given Σ . That is, the quantity $\mathbb{P}(B|A) \equiv \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)}$, which in elementary probability theory would be called "the conditional probability of *B* given *A*," appears here as the value that Kolmogorov's conditional expectation value of $\mathbf{1}_B$ given $\sigma(\{A\})$ takes on *A*. More generally, if Σ is generated by a finite or countable partition $\mathcal{P} \subseteq \mathcal{F}$ of Ω , then, for any non-negative or integrable random variable *X*,

$$\omega \leadsto \sum_{\{A \in \mathcal{P}: \mathbb{P}(A) > 0\}} \frac{\mathbb{E}[X, A]}{\mathbb{P}(A)} \mathbf{1}_A(\omega)$$

will be a conditional expectation of X given Σ .

Of course, Kolmogorov's definition brings up two essential questions: existence and uniqueness. A proof of existence in general can be done in any one of many ways. For instance, when $\mathbb{E}[X^2] < \infty$, one can easily see that, just as $\mathbb{E}[X]\mathbf{1}$ is the X' at which $\mathbb{E}[(X - X')^2]$ achieves its minimum among constant random variables X', so X_{Σ} is the Σ -measurable X' at which $\mathbb{E}[(X - X')^2]$ achieves its minimum value. In this way, the problem of existence can be related to a problem of orthogonal projection in the space of all square-integrable random variables, and, although they are outside the scope of this book, such projection results are familiar to anyone who has studied the theory of Hilbert spaces.

Uniqueness, on the other hand, is both easier and more subtle than existence. Specifically, there is no naïve uniqueness statement here, because, in general, there will be uncountably many ways to take X_{Σ} .⁹ On the other hand, every choice differs from any other choice on a set of measure at most 0. To see this, suppose that X'_{Σ} is a second non-negative random variable which satisfies (7.4.1). Then $A = \{X'_{\Sigma} > X_{\Sigma}\} \in \Sigma$, and so the only way that (7.4.1) can hold is if $\mathbb{P}(A) = 0$. Similarly, $\mathbb{P}(X_{\Sigma} > X'_{\Sigma}) = 0$, and so $\mathbb{P}(X_{\Sigma} \neq X'_{\Sigma}) = 0$.

In spite of the ambiguity caused by the sort of uniqueness problems just discussed, it is common to ignore, in so far as possible, this ambiguity and proceed as if a random variable possesses only one conditional expectation with respect to a given σ -algebra. In this connection, the standard notation for a conditional expectation of X given Σ is $\mathbb{E}[X|\Sigma]$ or, when $X = \mathbf{1}_B$, $\mathbb{P}(B|\Sigma)$, which is the notation that I adopted in the earlier chapters of this book.

⁹This non-uniqueness is the reason for my use of the article "a" instead of "the" in front of "conditional expectation."

7.4.1 Conditioning with Respect to Random Variables

In this book, essentially all conditioning is done when $\Sigma = \sigma(\mathfrak{F})$ (cf. Sect. 7.1.4) for some family \mathfrak{F} of measurable functions on $(\Omega, \mathcal{F}, \mathbb{P})$. When Σ has this form, the conditional expectation of a random variable X will be a measurable functions of the functions in \mathfrak{F} . For example, if $\mathfrak{F} = \{F_1, \ldots, F_n\}$ and the functions F_m all take their values in a countable space \mathbb{S} , a conditional expectation value $\mathbb{E}[X | \sigma(\mathfrak{F})]$ of X given $\sigma(\mathfrak{F})$ has the form $\Phi(F_1, \ldots, F_n)$, where $\Phi(i_1, \ldots, i_n)$ is equal to

$$\frac{\mathbb{E}[X, F_1 = i_1, \dots, F_n = i_n]}{\mathbb{P}(F_1 = i_1, \dots, F_n = i_n)} \quad \text{or} \quad 0$$

according to whether $\mathbb{P}(F_1 = i_1, ..., F_n = i_n)$ is positive or 0. In order to emphasize that conditioning with respect to $\sigma(\mathfrak{F})$ results in a function of \mathfrak{F} , I use the notation $\mathbb{E}[X|\mathfrak{F}]$ or $\mathbb{P}(B|\mathfrak{F})$ instead of $\mathbb{E}[X|\sigma(\mathfrak{F})]$ or $\mathbb{P}(B|\sigma(\mathfrak{F}))$.

To give more concrete examples of what we are talking about, first suppose that *X* and *Y* are independent random variables with values in some countable space \mathbb{S} , and set Z = F(X, Y), where $F : \mathbb{S}^2 \longrightarrow \mathbb{R}$ is bounded. Then

$$\mathbb{E}[Z|X] = v(X) \quad \text{where } v(i) = \mathbb{E}[F(i,Y)] \text{ for } i \in \mathbb{S}.$$
(7.4.2)

A less trivial example is provided by our discussion of Markov chains. In Chap. 2 we encoded the Markov property in equations like

$$\mathbb{P}(X_{n+1}=j \mid X_0, \dots, X_n) = (\mathbf{P})_{X_n j},$$

which displays this conditioning as a function, namely $(i_0, \ldots, i_n) \rightsquigarrow (\mathbf{P})_{i_n j}$, of the random variables in terms of which the condition is made. (Of course, the distinguishing feature of the Markov property is that the function depends only on i_n and not (i_0, \ldots, i_{n-1}) .) Similarly, when we discussed Markov processes with a continuous time parameter, we wrote

$$\mathbb{P}(X(t) = j \mid X(\sigma), \sigma \in [0, s]) = (\mathbf{P}(t - s))_{X(s)j},$$

which again makes it explicit that the conditioned quantity is a function of the random variables on which the condition is imposed.

References

- Diaconis, P., Stroock, D.: Geometric bounds for eigenvalues of Markov chains. Ann. Appl. Probab. 1(1), 36–61 (1991)
- Dunford, N., Schwartz, J.: Linear Operators, Part I. Wiley Classics Lib. Wiley-Interscience, New York (1988)
- Holley, R., Stroock, D.: Simulated annealing via Sobolev inequalities. Commun. Math. Phys. 115(4), 553–569 (1988)
- Karlin, S., Taylor, H.: A First Course in Stochastic Processes, 2nd edn. Academic Press, San Diego (1975)
- Norris, J.R.: Markov Chains. Cambridge Series in Statistical & Probabilistic Mathematics. Cambridge University Press, Cambridge (1997)
- 6. Revuz, D.: Markov Chains. Mathematical Library, vol. 11. North-Holland, Amsterdam (1984)
- 7. Riesz, F., Sz.-Nagy, B.: Functional Analysis. Dover, New York (1990). Translated from the French edition by F. Boron, reprint of 1955 original
- 8. Stroock, D.: Mathematics of Probability. GSM, vol. 149. AMS, Providence (2013)

Index

A

accessible, 39, 49 adjoint, 43, 132, 138 allowable path, 149, 157 almost everywhere, 187 aperiodic, 56

B

backward equation, 113 variable, 113 Bernoulli random variable, 1, 195 binomial coefficient, 2 generalized, 6 branching process, 44 extinction, 45

С

Cauchy convergent, 28 Chapman-Kolmogorov equation, 104 time-inhomogeneous, 165 cofactor matrix. 59 communicates, 49 communicating class, 49 property, 50 conditional expectation, 196 conditional probability, 2, 196 continuity properties of measures, 181 convex function, 170 set. 79 convolution. 103 cooling schedule, 163 countably additive, 180

D

detailed balance condition, 138, 152 Dirichlet form, 147 distribution function, 192, 196 initial, 26 of a random variable, 192 Doeblin's condition, 32 Doeblin's theorem, 30 basic, 31 sharpened, 44 spectral theory, 32 Doob's stopping time theorem, 53 doubly stochastic, 44

Е

edge, 64 eigenvalue, 59 simple, 59 empirical measure, 97, 134, 135 ergodic theory, 36 ergodic theorem empirical measure, 97, 134 individual. 96, 134 mean, 36, 84, 128 event, 192 exhaustion. 114 expected value, 192 explosion does not occur, 119 time, 115 extreme point, 79

F

Fatou's lemma, 189 finite measure space, 181 first passage time, 3

D.W. Stroock, *An Introduction to Markov Processes*, Graduate Texts in Mathematics 230, 201 DOI 10.1007/978-3-642-40523-5, © Springer-Verlag Berlin Heidelberg 2014

forward equation, 113 variable, 113 Fubini's theorem, 190

G

generalized binomial coefficient, 6 Gibbs state, 157, 158 Glauber dynamics, 157 graph, 64 adjacency matrix of, 65 complete, 64 connected, 64 degree of a vertex in, 65 directed path in, 64 loop in, 64 spanning tree, 64 greatest common divisor, 56 Gronwall's inequality, 114, 132

H

Doob's *h*-transform, 47 homogeneous increments, 99

I

independence of σ -algebras, 194 of random variables, 194 individual ergodic theorem, 42, 96, 134 initial distribution, 26, 27 inner product, 139 integer part $\lfloor s \rfloor$ of *s*, 33 integrable function, 186 irreducible, 49

J

Jensen's inequality, 171

K

Kirchhoff's matrix tree theorem, 69 Kolmogorov equation backward, 113, 116 forward, 113 time-inhomogeneous, 164 Kronecker symbol δ_{ij} , 2

L

Lebesgue's dominated convergence theorem, 189 loop erasure, 65

М

Markov chain, 25 initial distribution, 26 Markov process generated by \mathbf{Q} , 111 starting from i, 119with initial distribution μ , 119 Markov property, 25, 30 Markov's inequality, 187 measurable map, 183 space, 179 subsets, 179 measure, 179 measure space, 179 σ -finite, 181 finite, 181 Metropolis algorithm, 162 monotone convergence theorem, 189 monotonicity of integral, 185 mutually independent random variables, 194

N

non-explosion, 120 non-negative definite, 147 norm, 28 null recurrent, 79

P

partition function, 158 period of a state, 56 Poincaré inequality, 147 Poincaré constant, 147 Poisson process compound with jump distribution μ , and rate *R*, 102 simple, 99 positive recurrent, 79, 130 probability, 192 measure, 192 space, 192 vector, 26 **P**-stationary, 78

Q

Q-positive recurrent, 128 Q-null recurrent, 128 Q-accessible, 123 Q-communicates, 123 Q-irreducible, 123 Q-matrix, 111 rates determined by, 111 transition matrix determined by, 111 Q-recurrent, 123 Q-transient, 123 Index

queuing model, 22, 92 queuing theory, 22

R

random time change, 135 random variable, 192 Bernoulli, 195 random walk symmetric, 7, 10 rates determined by \mathbf{Q} , 111 recurrence, 7 recurrent, 9, 10, 21, 37, 50 null. 79 positive, 79 reflection principle, 4 renewal equation, 77 time, 9 reverse, 43, 132 reversible, 137

S

sample point, 192 space, 192 Schwarz's inequality, 17, 139, 171 semigroup, 104 set of measure 0, 187 σ -finite measure space, 181 σ -algebra, 179 Borel, 182 generated by, 182, 184 smallest containing, 182 signum, 5 simple eigenvalue, 59 simple function, 184 simple Poisson process, 99 simulated annealing algorithm, 162 spectral gap, 142 spectral theory, 32, 142 spin-flip systems, 175 state space, 25 stationary distribution or probability, 30, 78, 94 null recurrent case, 95, 133 uniqueness, 95, 134 Stirling's formula, 19 strong law of large numbers, 18 subadditivity of measures, 181 subgraph, 64 symmetric on L^2 , 141

Т

time of first return, 7, 10, 50 time-inhomogeneous, 164 Chapman–Kolmogorov equation, 165 Kolmogorov's forward equation, 164 transition probability, 164 transition probability matrix, 26 determined by **Q**, 111 time inhomogeneous, 164 tree, 64 triangle inequality, 28

U

uniform norm $\|\cdot\|_u$, 29

V

variation norm $\|\cdot\|_v$, 27 vertex, 64

W

weak law of large numbers, 17 Wilson runs, 64