

Universitext

Universitext

Series editors

Sheldon Axler

San Francisco State University

Carles Casacuberta

Universitat de Barcelona

Angus MacIntyre

Queen Mary, University of London

Kenneth Ribet

University of California, Berkeley

Claude Sabbah

École polytechnique, CNRS, Université Paris-Saclay, Palaiseau

Endre Süli

University of Oxford

Wojbor A. Woyczyński

Case Western Reserve University

Universitext is a series of textbooks that presents material from a wide variety of mathematical disciplines at master's level and beyond. The books, often well class-tested by their author, may have an informal, personal, even experimental approach to their subject matter. Some of the most successful and established books in the series have evolved through several editions, always following the evolution of teaching curricula, into very polished texts.

Thus as research topics trickle down into graduate-level teaching, first textbooks written for new, cutting-edge courses may make their way into *Universitext*.

More information about this series at <http://www.springer.com/series/223>

Jean-François Collet

Discrete Stochastic Processes and Applications

 Springer

Jean-François Collet
Laboratoire J.A. Dieudonné
Université de Nice Sophia-Antipolis
Nice Cedex 02
France

ISSN 0172-5939 ISSN 2191-6675 (electronic)

Universitext

ISBN 978-3-319-74017-1 ISBN 978-3-319-74018-8 (eBook)

<https://doi.org/10.1007/978-3-319-74018-8>

Library of Congress Control Number: 2017964594

Mathematics Subject Classification (2010): Primary: 60J10, 60J27, 60J28, 60J75, 94A17;
Secondary: 26B25, 60J80, 60J85, 94A24

© Springer International Publishing AG, part of Springer Nature 2018

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.

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.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Springer imprint is published by the registered company Springer International Publishing AG
part of Springer Nature

The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Pour vous quatre

Preface

Stochastic processes are quite prevalent in scientific modeling, and the fine analysis of their mathematical properties relies on tools coming from such diverse branches of mathematics as (aside from probability theory of course) linear algebra, convex analysis, and information theory.

The aim of this book is to give a self-contained introduction to this extensive toolbox at a rather elementary level and to discuss its use in the study of Markov processes as well as in some other applications such as coding theory, population dynamics, and the design of search engines.

The first part of the book focuses on the rigorous theory of Markov chains and could provide a basis for a one-semester introductory course on this topic.

The main goal is to make it possible for the reader to develop a good intuition for the probabilistic concepts relevant to Markov processes without having to digest measure theory first. This is why we restrict ourselves to Markov processes on countable spaces, a.k.a. Markov chains. In this case, the mathematical form of the idea of absence of long-term memory is rather intuitive and does not require the extra technical tools (such as stopping times and filtrations) needed in the continuous case. The approach will be very gradual, starting with the discrete-time case, then moving on to continuous time. Along the way, new technicalities arise but the underlying probabilistic concepts are the same.

One of the most important results from the theory of Markov chains is the convergence theorem for ergodic chains (Theorem 1.39 in the text). This will be proved twice, first with probabilistic techniques in Chapter 1 (the coupling method is fundamentally probabilistic in nature) and second (for finite spaces) by pure linear algebra in Chapter 2. The fact that such a general result may be proved by attacking it from such radically different directions should cause readers to marvel once they have reached the end of Chapter 2, and it says a lot about the richness of the field and the beauty of mathematics. Chapter 2 ends with a popular application of the convergence result for Markov chains, the PageRank algorithm used in search engines.

The next step is to move on to continuous-time processes. The study of the Poisson process in Chapter 3 will serve as a pretext to introduce some of the relevant concepts without getting too technical. This will begin with a handmade construction of the Poisson process, before a review of the classical definitions and proof of their equivalence. Chapter 4 then deals with the rigorous theory of general continuous-time Markov processes. Since this text is intended as only a first exposition to stochastic processes, in order to avoid technicalities (which would mean introducing stopping times and filtrations) we do not cover the strong Markov property. Admittedly, the price to pay is that the rigorous proofs based on the use of skeletons may get a bit technical in places, so the reader who is interested mostly in applications might skip the detailed proofs.

Chapter 5 should provide some necessary relief after so much theory. It explores two kinds of examples: random walks and birth–death processes. The approach taken is much more applied than in the previous chapters and focuses on actual computations rather than theoretical results.

The first chapters provide two approaches to the study of the asymptotic behavior of Markov chains; a third route (and probably the least represented in textbooks on stochastic processes) is the use of entropy.

The second part of the text provides an introduction to various uses of convexity in probability as well as to the concept of entropy. It is more applied in nature (although mathematically rigorous) and might serve as a one-semester gentle introduction to coding and information theory.

At the heart of the uses of entropy is a certain number of inequalities that are in fact more or less all convexity inequalities. Assuming no prior knowledge of convexity whatsoever, Chapter 6 is an introduction to convex sets and convex maps, and to the relevant inequalities. In particular, it includes a discussion of Bregman divergences, which although they are now used in data clustering, are not often found in basic probability textbooks.

Chapter 7 collects the main quantities used in information theory; according to the expositional philosophy of the first part we present everything in the framework of discrete distributions, but the extension to absolutely continuous distributions is trivial.

The main theme here is that many quantities (such as Bregman divergences and Φ -divergences and in particular the Kullback–Leibler divergence) defined on the space of probability vectors may be used as distances even though they are not distances in the true mathematical sense of the term. What is meant here by “may be used as distances” becomes clear once the main inequalities, Jensen’s and Pinsker’s, are understood. In Section 7.4 we show how to use entropy to give a third proof of the convergence theorem for ergodic chains (at which point the above comment on the beauty and unity of mathematics should be reiterated). This is a use in the setting of Markov chains of what is known as the “entropy dissipation method” in partial differential equations and dynamical systems. The chapter ends with a detailed

introduction to exponential families, which are parametric families of probability distributions that under some constraints maximize entropy. Finally, Chapter 8 explores the connection with binary coding; essentially, it is shown that the entropy of a random source provides an incompressible lower bound for the length of any uniquely decipherable code. To make a long story short, loss of memory in a Markov chain or complexity of a random source may be quantified using the same tool, entropy.

The prerequisites for this text are as follows.

- From linear algebra: multiplication of matrices, scalar product, and the concept of eigenvalue (no knowledge of diagonalization results is needed); the section on the Perron–Frobenius theorem is totally self-contained.
- From calculus: sequences and proofs by induction, manipulation of convergent series, the chain rule, computation of one-dimensional integrals by change of variables and integration by parts. From calculus in several variables, the notions of gradient and Hessian.
- From probability theory: discrete and absolutely continuous random variables; conditional probability, computations of expectations, moments, Gaussian variables.

Acknowledgments. Now comes the most pleasant and easiest paragraph of this book to write. It is indeed a pleasure to thank all the people without whom this book would not have been written: Didier Auroux, Florent Berthelin, Jaques Blum, Cédric Boulbe, David Chiron, Jean-Antoine Désidéri, Frédéric Gruy, Bernard Guy, David Hoff, Stéphane Junca, Frédéric Poupaud, Francesca Rapetti, Vitaly Volpert.

Touët-sur-Var
Jean-François Collet July 2017

Contents

Part I Markov processes

1	Discrete time, countable space	3
1.1	Conditional probability on a discrete space	3
1.2	Formal definition of a Markov chain on a countable space	4
1.3	Homogeneous chains and transition matrices	5
1.4	The Chapman–Kolmogorov equation	6
1.5	Transient and recurrent states	8
1.6	Hitting times	15
1.7	Closed sets and state space decomposition	16
1.8	Asymptotic behavior	17
1.8.1	Irreducibility	17
1.8.2	Positive recurrence	20
1.8.3	Periodicity and ergodicity	22
1.9	Finite state space	26
1.10	Problems	29
2	Linear algebra and search engines	33
2.1	Spectra of Markov matrices	33
2.2	More on the location of eigenvalues for general square matrices	37
2.3	Positive matrices and Perron’s theorem	39
2.4	Irreducibility and the theorem of Frobenius	43
2.5	Primitivity and the power method	49
2.6	The PageRank algorithm	52
2.6.1	Formulation as an eigenvalue problem	52
2.6.2	Formulation as a linear system	54
2.7	Problems	54

3	The Poisson process	57
3.1	Some heuristics on jumps and consequences	57
3.2	The Poisson process	60
3.2.1	A handmade construction	60
3.2.2	Formal definition	60
3.2.3	Jump times	63
3.2.4	The bus paradox	68
3.3	Problems	71
4	Continuous time, discrete space	75
4.1	Continuous-time Markov chains and their discretizations	75
4.1.1	The Markov property and the transition function	75
4.1.2	Discretizing a continuous-time Markov chain	77
4.1.3	The distribution of the first jump time	79
4.2	The semigroup approach	80
4.2.1	Markov semigroups	80
4.2.2	Continuity and differentiability of the semigroup	82
4.2.3	The rate matrix and Kolmogorov's equation	86
4.3	Communication classes	93
4.3.1	The embedded chain	93
4.3.2	Communication and recurrence	94
4.4	Stationary distributions and convergence to equilibrium	96
4.5	Defining a continuous-time Markov chain by its generator ...	98
4.6	Problems	100
5	Examples	103
5.1	Random walks	103
5.1.1	Random walks on \mathbb{Z}	103
5.1.2	Symmetric random walks on \mathbb{Z}^d	105
5.1.3	The Reflecting walk on \mathbb{N}	108
5.1.4	Extinction of a discrete population	110
5.2	Birth–death processes	110
5.2.1	Stationary distribution	111
5.2.2	Solution by the Laplace transform method	112
5.3	Problems	115

Part II Entropy and applications

6	Prelude: a user's guide to convexity	121
6.1	Why are convex sets and convex maps so important in probability?	121
6.2	The Legendre transform	127
6.3	Bregman divergences	129
6.4	Problems	132

7	The basic quantities of information theory	139
7.1	Entropy discrete random variables	139
7.1.1	Entropy of one variable	139
7.1.2	Entropy of a pair of random variables, conditional entropy	140
7.2	Using entropy to count things: Shearer’s lemma	142
7.3	Φ -divergences	144
7.3.1	General form	144
7.3.2	The most famous Φ -Divergence: Kullback–Leibler divergence	145
7.4	Entropy dissipation in Markov processes	149
7.5	Exponential families	151
7.6	Problems	156
8	Application: binary coding	161
8.1	The basic vocabulary of binary coding	161
8.2	Instantaneous codes and binary trees	163
8.3	The Kraft–McMillan inequality	164
8.3.1	The instantaneous case	164
8.3.2	The uniquely decipherable case	166
8.4	Code length and entropy	167
8.5	Problems	169
A	Some useful facts from calculus	171
A.1	Croft’s lemma	171
A.2	A characterization of exponential functions and distributions	173
A.3	Countable sums	175
A.4	Right continuous and right constant functions	177
A.5	The Gamma function	178
A.6	The Laplace transform	178
B	Some useful facts from probability	181
B.1	Some limit theorems in probability theory	181
B.1.1	Continuity of probability from above and below	181
B.1.2	Three notions of convergence of sequences of random variables	183
B.2	Exponential and related distributions	183
B.2.1	Some properties of the exponential distribution	183
B.2.2	The Gamma distribution	184
B.2.3	The truncated exponential distribution	185
B.2.4	Binomial coefficients, binomial and related distributions	186
B.2.5	The Vandermonde convolution identity	186
B.2.6	Obtaining the Poisson distribution from the Binomial distribution	186

Contents

B.2.7	The negative binomial distribution	187
B.3	Order statistics.....	188
C	Some useful facts from linear algebra	189
C.1	Matrix norms	189
C.2	Eigenvalues and spectral radius	190
C.3	Monotone matrices and M matrices.....	192
C.4	Permutation matrices	194
C.5	Matrix exponentials	195
D	An arithmetic lemma	201
E	Table of exponential families	203
	References	213
	Index	217

Terminology, Abbreviations, and Typographical Conventions

Terminology and Abbreviations

Throughout the text, the words “positive” and “negative” are understood strictly; the components of a vector q are interchangeably denoted by q_i or $q(i)$, and similarly, the components of a matrix p are interchangeably denoted by $p_{i,j}$ or $p(i,j)$. We refer to the probability mass function of a discrete random variable X simply as its probability vector q , and then write $q(x)$ for $P(X = x)$, even in the case that X takes its values in an infinite countable set.

Unless otherwise stated, vectors appearing in products are understood as row vectors; therefore, in a product of a vector by a matrix, the vector will appear to the left of the matrix.

The end of a rigorous proof is indicated by the symbol \square . The notation $:=$ indicates that an equality should be understood as a definition of its left-hand-side (as opposed to equalities resulting from a computation).

Notation

- A^t : the transpose matrix of a matrix A , meaning $A^t(i,j) := A(j,i)$. This notation includes vectors, so for example, if x is a column vector, then x^t is the corresponding row vector.
- I_n : the $n \times n$ identity matrix; sometimes denoted by I when n is implicit.
- $\delta_{x,y}$ is the well-known Kronecker delta : it is 1 if $x = y$, and 0 otherwise.
- \mathbb{N}^* , \mathbb{R}^* : positive integers, positive real numbers.
- $[x,y]$: the segment joining two points $x, t \in \mathbb{R}^n$, which means the set of all points of the form $\lambda x + (1 - \lambda)y$ for some $0 \leq \lambda \leq 1$. In particular, if $n = 1$, this is just the familiar *interval* $[x,y]$. If we wish to exclude (for instance) y , we write $[x,y)$.

Terminology, Abbreviations, and Typographical Conventions

- $\nabla f(x)$: the gradient of a scalar-valued map (in this text we never have to use the derivatives of a vector-valued map) $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\nabla f(x)(i) := \frac{\partial f}{\partial x_i}(x).$$

- $\nabla^2 f(x)$: the Hessian matrix of a map $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\nabla^2 f(x)(i, j) := \frac{\partial^2 f}{\partial x_i \partial x_j}(x).$$

- $H(X), H(p)$: entropy of a random variable, of a probability vector
- $\Gamma(\alpha), B(\alpha, \beta)$: the Euler gamma and beta functions:

$$\Gamma(\alpha) := \int_0^\infty e^{-t} t^{\alpha-1} dt, \quad B(\alpha, \beta) := \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt.$$

- $I_0(z)$: the modified Bessel function of order zero :

$$I_0(z) := \frac{1}{\pi} \int_0^\pi \exp(x \cos t) dt.$$

- $D(f), R(f)$: the domain and range of a map f . The notation $f : X \rightarrow Y$ means that the function f is defined on some subset of X , the domain $D(f)$. This means that what we call a function is what some people call a *partial function*. The raison d'être for this choice is that we use X just as a way to specify the nature of the variable (e.g., scalar or vector), without having to make assumptions about $D(f)$.
- \mathcal{P}_n : the set of n -component probability vectors:

$$\mathcal{P}_n := \{p \in \mathbb{R}^n : p_i \geq 0, p_1 + \cdots + p_n = 1\}.$$

- \mathcal{P}_n^* : the set of n -component probability vectors with positive components:

$$\mathcal{P}_n^* := \{p \in \mathbb{R}^n : p_i > 0, p_1 + \cdots + p_n = 1\}.$$

- $\log_2 x$: the logarithm to the base two; if a computation or statement is valid independently of the base, then the notation \log is used. In this case, we assume only that the base is strictly greater than 1 (so that the corresponding logarithm function is increasing).
- $E_p(\phi)$: the expectation of ϕ under the probability distribution p . In other words, $E_p(\phi) := E(\phi(X))$, where X is a random variable having probability distribution p .
- \mathcal{A}^+ : the set of all finite strings over a finite set \mathcal{A} .
- c^+ : the extension to \mathcal{A}^+ of a binary code c defined on \mathcal{A} .
- $l_c(x)$: the code length of the symbol x , that is, the length of the codeword $c(x)$.

Terminology, Abbreviations, and Typographical Conventions

- $L(c)$: the average code length of the code c .
- $d_H(x, y)$: the Hamming distance between two binary words x and y of the same length.
- $\stackrel{d}{=}$: equality in distribution.
- $\xrightarrow{d}, \xrightarrow{p}, \xrightarrow{a.s.}$: convergence in distribution, in probability, almost sure.
- \sim : distributed as:
 1. $X \sim \mathcal{P}(\lambda)$ means that X is a Poisson variable of parameter λ ;
 2. $X \sim \mathcal{B}(n, p)$ means that X is a binomial variable of parameters n and p ;
 3. $X \sim \mathcal{NB}(k, p)$ means that X is a negative binomial variable of parameters k and p ;
 4. $X \sim \mathcal{G}(p)$ means that X is a geometric variable of parameter p ;
 5. $X \sim \mathcal{E}(\lambda)$ means that X is an exponential variable of parameter λ ;
 6. $X \sim \mathcal{TE}(\lambda)$ means that X is a truncated exponential variable of parameters λ and n ;
 7. $X \sim \Gamma(k, \lambda)$ means that X is a Γ variable of parameters λ and k .