

Lecture Notes in Mathematics

1837

Editors:

J.-M. Morel, Cachan

F. Takens, Groningen

B. Teissier, Paris

Springer

Berlin

Heidelberg

New York

Hong Kong

London

Milan

Paris

Tokyo

Simon Tavaré Ofer Zeitouni

Lectures on Probability Theory and Statistics

Ecole d'Été de Probabilités
de Saint-Flour XXXI - 2001

Editor: Jean Picard



Springer

Authors

Simon Tavaré
Program in Molecular and
Computational Biology
Department of Biological Sciences
University of Southern California
Los Angeles, CA 90089-1340
USA

e-mail: stavare@usc.edu

Ofer Zeitouni
Departments of Electrical Engineering
and of Mathematics
Technion - Israel Institute of Technology
Haifa 32000, Israel
and
Department of Mathematics
University of Minnesota
206 Church St. SE
Minneapolis, MN 55455
USA

e-mail: zeitouni@ee.technion.ac.il
zeitouni@math.umn.edu

Editor

Jean Picard
Laboratoire de Mathématiques Appliquées
UMR CNRS 6620
Université Blaise Pascal Clermont-Ferrand
63177 Aubière Cedex, France

e-mail: Jean.Picard@math.univ-bpclermont.fr

Cover illustration: Blaise Pascal (1623-1662)

Cataloging-in-Publication Data applied for

Bibliographic information published by Die Deutsche Bibliothek

Die Deutsche Bibliothek lists this publication in the Deutsche Nationalbibliografie;
detailed bibliographic data is available in the Internet at <http://dnb.ddb.de>

Mathematics Subject Classification (2001):
60-01, 60-06, 62-01, 62-06, 92D10, 60K37, 60F05, 60F10

ISSN 0075-8434 Lecture Notes in Mathematics
ISSN 0721-5363 Ecole d'Été des Probabilités de St. Flour
ISBN 3-540-20832-1 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

Springer-Verlag Berlin Heidelberg New York a member of BertelsmannSpringer
Science + Business Media GmbH

<http://www.springer.de>

© Springer-Verlag Berlin Heidelberg 2004
Printed in Germany

The use of general descriptive names, registered names, trademarks, 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.

Typesetting: Camera-ready \TeX output by the authors

SPIN: 10981573 41/3142/du - 543210 - Printed on acid-free paper

Preface

Three series of lectures were given at the 31st Probability Summer School in Saint-Flour (July 8–25, 2001), by the Professors Catoni, Tavaré and Zeitouni. In order to keep the size of the volume not too large, we have decided to split the publication of these courses into two parts. This volume contains the courses of Professors Tavaré and Zeitouni. The course of Professor Catoni entitled “Statistical Learning Theory and Stochastic Optimization” will be published in the *Lecture Notes in Statistics*. We thank all the authors warmly for their important contribution.

55 participants have attended this school. 22 of them have given a short lecture. The lists of participants and of short lectures are enclosed at the end of the volume.

Finally, we give the numbers of volumes of Springer *Lecture Notes* where previous schools were published.

Lecture Notes in Mathematics

1971: vol 307	1973: vol 390	1974: vol 480	1975: vol 539
1976: vol 598	1977: vol 678	1978: vol 774	1979: vol 876
1980: vol 929	1981: vol 976	1982: vol 1097	1983: vol 1117
1984: vol 1180	1985/86/87: vol 1362	1988: vol 1427	1989: vol 1464
1990: vol 1527	1991: vol 1541	1992: vol 1581	1993: vol 1608
1994: vol 1648	1995: vol 1690	1996: vol 1665	1997: vol 1717
1998: vol 1738	1999: vol 1781	2000: vol 1816	

Lecture Notes in Statistics

1986: vol 50	2003: vol 179
--------------	---------------

Contents

Part I Simon Tavaré: Ancestral Inference in Population Genetics

Contents	3
1 Introduction	6
2 The Wright-Fisher model	9
3 The Ewens Sampling Formula	30
4 The Coalescent	44
5 The Infinitely-many-sites Model	54
6 Estimation in the Infinitely-many-sites Model	79
7 Ancestral Inference in the Infinitely-many-sites Model	94
8 The Age of a Unique Event Polymorphism	111
9 Markov Chain Monte Carlo Methods	120
10 Recombination	151
11 ABC: Approximate Bayesian Computation	169
12 Afterwords	179
References	180

Part II Ofer Zeitouni: Random Walks in Random Environment

Contents	191
1 Introduction	193
2 RWRE – $d=1$	195
3 RWRE – $d > 1$	258
References	308

List of Participants 313

List of Short Lectures 315

**Simon Tavaré: Ancestral Inference in
Population Genetics**

Ancestral Inference in Population Genetics

Simon Tavaré

Departments of Biological Sciences, Mathematics and Preventive Medicine
University of Southern California.

1	Introduction	6
1.1	Genealogical processes	6
1.2	Organization of the notes	7
1.3	Acknowledgements	8
2	The Wright-Fisher model	9
2.1	Random drift	9
2.2	The genealogy of the Wright-Fisher model	12
2.3	Properties of the ancestral process	19
2.4	Variable population size	23
3	The Ewens Sampling Formula	30
3.1	The effects of mutation	30
3.2	Estimating the mutation rate	32
3.3	Allozyme frequency data	33
3.4	Simulating an infinitely-many alleles sample	34
3.5	A recursion for the ESF	35
3.6	The number of alleles in a sample	37
3.7	Estimating θ	38
3.8	Testing for selective neutrality	41
4	The Coalescent	44
4.1	Who is related to whom?	44
4.2	Genealogical trees	47
4.3	Robustness in the coalescent	47
4.4	Generalizations	52
4.5	Coalescent reviews	53
5	The Infinitely-many-sites Model	54
5.1	Measures of diversity in a sample	56

5.2	Pairwise difference curves	59
5.3	The number of segregating sites	59
5.4	The infinitely-many-sites model and the coalescent	64
5.5	The tree structure of the infinitely-many-sites model	65
5.6	Rooted genealogical trees	67
5.7	Rooted genealogical tree probabilities	68
5.8	Unrooted genealogical trees	71
5.9	Unrooted genealogical tree probabilities	73
5.10	A numerical example	74
5.11	Maximum likelihood estimation	77
6	Estimation in the Infinitely-many-sites Model	79
6.1	Computing likelihoods	79
6.2	Simulating likelihood surfaces	81
6.3	Combining likelihoods	82
6.4	Unrooted tree probabilities	83
6.5	Methods for variable population size models	84
6.6	More on simulating mutation models	86
6.7	Importance sampling	87
6.8	Choosing the weights	90
7	Ancestral Inference in the Infinitely-many-sites Model	94
7.1	Samples of size two	94
7.2	No variability observed in the sample	95
7.3	The rejection method	96
7.4	Conditioning on the number of segregating sites	97
7.5	An importance sampling method	101
7.6	Modeling uncertainty in N and μ	101
7.7	Varying mutation rates	104
7.8	The time to the MRCA of a population given data from a sample	105
7.9	Using the full data	108
8	The Age of a Unique Event Polymorphism	111
8.1	UEP trees	111
8.2	The distribution of T_Δ	114
8.3	The case $\mu = 0$	116
8.4	Simulating the age of an allele	118
8.5	Using intra-allelic variability	118
9	Markov Chain Monte Carlo Methods	120
9.1	K -Allele models	121
9.2	A biomolecular sequence model	124
9.3	A recursion for sampling probabilities	125
9.4	Computing probabilities on trees	126
9.5	The MCMC approach	127

9.6	Some alternative updating methods	132
9.7	Variable population size	137
9.8	A Nuu Chah Nulth data set	138
9.9	The age of a UEP	142
9.10	A Yakima data set	145
10	Recombination	151
10.1	The two locus model	151
10.2	The correlation between tree lengths	157
10.3	The continuous recombination model	160
10.4	Mutation in the ARG	163
10.5	Simulating samples	165
10.6	Linkage disequilibrium and haplotype sharing	167
11	ABC: Approximate Bayesian Computation	169
11.1	Rejection methods	169
11.2	Inference in the fossil record	170
11.3	Using summary statistics	175
11.4	MCMC methods	176
11.5	The genealogy of a branching process	177
12	Afterwords	179
12.1	The effects of selection	179
12.2	The combinatorics connection	179
12.3	Bugs and features	180
	References	180

1 Introduction

One of the most important challenges facing modern biology is how to make sense of genetic variation. Understanding how genotypic variation translates into phenotypic variation, and how it is structured in populations, is fundamental to our understanding of evolution. Understanding the genetic basis of variation in phenotypes such as disease susceptibility is of great importance to human geneticists. Technological advances in molecular biology are making it possible to survey variation in natural populations on an enormous scale. The most dramatic examples to date are provided by Perlegen Sciences Inc., who resequenced 20 copies of chromosome 21 (Patil *et al.*, 2001) and by Genaissance Pharmaceuticals Inc., who studied haplotype variation and linkage disequilibrium across 313 human genes (Stephens *et al.*, 2001). These are but two of the large number of variation surveys now underway in a number of organisms. The amount of data these studies will generate is staggering, and the development of methods for their analysis and interpretation has become central. In these notes I describe the basics of *coalescent theory*, a useful quantitative tool in this endeavor.

1.1 Genealogical processes

These Saint Flour lectures concern *genealogical processes*, the stochastic models that describe the ancestral relationships among samples of individuals. These individuals might be species, humans or cells – similar methods serve to analyze and understand data on very disparate time scales. The main theme is an account of methods of statistical inference for such processes, based primarily on stochastic computation methods. The notes do not claim to be even-handed or comprehensive; rather, they provide a personal view of some of the theoretical and computational methods that have arisen over the last 20 years. A comprehensive treatment is impossible in a field that is evolving as fast as this one. Nonetheless I think the notes serve as a useful starting point for accessing the extensive literature.

Understanding molecular variation data

The first lecture in the Saint Flour Summer School series reviewed some basic molecular biology and outlined some of the problems faced by computational molecular biologists. This served to place the problems discussed in the remaining lectures into a broader perspective. I have found the books of Hartl and Jones (2001) and Brown (1999) particularly useful.

It is convenient to classify evolutionary problems according to the time scale involved. On long time scales, think about trying to reconstruct the molecular phylogeny of a collection of species using DNA sequence data taken

from a homologous region in each species. Not only is the phylogeny, or branching order, of the species of interest but so too might be estimation of the divergence time between pairs of species, of aspects of the mutation process that gave rise to the observed differences in the sequences, and questions about the nature of the common ancestor of the species. A typical population genetics problem involves the use of patterns of variation observed in a sample of humans to locate disease susceptibility genes. In this example, the time scale is of the order of thousands of years. Another example comes from cancer genetics. In trying to understand the evolution of tumors we might extract a sample of cells, type them for microsatellite variation at a number of loci and then use the observed variability to infer the time since a checkpoint in the tumor's history. The time scale in this example is measured in years.

The common feature that links these examples is the dependence in the data generated by common ancestral history. Understanding the way in which ancestry produces dependence in the sample is the key principle of these notes. Typically the ancestry is never known over the whole time scale involved. To make any progress, the ancestry has to be modelled as a stochastic process. Such processes are the subject of these notes.

Backwards or Forwards?

The theory of population genetics developed in the early years of the last century focused on a *prospective* treatment of genetic variation (see Provine (2001) for example). Given a stochastic or deterministic model for the evolution of gene frequencies that allows for the effects of mutation, random drift, selection, recombination, population subdivision and so on, one can ask questions like ‘How long does a new mutant survive in the population?’, or ‘What is the chance that an allele becomes fixed in the population?’. These questions involve the analysis of the future behavior of a system given initial data. Most of this theory is much easier to think about if the focus is *retrospective*. Rather than ask where the population will go, ask where it has been. This changes the focus to the study of ancestral processes of various sorts. While it might be a truism that genetics is all about ancestral history, this fact has not pervaded the population genetics literature until relatively recently. We shall see that this approach makes most of the underlying methodology easier to derive – essentially all classical prospective results can be derived more simply by this dual approach – and in addition provides methods for analyzing modern genetic data.

1.2 Organization of the notes

The notes begin with forwards and backwards descriptions of the Wright-Fisher model of gene frequency fluctuation in Section 2. The ancestral process that records the number of distinct ancestors of a sample back in time is described, and a number of its basic properties derived. Section 3 introduces

the effects of mutation in the history of a sample, introduces the genealogical approach to simulating samples of genes. The main result is a derivation of the Ewens sampling formula and a discussion of its statistical implications. Section 4 introduces Kingman's coalescent process, and discusses the robustness of this process for different models of reproduction.

Methods more suited to the analysis of DNA sequence data begin in Section 5 with a theoretical discussion of the infinitely-many-sites mutation model. Methods for finding probabilities of the underlying reduced genealogical trees are given. Section 6 describes a computational approach based on importance sampling that can be used for maximum likelihood estimation of population parameters such as mutation rates. Section 7 introduces a number of problems concerning inference about properties of coalescent trees conditional on observed data. The motivating example concerns inference about the time to the most recent common ancestor of a sample. Section 8 develops some theoretical and computational methods for studying the ages of mutations. Section 9 discusses Markov chain Monte Carlo approaches for Bayesian inference based on sequence data. Section 10 introduces Hudson's coalescent process that models the effects of recombination. This section includes a discussion of ancestral recombination graphs and their use in understanding linkage disequilibrium and haplotype sharing.

Section 11 discusses some alternative approaches to inference using approximate Bayesian computation. The examples include two at opposite ends of the evolutionary time scale: inference about the divergence time of primates and inference about the age of a tumor. This section includes a brief introduction to computational methods of inference for samples from a branching process. Section 12 concludes the notes with pointers to some topics discussed in the Saint Flour lectures, but not included in the printed version. This includes models with selection, and the connection between the stochastic structure of certain decomposable combinatorial models and the Ewens sampling formula.

1.3 Acknowledgements

Paul Marjoram, John Molitor, Duncan Thomas, Vincent Plagnol, Darryl Shibata and Oliver Will were involved with aspects of the unpublished research described in Section 11. I thank Lada Markovtsova for permission to use some of the figures from her thesis (Markovtsova (2000)) in Section 9. I thank Magnus Nordborg for numerous discussions about the mysteries of recombination. Above all I thank Warren Ewens and Bob Griffiths, collaborators for over 20 years. Their influence on the statistical development of population genetics has been immense; it is clearly visible in these notes.

Finally I thank Jean Picard for the invitation to speak at the summer school, and the Saint-Flour participants for their comments on the earlier version of the notes.

2 The Wright-Fisher model

This section introduces the Wright-Fisher model for the evolution of gene frequencies in a finite population. It begins with a prospective treatment of a population in which each individual is one of two types, and the effects of mutation, selection, . . . are ignored. A genealogical (or retrospective) description follows. A number of properties of the ancestral relationships among a sample of individuals are given, along with a genealogical description in the case of variable population size.

2.1 Random drift

The simplest Wright-Fisher model (Fisher (1922), Wright (1931)) describes the evolution of a two-allele locus in a population of constant size undergoing random mating, ignoring the effects of mutation or selection. This is the so-called ‘random drift’ model of population genetics, in which the fundamental source of “randomness” is the reproductive mechanism.

A Markov chain model

We assume that the population is of constant size N in each non-overlapping generation n , $n = 0, 1, 2, \dots$. At the locus in question there are two alleles, denoted by A and B . X_n counts the number of A alleles in generation n . We assume first that there is no mutation between the types. The population at generation $r + 1$ is derived from the population at time r by binomial sampling of N genes from a gene pool in which the fraction of A alleles is its current frequency, namely $\pi_i = i/N$. Hence given $X_r = i$, the probability that $X_{r+1} = j$ is

$$p_{ij} = \binom{N}{j} \pi_i^j (1 - \pi_i)^{N-j}, \quad 0 \leq i, j \leq N. \quad (2.1.1)$$

The process $\{X_r, r = 0, 1, \dots\}$ is a time-homogeneous Markov chain. It has transition matrix $P = (p_{ij})$, and state space $\mathcal{S} = \{0, 1, \dots, N\}$. The states 0 and N are absorbing; if the population contains only one allele in some generation, then it remains so in every subsequent generation. In this case, we say that the population is *fixed* for that allele.

The binomial nature of the transition matrix makes some properties of the process easy to calculate. For example,

$$\mathbb{E}(X_r | X_{r-1}) = N \frac{X_{r-1}}{N} = X_{r-1},$$

so that by averaging over the distribution of X_{r-1} we get $\mathbb{E}(X_r) = \mathbb{E}(X_{r-1})$, and

$$\mathbb{E}(X_r) = \mathbb{E}(X_0), \quad r = 1, 2, \dots \quad (2.1.2)$$

The result in (2.1.2) can be thought of as the analog of the Hardy-Weinberg law: in an infinitely large random mating population, the relative frequency of the alleles remains constant in every generation. Be warned though that average values in a stochastic process do not tell the whole story! While on average the number of A alleles remains constant, variability must eventually be lost. That is, eventually the population contains all A alleles or all B alleles.

We can calculate the probability a_i that eventually the population contains only A alleles, given that $X_0 = i$. The standard way to find such a probability is to derive a system of equations satisfied by the a_i . To do this, we condition on the value of X_1 . Clearly, $a_0 = 0$, $a_N = 1$, and for $1 \leq i \leq N - 1$, we have

$$a_i = p_{i0} \cdot 0 + p_{iN} \cdot 1 + \sum_{j=1}^{N-1} p_{ij} a_j. \quad (2.1.3)$$

This equation is derived by noting that if $X_1 = j \in \{1, 2, \dots, N - 1\}$, then the probability of reaching N before 0 is a_j . The equation in (2.1.3) can be solved by recalling that $\mathbb{E}(X_1 | X_0 = i) = i$, or

$$\sum_{j=0}^N p_{ij} j = i.$$

It follows that $a_i = Ci$ for some constant C . Since $a_N = 1$, we have $C = 1/N$, and so $a_i = i/N$. Thus the probability that an allele will fix in the population is just its initial frequency.

The variance of X_r can also be calculated from the fact that

$$\text{Var}(X_r) = \mathbb{E}(\text{Var}(X_r | X_{r-1})) + \text{Var}(\mathbb{E}(X_r | X_{r-1})).$$

After some algebra, this leads to

$$\text{Var}(X_r) = \mathbb{E}(X_0)(N - \mathbb{E}(X_0))(1 - \lambda^r) + \lambda^r \text{Var}(X_0), \quad (2.1.4)$$

where

$$\lambda = 1 - 1/N.$$

We have noted that genetic variability in the population is eventually lost. It is of some interest to assess how fast this loss occurs. A simple calculation shows that

$$\mathbb{E}(X_r(N - X_r)) = \lambda^r \mathbb{E}(X_0(N - X_0)). \quad (2.1.5)$$

Multiplying both sides by $2N^{-2}$ shows that the probability $h(r)$ that two genes chosen at random with replacement in generation r are different is

$$h(r) = \lambda^r h(0). \quad (2.1.6)$$

The quantity $h(r)$ is called the *heterozygosity* of the population in generation r , and it measures the genetic variability surviving in the population. Equation

(2.1.6) shows that the heterozygosity decays geometrically quickly as $r \rightarrow \infty$. Since fixation must occur, we have $h(r) \rightarrow 0$.

We have seen that variability is lost from the population. How long does this take? First we find an equation satisfied by m_i , the mean time to fixation starting from $X_0 = i$. To do this, notice first that $m_0 = m_N = 0$, and, by conditioning on the first step once more, we see that for $1 \leq i \leq N - 1$

$$\begin{aligned} m_i &= p_{i0} \cdot 1 + p_{iN} \cdot 1 + \sum_{j=1}^{N-1} p_{ij}(1 + m_j) \\ &= 1 + \sum_{j=0}^N p_{ij}m_j. \end{aligned} \tag{2.1.7}$$

Finding an explicit expression for m_i is difficult, and we resort instead to an approximation when N is large and time is measured in units of N generations.

Diffusion approximations

This takes us into the world of diffusion theory. It is usual to consider not the total number $X_r \equiv X(r)$ of A alleles but rather the proportion X_r/N . To get a non-degenerate limit we must also rescale time, in units of N generations. This leads us to study the rescaled process

$$Y_N(t) = N^{-1}X(\lfloor Nt \rfloor), \quad t \geq 0, \tag{2.1.8}$$

where $\lfloor x \rfloor$ is the integer part of x . The idea is that as $N \rightarrow \infty$, $Y_N(\cdot)$ should converge in distribution to a process $Y(\cdot)$. The fraction $Y(t)$ of A alleles at time t evolves like a continuous-time, continuous state-space process in the interval $\mathcal{S} = [0, 1]$. $Y(\cdot)$ is an example of a diffusion process. Time scalings in units proportional to N generations are typical for population genetics models appearing in these notes.

Diffusion theory is the basic tool of classical population genetics, and there are several good references. Crow and Kimura (1970) has a lot of the ‘old style’ references to the theory. Ewens (1979) and Kingman (1980) introduce the sampling theory ideas. Diffusions are also discussed by Karlin and Taylor (1980) and Ethier and Kurtz (1986), the latter in the measure-valued setting. A useful modern reference is Neuhauser (2001).

The properties of a one-dimensional diffusion $Y(\cdot)$ are essentially determined by the infinitesimal mean and variance, defined in the time-homogeneous case by

$$\begin{aligned} \mu(y) &= \lim_{h \rightarrow 0} h^{-1} \mathbb{E}(Y(t+h) - Y(t) \mid Y(t) = y), \\ \sigma^2(y) &= \lim_{h \rightarrow 0} h^{-1} \mathbb{E}((Y(t+h) - Y(t))^2 \mid Y(t) = y). \end{aligned}$$

For the discrete Wright-Fisher model, we know that given $X_r = i$, X_{r+1} is binomially distributed with number of trials N and success probability i/N . Hence

$$\begin{aligned}\mathbb{E}(X(r+1)/N - X(r)/N \mid X(r)/N = i/N) &= 0, \\ \mathbb{E}((X(r+1)/N - X(r)/N)^2 \mid X(r)/N = i/N) &= \frac{1}{N} \frac{i}{N} \left(1 - \frac{i}{N}\right),\end{aligned}$$

so that for the process $Y(\cdot)$ that gives the proportion of allele A in the population at time t , we have

$$\mu(y) = 0, \quad \sigma^2(y) = y(1-y), \quad 0 < y < 1. \quad (2.1.9)$$

Classical diffusion theory shows that the mean time $m(x)$ to fixation, starting from an initial fraction $x \in (0, 1)$ of the A allele, satisfies the differential equation

$$\frac{1}{2}x(1-x)m''(x) = -1, \quad m(0) = m(1) = 0. \quad (2.1.10)$$

This equation, the analog of (2.1.7), can be solved using partial fractions, and we find that

$$m(x) = -2(x \log x + (1-x) \log(1-x)), \quad 0 < x < 1. \quad (2.1.11)$$

In terms of the underlying discrete model, the approximation for the expected number m_i of generations to fixation, starting from i A alleles, is $m_i \approx Nm(i/N)$. If $i/N = 1/2$,

$$Nm(1/2) = (-2 \log 2)N \approx 1.39N \text{ generations,}$$

whereas if the A allele is introduced at frequency $1/N$,

$$Nm(1/N) = 2 \log N \text{ generations.}$$

2.2 The genealogy of the Wright-Fisher model

In this section we consider the Wright-Fisher model from a genealogical perspective. In the absence of recombination, the DNA sequence representing the gene of interest is a copy of a sequence in the previous generation, that sequence is itself a copy of a sequence in the generation before that and so on. Thus we can think of the DNA sequence as an ‘individual’ that has a ‘parent’ (namely the sequence from which it was copied), and a number of ‘offspring’ (namely the sequences that originate as a copy of it in the next generation).

To study this process either forwards or backwards in time, it is convenient to label the individuals in a given generation as $1, 2, \dots, N$, and let ν_i denote the number of offspring born to individual i , $1 \leq i \leq N$. We suppose that individuals have independent Poisson-distributed numbers of offspring,

subject to the requirement that the total number of offspring is N . It follows that (ν_1, \dots, ν_N) has a symmetric multinomial distribution, with

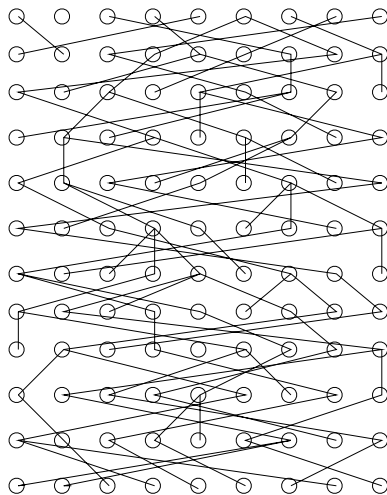
$$\mathbb{P}(\nu_1 = m_1, \dots, \nu_N = m_N) = \frac{N!}{m_1! \cdots m_N!} \left(\frac{1}{N}\right)^N \quad (2.2.1)$$

provided $m_1 + \cdots + m_N = N$. We assume that offspring numbers are independent from generation to generation, with distribution specified by (2.2.1).

To see the connection with the earlier description of the Wright-Fisher model, imagine that each individual in a given generation carries either an A allele or a B allele, i of the N individuals being labelled A . Since there is no mutation, all offspring of type A individuals are also of type A . The distribution of the number of type A in the offspring therefore has the distribution of $\nu_1 + \cdots + \nu_i$ which (from elementary properties of the multinomial distribution) has the binomial distribution with parameters N and success probability $p = i/N$. Thus the number of A alleles in the population does indeed evolve according to the Wright-Fisher model described in (2.1.1).

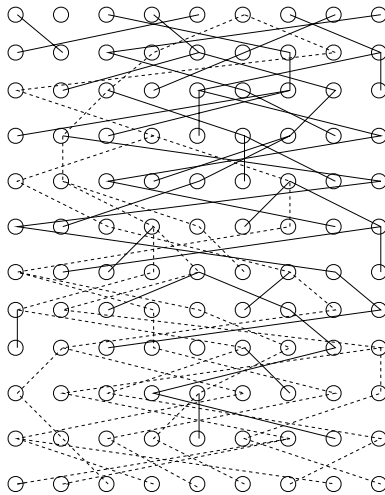
This specification shows how to simulate the offspring process from parents to children to grandchildren and so on. A realization of such a process for $N = 9$ is shown in Figure 2.1. Examination of Figure 2.1 shows that individuals 3 and 4 have their most recent common ancestor (MRCA) 3 generations ago, whereas individuals 2 and 3 have their MRCA 11 generations ago. More

Fig. 2.1. Simulation of a Wright-Fisher model of $N = 9$ individuals. Generations are evolving down the figure. The individuals in the last generation should be labelled 1, 2, ..., 9 from left to right. Lines join individuals in two generations if one is the offspring of the other



generally, for any population size N and sample of size n taken from the present generation, what is the structure of the ancestral relationships linking the members of the sample? The crucial observation is that if we view the process from the present generation back into the past, then individuals choose their parents independently and at random from the individuals in the previous generation, and successive choices are independent from generation to generation. Of course, not all members of the previous generations are ancestors of individuals in the present-day sample. In Figure 2.2 the ancestry of those individuals who are ancestral to the sample is highlighted with broken lines, and in Figure 2.3 those lineages that are not connected to the sample are removed, the resulting figure showing just the successful ancestors. Finally, Figure 2.3 is untangled in Figure 2.4. This last figure shows the tree-like nature of the genealogy of the sample.

Fig. 2.2. Simulation of a Wright-Fisher model of $N = 9$ individuals. Lines indicate ancestors of the sampled individuals. Individuals in the last generation should be labelled $1, 2, \dots, 9$ from left to right. Dashed lines highlight ancestry of the sample.



Understanding the genealogical process provides a direct way to study gene frequencies in a model with no mutation (Felsenstein (1971)). We content ourselves with a genealogical derivation of (2.1.6). To do this, we ask how long it takes for a sample of two genes to have their first common ancestor. Since individuals choose their parents at random, we see that

$$\mathbb{P}(\text{2 individuals have 2 distinct parents}) = \lambda = \left(1 - \frac{1}{N}\right).$$

Fig. 2.3. Simulation of a Wright-Fisher model of $N = 9$ individuals. Individuals in the last generation should be labelled 1,2,...,9 from left to right. Dashed lines highlight ancestry of the sample. Ancestral lineages not ancestral to the sample are removed.

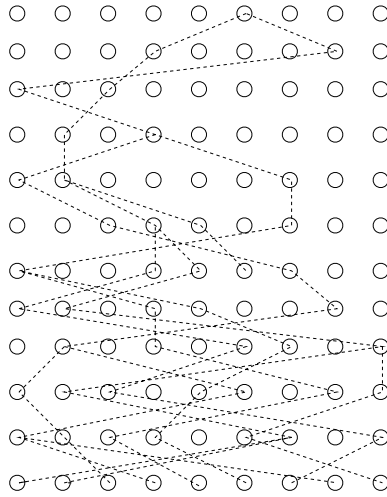
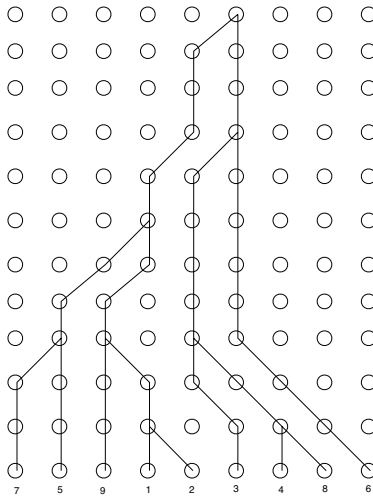


Fig. 2.4. Simulation of a Wright-Fisher model of $N = 9$ individuals. This is an untangled version of Figure 2.3.



Since those parents are themselves a random sample from their generation, we may iterate this argument to see that

$$\begin{aligned} & \mathbb{P}(\text{First common ancestor more than } r \text{ generations ago}) \\ &= \lambda^r = \left(1 - \frac{1}{N}\right)^r. \end{aligned} \tag{2.2.2}$$

Now consider the probability $h(r)$ that two individuals chosen with replacement from generation r carry distinct alleles. Clearly if we happen to choose the same individual twice (probability $1/N$) this probability is 0. In the other case, the two individuals are different if and only if their common ancestor is more than r generations ago, and the ancestors at time 0 are distinct. The probability of this latter event is the chance that 2 individuals chosen without replacement at time 0 carry different alleles, and this is just $\mathbb{E}2X_0(N - X_0)/N(N - 1)$. Combining these results gives

$$h(r) = \lambda^r \frac{(N - 1)}{N} \frac{\mathbb{E}2X_0(N - X_0)}{N(N - 1)} = \lambda^r h(0),$$

just as in (2.1.6).

When the population size is large and time is measured in units of N generations, the distribution of the time to the MRCA of a sample of size 2 has approximately an exponential distribution with mean 1. To see this, rescale time so that $r = Nt$, and let $N \rightarrow \infty$ in (2.2.2). We see that this probability is

$$\left(1 - \frac{1}{N}\right)^{Nt} \rightarrow e^{-t}.$$

This time scaling is the same as used to derive the diffusion approximation earlier. This should be expected, as the forward and backward approaches are just alternative views of the same underlying process.

The ancestral process in a large population

What can be said about the number of ancestors in larger samples? The probability that a sample of size three has distinct parents is

$$\left(1 - \frac{1}{N}\right) \left(1 - \frac{2}{N}\right)$$

and the iterative argument above can be applied once more to see that the sample has three distinct ancestors for more than r generations with probability

$$\left[\left(1 - \frac{1}{N}\right) \left(1 - \frac{2}{N}\right)\right]^r = \left(1 - \frac{3}{N} + \frac{2}{N^2}\right)^r.$$

Rescaling time once more in units of N generations, and taking $r = Nt$, shows that for large N this probability is approximately e^{-3t} , so that on the new time scale the time taken to find the first common ancestor in the sample of three genes is exponential with parameter 3. What happens when a common ancestor is found? Note that the chance that three distinct individuals have at most two distinct parents is

$$\frac{3(N-1)}{N^2} + \frac{1}{N^2} = \frac{3N-2}{N^2}.$$

Hence, given that a first common ancestor is found in generation r , the conditional probability that the sample has two distinct ancestors in generation r is

$$\frac{3N-3}{3N-2},$$

which tends to 1 as N increases. Thus in our approximating process the number of distinct ancestors drops by precisely 1 when a common ancestor is found.

We can summarize the discussion so far by noting that in our approximating process a sample of three genes waits an exponential amount of time T_3 with parameter 3 until a common ancestor is found, at which point the sample has two distinct ancestors for a further amount of time T_2 having an exponential distribution with parameter 1. Furthermore, T_3 and T_2 are independent random variables.

More generally, the number of distinct parents of a sample of size k individuals can be thought of as the number of occupied cells after k balls have been dropped (uniformly and independently) into N cells. Thus

$$g_{kj} \equiv \mathbb{P}(k \text{ individuals have } j \text{ distinct parents}) \tag{2.2.3}$$

$$= N(N-1)\cdots(N-j+1)\mathfrak{S}_k^{(j)}N^{-k} \quad j = 1, 2, \dots, k$$

where $\mathfrak{S}_k^{(j)}$ is a Stirling number of the second kind; that is, $\mathfrak{S}_k^{(j)}$ is the number of ways of partitioning a set of k elements into j nonempty subsets. The terms in (2.2.3) arise as follows: $N(N-1)\cdots(N-j+1)$ is the number of ways to choose j distinct parents; $\mathfrak{S}_k^{(j)}$ is the number of ways assigning k individuals to these j parents; and N^k is the total number of ways of assigning k individuals to their parents.

For fixed values of N , the behavior of this ancestral process is difficult to study analytically, but we shall see that the simple approximation derived above for samples of size two and three can be developed for any sample size n . We first define an ancestral process $\{A_n^N(t) : t = 0, 1, \dots\}$ where

$$A_n^N(t) \equiv \begin{array}{l} \text{number of distinct ancestors in generation } t \text{ of a} \\ \text{sample of size } n \text{ at time 0.} \end{array}$$

It is evident that $A_n^N(\cdot)$ is a Markov chain with state space $\{1, 2, \dots, n\}$, and with transition probabilities given by (2.2.3):

$$\mathbb{P}(A_n^N(t+1) = j | A_n^N(t) = k) = g_{kj}.$$

For fixed sample size n , as $N \rightarrow \infty$,

$$\begin{aligned} g_{k,k-1} &= \mathfrak{S}_k^{(k-1)} \frac{N(N-1) \cdots (N-k+2)}{N^k} \\ &= \binom{k}{2} \frac{1}{N} + O(N^{-2}), \end{aligned}$$

since $\mathfrak{S}_k^{(k-1)} = \binom{k}{2}$. For $j < k-1$, we have

$$g_{k,j} = \mathfrak{S}_k^{(j)} \frac{N(N-1) \cdots (N-j+1)}{N^k} = O(N^{-2})$$

and

$$\begin{aligned} g_{k,k} &= N^{-k} N(N-1) \cdots (N-k+1) \\ &= 1 - \binom{k}{2} \frac{1}{N} + O(N^{-2}). \end{aligned}$$

Writing G_N for the transition matrix with elements g_{kj} , $1 \leq j \leq k \leq n$. Then

$$G_N = I + N^{-1}Q + O(N^{-2}),$$

where I is the identity matrix, and Q is a lower diagonal matrix with non-zero entries given by

$$q_{kk} = -\binom{k}{2}, \quad q_{k,k-1} = \binom{k}{2}, \quad k = n, n-1, \dots, 2. \quad (2.2.4)$$

Hence with time rescaled for units of N generations, we see that

$$G_N^{Nt} = (I + N^{-1}Q + O(N^{-2}))^{Nt} \rightarrow e^{Qt}$$

as $N \rightarrow \infty$. Thus the number of distinct ancestors in generation Nt is approximated by a Markov chain $A_n(t)$ whose behavior is determined by the matrix Q in (2.2.4). $A_n(\cdot)$ is a pure death process that starts from $A_n(0) = n$, and decreases by jumps of size one only. The waiting time T_k in state k is exponential with parameter $\binom{k}{2}$, the T_k being independent for different k .

Remark. We call the process $A_n(t)$, $t \geq 0$ the *ancestral process* for a sample of size n .

Remark. The ancestral process of the Wright-Fisher model has been studied in several papers, including Karlin and McGregor (1972), Cannings (1974), Watterson (1975), Griffiths (1980), Kingman (1980) and Tavaré (1984).

2.3 Properties of the ancestral process

Calculation of the distribution of $A_n(t)$ is an elementary exercise in Markov chains. One way to do this is to diagonalize the matrix Q by writing $Q = RDL$, where D is the diagonal matrix of eigenvalues $\lambda_k = -\binom{k}{2}$ of Q , and R and L are matrices of right and left eigenvalues of Q , normalized so that $RL = LR = I$. From this approach we get, for $j = 1, 2, \dots, n$,

$$g_{nj}(t) \equiv \mathbb{P}(A_n(t) = j) = \sum_{k=j}^n e^{-k(k-1)t/2} \frac{(2k-1)(-1)^{k-j} j_{(k-1)} n_{[k]}}{j!(k-j)!n_{(k)}} \quad (2.3.1)$$

where

$$\begin{aligned} a_{(n)} &= a(a+1) \cdots (a+n-1) \\ a_{[n]} &= a(a-1) \cdots (a-n+1) \\ a_{(0)} &= a_{[0]} = 1. \end{aligned}$$

The mean number of ancestors at time t is given by

$$\mathbb{E}A_n(t) = \sum_{k=1}^n e^{-k(k-1)t/2} \frac{(2k-1)n_{[k]}}{n_{(k)}}, \quad (2.3.2)$$

and the falling factorial moments are given by

$$\mathbb{E}(A_n(t))_{[r]} = \sum_{k=r}^n \frac{n_{[k]}}{n_{(k)}} e^{-k(k-1)t/2} (2k-1) \frac{(r+k-2)!}{(r-1)!(k-r)!},$$

for $r = 2, \dots, n$. In Figure 2.5 $\mathbb{E}A_n(t)$ is plotted as a function of t for $n = 5, 10, 20, 50$.

The process $A_n(\cdot)$ is eventually absorbed at 1, when the sample is traced back to its most recent common ancestor (MRCA). The time it takes the sample to reach its MRCA is of some interest to population geneticists. We study this time in the following section.

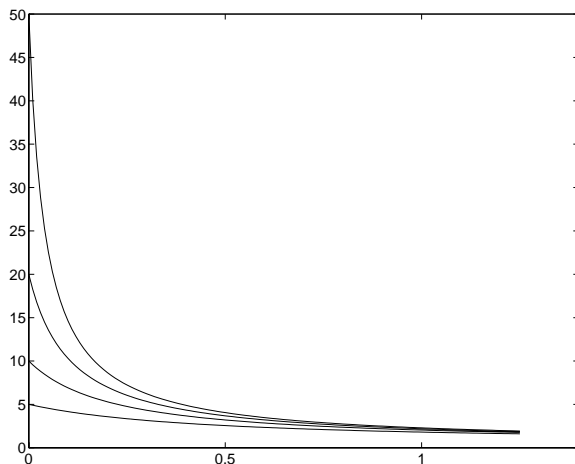
The time to the most recent common ancestor

Many quantities of genetic interest depend on the time W_n taken to trace a sample of size n back to its MRCA. Remember that time here is measured in units of N generations, and that

$$W_n = T_n + T_{n-1} + \cdots + T_2 \quad (2.3.3)$$

where T_k are independent exponential random variables with parameter $\binom{k}{2}$. It follows that

Fig. 2.5. The mean number of ancestors at time t (x axis) for samples of size $n = 5, 10, 20, 50$, from (2.3.2).



$$\mathbb{E}W_n = \sum_{k=2}^n \mathbb{E}T_k = \sum_{k=2}^n \frac{2}{k(k-1)} = 2 \sum_{k=2}^n \left(\frac{1}{k-1} - \frac{1}{k} \right) = 2 \left(1 - \frac{1}{n} \right).$$

Therefore

$$1 = \mathbb{E}W_2 \leq \mathbb{E}W_n \leq \mathbb{E}W_N < 2,$$

where W_N is thought of as the time until the whole population has a single common ancestor. Note that $\mathbb{E}W_n$ is close to 2 even for moderate n . Also

$$\mathbb{E}(W_N - W_n) = 2 \left(\frac{1}{n} - \frac{1}{N} \right) < \frac{2}{n}$$

so the mean difference between the time for a sample to reach its MRCA, and the time for the whole population to reach its MRCA, is small.

Note that T_2 makes a substantial contribution to the sum (2.3.3) defining W_n . For example, on average for over half the time since its MRCA, the sample will have exactly two ancestors. Further, using the independence of the T_k ,

$$\begin{aligned} \text{Var}W_n &= \sum_{k=2}^n \text{Var}T_k = \sum_{k=2}^n \binom{k}{2}^{-2} \\ &= 8 \sum_{k=1}^{n-1} \frac{1}{k^2} - 4 \left(1 - \frac{1}{n} \right) \left(3 + \frac{1}{n} \right) \end{aligned}$$

It follows that

$$1 = \text{Var}W_2 \leq \text{Var}W_n \leq \lim_{n \rightarrow \infty} \text{Var}W_n = 8 \frac{\pi^2}{6} - 12 \approx 1.16.$$

We see that T_2 also contributes most to the variance.

The distribution of W_n can be obtained from (2.3.1):

$$\mathbb{P}(W_n \leq t) = \mathbb{P}(A_n(t) = 1) = \sum_{k=1}^n e^{-k(k-1)t/2} \frac{(2k-1)(-1)^{k-1}n_{[k]}}{n_{(k)}}. \quad (2.3.4)$$

From this it follows that

$$\mathbb{P}(W_n > t) = 3 \frac{n-1}{n+1} e^{-t} + O(e^{-3t}) \text{ as } t \rightarrow \infty.$$

Now focus on two particular individuals in the sample and observe that if these two individuals do not have a common ancestor at t , the whole sample cannot have a common ancestor. Since the two individuals are themselves a random sample of size two from the population, we see that

$$\mathbb{P}(W_n > t) \geq \mathbb{P}(W_2 > t) = e^{-t},$$

an inequality that also follows from (2.3.3). A simple Markov chain argument shows that

$$\mathbb{P}(W_n > t) \leq \frac{3(n-1)e^{-t}}{n+1},$$

so that

$$e^{-t} \leq \mathbb{P}(W_n > t) \leq 3e^{-t}$$

for all n and t (see Kingman (1980), (1982c)).

The density function of W_n follows immediately from (2.3.4) by differentiating with respect to t :

$$f_{W_n}(t) = \sum_{k=2}^n (-1)^k e^{-k(k-1)t/2} \frac{(2k-1)k(k-1)n_{[k]}}{2n_{(k)}}. \quad (2.3.5)$$

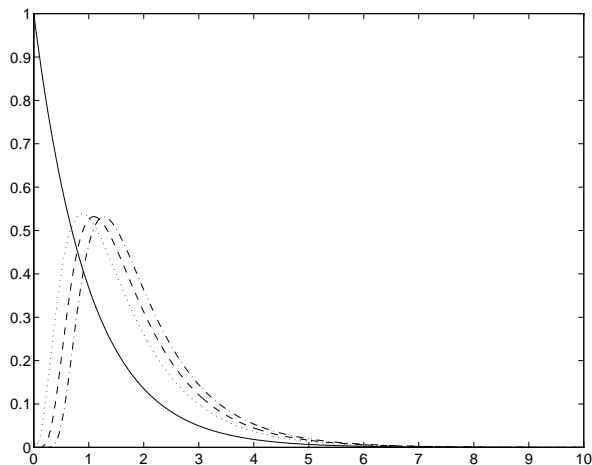
In Figure 2.6, this density is plotted for values of $n = 2, 10, 100, 500$. The shape of the densities reflects the fact that most of the contribution to the density comes from T_2 .

The tree length

In contrast to the distribution of W_n , the distribution of the total length $L_n = 2T_2 + \dots + nT_n$ is easy to find. As we will see, L_n is the total length of the branches in the genealogical tree linking the individuals in the sample. First of all,

$$\mathbb{E}L_n = 2 \sum_{j=1}^{n-1} \frac{1}{j} \sim 2 \log n,$$

Fig. 2.6. Density functions for the time W_n to most recent common ancestor of a sample of n individuals, from (2.3.5). — $n = 2$; $\cdots \cdots$ $n = 10$; — — — $n = 100$; — · — · $n = 500$.



and

$$\text{Var}L_n = 4 \sum_{j=1}^{n-1} \frac{1}{j^2} \sim 2\pi^2/3.$$

To find the distribution of L_n , let $E(\lambda)$ denote an exponential random variable with mean $1/\lambda$, all occurrences being independent of each other, and write $=_d$ for equality in distribution. Then

$$\begin{aligned} L_n &= \sum_{j=2}^n jT_j =_d \sum_{j=2}^n E((j-1)/2) \\ &=_d \sum_{j=1}^{n-1} \min_{1 \leq k \leq j} E_{jk}(1/2) \\ &=_d \max_{1 \leq j \leq n-1} E_j(1/2), \end{aligned}$$

the last step following by a coupling argument (this is one of many proofs of Feller's representation of the distribution of the maximum of independent and identically distributed exponential random variables as a sum of independent random variables). Thus

$$\mathbb{P}(L_n \leq t) = \left(1 - e^{-t/2}\right)^{n-1}, \quad t \geq 0.$$

It follows directly that $L_n - 2 \log n$ has a limiting extreme value distribution with distribution function $\exp(-\exp(-t/2))$, $-\infty < t < \infty$.

2.4 Variable population size

In this section we discuss the behavior of the ancestral process in the case of deterministic fluctuations in population size. For convenience, suppose the model evolves in discrete generations and label the current generation as 0. Denote by $N(j)$ the number of sequences in the population j generations before the present. We assume that the variation in population size is due to either external constraints *e.g.* changes in the environment, or random variation which depends only on the total population size *e.g.* if the population grows as a branching process. This excludes so-called density dependent cases in which the variation depends on the genetic composition of the population, but covers many other settings. We continue to assume neutrality and random mating.

Here we develop the theory for a particular class of population growth models in which, roughly speaking, all the population sizes are large. Time will be scaled in units of $N \equiv N(0)$ generations. To this end, define the relative size function $f_N(x)$ by

$$\begin{aligned} f_N(x) &= \frac{N(\lceil Nx \rceil)}{N} \\ &= \frac{N(j)}{N}, \quad \frac{j-1}{N} < x \leq \frac{j}{N}, \quad j = 1, 2, \dots \end{aligned} \quad (2.4.1)$$

We are interested in the behavior of the process when the size of each generation is large, so we suppose that

$$\lim_{N \rightarrow \infty} f_N(x) = f(x) \quad (2.4.2)$$

exists and is strictly positive for all $x \geq 0$.

Many demographic scenarios can be modelled in this way. For an example of geometric population growth, suppose that for some constant $\rho > 0$

$$N(j) = \lfloor N(1 - \rho/N)^j \rfloor.$$

Then

$$\lim_{N \rightarrow \infty} f_N(x) = e^{-\rho x} \equiv f(x), \quad x > 0.$$

A commonly used model is one in which the population has constant size prior to generation V , and geometric growth from then to the present time. Thus for some $\alpha \in (0, 1)$

$$N(j) = \begin{cases} \lfloor N\alpha \rfloor, & j \geq V \\ \lfloor N\alpha^{j/V} \rfloor, & j = 0, \dots, V \end{cases}$$

If we suppose that $V = \lfloor Nv \rfloor$ for some $v > 0$, so that the expansion started v time units ago, then

$$f_N(x) \rightarrow f(x) = \alpha^{\min(x/v, 1)}.$$

The ancestral process

In a Wright-Fisher model of reproduction, note that the probability that two individuals chosen at time 0 have distinct ancestors s generations ago is

$$\mathbb{P}(T_2(N) > s) = \prod_{j=1}^s \left(1 - \frac{1}{N(j)}\right),$$

where $T_2(N)$ denotes the time to the common ancestor of the two individuals. Recalling the inequality

$$x \leq -\log(1-x) \leq \frac{x}{1-x}, \quad x < 1,$$

we see that

$$\sum_{j=1}^s \frac{1}{N(j)} \leq -\sum_{j=1}^s \log\left(1 - \frac{1}{N(j)}\right) \leq \sum_{j=1}^s \frac{1}{N(j)-1}.$$

It follows that

$$\lim_{N \rightarrow \infty} -\sum_{j=1}^{\lfloor Nt \rfloor} \log\left(1 - \frac{1}{N(j)}\right) = \lim_{N \rightarrow \infty} \sum_{j=1}^{\lfloor Nt \rfloor} \frac{1}{N(j)}.$$

Since

$$\sum_{j=1}^s \frac{1}{N(j)} = \int_0^{s/N} \frac{dx}{f_N(x)},$$

we can use (2.4.2) to see that for $t > 0$, with time rescaled in units of N generations,

$$\lim_{N \rightarrow \infty} \mathbb{P}(T_2(N) > \lfloor Nt \rfloor) = \exp\left(-\int_0^t \lambda(u) du\right),$$

where $\lambda(\cdot)$ is the intensity function defined by

$$\lambda(u) = \frac{1}{f(u)}, \quad u \geq 0. \tag{2.4.3}$$

If we define

$$\Lambda(t) = \int_0^t \lambda(u) du,$$

the integrated intensity function, then (2.4.2) shows that as $N \rightarrow \infty$

$$N^{-1}T_2(N) \Rightarrow T_2,$$

where

$$\mathbb{P}(T_2 > t) = \exp(-\Lambda(t)), \quad t \geq 0. \tag{2.4.4}$$

We expect the two individuals to have a common ancestor with probability one, this corresponding to the requirement that

$$\lim_{t \rightarrow \infty} \Lambda(t) = \infty,$$

which we assume from now on. When the population size is constant, $\Lambda(t) = t$ and the time to the MRCA has an exponential distribution with mean 1. From (2.4.4) we see that

$$\mathbb{E}T_2 = \int_0^\infty \mathbb{P}(T_2 > t) dt = \int_0^\infty e^{-\Lambda(t)} dt.$$

If the population has been expanding, so that $f(t) \leq 1$ for all t , then $\Lambda(t) \geq t$, and therefore

$$\mathbb{P}(T_2 > t) \leq \mathbb{P}(T_2^c > t), \quad t \geq 0,$$

where T_2^c denotes the corresponding time in the constant population size case. We say that T_2^c is *stochastically larger* than T_2 , so that in particular $\mathbb{E}T_2 \leq \mathbb{E}T_2^c = 1$. This corresponds to the fact that if the population size has been shrinking into the past, it should be possible to find the MRCA sooner than if the population size had been constant.

In the varying environment setting, the ancestral process satisfies

$$\begin{aligned} \mathbb{P}(A_2(t+s) = 1 | A_2(t) = 2) &= \mathbb{P}(T_2 \leq t+s | T_2 > t) \\ &= \mathbb{P}(t < T_2 \leq t+s) / \mathbb{P}(T_2 > t) \\ &= 1 - \exp(-(\Lambda(t+s) - \Lambda(t))), \end{aligned}$$

so that

$$\mathbb{P}(A_2(t+h) = 1 | A_2(t) = 2) = \lambda(t)h + o(h), \quad h \downarrow 0.$$

We see that $A_2(\cdot)$ is a non-homogeneous Markov chain. What is the structure of $A_n(\cdot)$?

Define $T_k(N)$ to be the number of generations for which the sample has k distinct ancestors. In the event that the sample never has exactly k distinct ancestors, define $T_k(N) = \infty$. We calculate first the joint distribution of $T_3(N)$ and $T_2(N)$. The probability that $T_3(N) = k, T_2(N) = l$ is the probability that the sample of size 3 has 3 distinct ancestors in generations 1, 2, ..., $k-1$, 2 distinct ancestors in generations $k, \dots, k+l-1$, and 1 in generation $l+k$. The probability that a sample of three individuals taken in generation

$j - 1$ has three distinct parents is $N(j)(N(j) - 1)(N(j) - 2)/N(j)^3$, and the probability that three individuals in generation $k - 1$ have two distinct parents is $3N(k)(N(k) - 1)/N(k)^3$. Hence

$$\begin{aligned} & \mathbb{P}(T_3(N) = k, T_2(N) = l) \\ &= \left\{ \prod_{j=1}^{k-1} \frac{(N(j) - 1)(N(j) - 2)}{N(j)^3} \right\} \frac{3(N(k) - 1)}{N(k)^2} \left\{ \prod_{j=k+1}^{k+l-1} \frac{N(j) - 1}{N(j)} \right\} \frac{1}{N(k+l)}. \end{aligned}$$

For the size fluctuations we are considering, the first term in brackets is

$$\prod_{j=1}^{k-1} \left(1 - \frac{3}{N(j)} + \frac{2}{N(j)^2} \right) \sim \exp \left(-3 \int_0^{k/N} \frac{dx}{f_N(x)} \right),$$

while the second term in brackets is

$$\prod_{j=k+1}^{k+l-1} \left(1 - \frac{1}{N(j)} \right) \sim \exp \left(- \int_{k/N}^{(k+l)/N} \frac{dx}{f_N(x)} \right).$$

For $k \sim Nt_3, l \sim Nt_2$ with $t_3 > 0, t_2 > 0$, we see via (2.4.2) that $N^2 \mathbb{P}(T_3(N) = k, T_2(N) = l)$ converges to

$$f(t_3, t_2) := e^{-3\Lambda(t_3)} 3\lambda(t_3) e^{-(\Lambda(t_2+t_3) - \Lambda(t_3))} \lambda(t_3 + t_2) \quad (2.4.5)$$

as $N \rightarrow \infty$. It follows that

$$N^{-1}(T_3(N), T_2(N)) \Rightarrow (T_3, T_2),$$

where (T_3, T_2) have joint probability density $f(t_3, t_2)$ given in (2.4.5).

This gives the joint law of the times spent with different numbers of ancestors, and shows that in the limit the number of ancestors decreases by one at each jump. Just as in the constant population-size case, the ancestral process for the Wright-Fisher model is itself a Markov chain, since the distribution of the number of distinct ancestors in generation r is determined just by the number in generation $r - 1$. The Markov property is inherited in the limit, and we conclude that $\{A_3(t), t \geq 0\}$ is a Markov chain on the set $\{3, 2, 1\}$. Its transition intensities can be calculated as a limit from the Wright-Fisher model. We see that

$$\mathbb{P}(A_3(t+h) = j | A_3(t) = i) = \begin{cases} \binom{i}{2} \lambda(t) h + o(h), & j = i - 1 \\ 1 - \binom{i}{2} \lambda(t) h + o(h), & j = i \\ 0, & \text{otherwise} \end{cases}$$

We can now establish the general case in a similar way. The random variables $T_n(N), \dots, T_2(N)$ have a joint limit law when rescaled:

$$N^{-1}(T_n(N), \dots, T_2(N)) \Rightarrow (T_n, \dots, T_2)$$

for each fixed n as $N \rightarrow \infty$, and the joint density $f(t_n, \dots, t_2)$ of T_n, \dots, T_2 is given by

$$f(t_n, \dots, t_2) = \prod_{j=2}^n \binom{j}{2} \lambda(s_j) \exp \left\{ - \binom{j}{2} (\Lambda(s_j) - \Lambda(s_{j+1})) \right\}, \quad (2.4.6)$$

for $0 \leq t_n, \dots, t_2 < \infty$, where $s_{n+1} = 0, s_n = t_n, s_j = t_j + \dots + t_n, j = 2, \dots, n-1$.

Remark. The joint density in (2.4.6) should really be denoted by $f_n(t_n, \dots, t_2)$, and the limiting random variables T_{n1}, \dots, T_{n2} , but we keep the simpler notation. This should not cause any confusion.

From this it is elementary to show that if $S_j \equiv T_n + \dots + T_j$, then the joint density of (S_n, \dots, S_2) is given by

$$g(s_n, \dots, s_2) = \prod_{j=2}^n \binom{j}{2} \lambda(s_j) \exp \left(- \binom{j}{2} (\Lambda(s_j) - \Lambda(s_{j+1})) \right),$$

for $0 \leq s_n < s_{n-1} < \dots < s_2$. This parlays immediately into the distribution of the time the sample spends with j distinct ancestors, given that $S_{j+1} = s$:

$$\mathbb{P}(T_j > t | S_{j+1} = s) = \exp \left(- \binom{j}{2} (\Lambda(s+t) - \Lambda(s)) \right).$$

Note that the sequence $S_{n+1} := 0, S_n, S_{n-1}, \dots, S_2$ is a Markov chain. The approximating ancestral process $\{A_n(t), t \geq 0\}$ is a non-homogeneous pure death process on $[n]$ with $A_n(0) = n$ whose transition rates are determined by

$$\mathbb{P}(A_n(t+h) = j | A_n(t) = i) = \begin{cases} \binom{i}{2} \lambda(t)h + o(h), & j = i-1 \\ 1 - \binom{i}{2} \lambda(t)h + o(h), & j = i \\ 0, & \text{otherwise} \end{cases} \quad (2.4.7)$$

The time change representation

Denote the process that counts the number of ancestors at time t of a sample of size n taken at time 0 by $\{A_n^v(t), t \geq 0\}$, the superscript v denoting variable population size. We have seen that $A_n^v(\cdot)$ is now a time-inhomogeneous Markov process. Given that $A_n^v(t) = j$, it jumps to $j-1$ at rate $j(j-1)\lambda(t)/2$. A useful way to think of the process $A_n^v(\cdot)$ is to notice that a realization may be constructed via

$$A_n^v(t) = A_n(\Lambda(t)), \quad t \geq 0, \quad (2.4.8)$$

where $A_n(\cdot)$ is the corresponding ancestral process for the constant population size case. This may be verified immediately from (2.4.7). We see that the variable population size model is just a deterministic time change of the constant

population size model. Some of the properties of $A_n^v(\cdot)$ follow immediately from this representation. For example,

$$\mathbb{P}(A_n^v(t) = j) = g_{nj}(\Lambda(t)), \quad j = 1, \dots, n$$

where $g_{nj}(t)$ is given in (2.3.1), and so

$$\mathbb{E}A_n^v(t) = \sum_{j=1}^n e^{-j(j-1)\Lambda(t)/2} \frac{(2l-1)n[j]}{n(j)}, t \geq 0.$$

It follows from (2.4.8) that $A_n(s) = A_n^v(\Lambda^{-1}(s))$, $s > 0$. Hence if $A_n(\cdot)$ has a jump at time s , then $A_n^v(\cdot)$ has one at time $\Lambda^{-1}(s)$. Since $A_n(\cdot)$ has jumps at $S_n = T_n, S_{n-1} = T_n + T_{n-1}, \dots, S_2 = T_n + \dots + T_2$, it follows that the jumps of $A_n^v(\cdot)$ occur at $\Lambda^{-1}(S_n), \dots, \Lambda^{-1}(S_2)$. Thus, writing T_j^v for the time the sample from a variable-size population spends with j ancestors, we see that

$$\begin{aligned} T_n^v &= \Lambda^{-1}(S_n) \\ T_j^v &= \Lambda^{-1}(S_j) - \Lambda^{-1}(S_{j+1}), \quad j = n-1, \dots, 2. \end{aligned} \tag{2.4.9}$$

This result provides a simple way to simulate the times $T_n^v, T_{n-1}^v, \dots, T_2^v$. Let U_n, \dots, U_2 be independent and identically distributed random variables having the uniform distribution on $(0,1)$.

Algorithm 2.1 Algorithm to generate T_n^v, \dots, T_2^v for a variable size process with intensity function Λ :

1. Generate $t_j = -\frac{2 \log(U_j)}{j(j-1)}$, $j = 2, 3, \dots, n$
2. Form $s_n = t_n, s_j = t_j + \dots + t_n$, $j = 2, \dots, n-1$
3. Compute $t_n^v = \Lambda^{-1}(s_n), t_j^v = \Lambda^{-1}(s_j) - \Lambda^{-1}(s_{j+1})$, $j = n-1, \dots, 2$.
4. Return $T_j^v = t_j^v, j = 2, \dots, n$.

There is also a sequential version of the algorithm, essentially a restatement of the last one:

Algorithm 2.2 Step-by-step version of Algorithm 2.1.

1. Set $t = 0, j = n$
2. Generate $t_j = -\frac{2 \log(U_j)}{j(j-1)}$
3. Solve for s the equation

$$\Lambda(t+s) - \Lambda(t) = t_j \tag{2.4.10}$$

4. Set $t_j^v = s, t = t + s, j = j - 1$. If $j \geq 2$, go to 2. Else return $T_n^v = t_n^v, \dots, T_2^v = t_2^v$.

Note that t_j generated in step 2 above has an exponential distribution with parameter $j(j-1)/2$. If the population size is constant then $\Lambda(t) = t$, and so $t_j^v = t_j$, as it should.

Example For an exponentially growing population $f(x) = e^{-\rho x}$, so that $\Lambda(t) = (e^{\rho t} - 1)/\rho$. It follows that $\Lambda^{-1}(y) = \rho^{-1} \log(1 + \rho y)$, and

$$T_n^v = \rho^{-1} \log(1 + \rho T_n), \quad T_j^v = \frac{1}{\rho} \left(\frac{1 + \rho S_j}{1 + \rho S_{j+1}} \right), \quad j = 2, \dots, n-1. \quad (2.4.11)$$

In an exponentially growing population, most of the coalescence events occur near the root of the tree, and the resulting genealogy is then star-like; it is harder to find common ancestors when the population size is large. See Section 4.2 for further illustrations.

3 The Ewens Sampling Formula

In this section we bring mutation into the picture, and show how the genealogical approach can be used to derive the classical Ewens sampling formula. This serves as an introduction to statistical inference for molecular data based obtained from samples.

3.1 The effects of mutation

In Section 2.1 we looked briefly at the process of random drift, the mechanism by which genetic variability is lost through the effects of random sampling. In this section, we study the effect of mutation on the evolution of gene frequencies at a locus with two alleles.

Now we suppose there is a probability $\mu_A > 0$ that an A allele mutates to a B allele in a single generation, and a probability $\mu_B > 0$ that a B allele mutates to an A . The stochastic model for the frequency X_n of the A allele in generation n is described by the transition matrix in (2.1.1), but where

$$\pi_i = \frac{i}{N}(1 - \mu_A) + \left(1 - \frac{i}{N}\right)\mu_B. \quad (3.1.1)$$

The frequency π_i reflects the effects of mutation in the gene pool. In this model, it can be seen that $p_{ij} > 0$ for all $i, j \in \mathcal{S}$. It follows that the Markov chain $\{X_n\}$ is irreducible; it is possible to get from any state to any other state. An irreducible finite Markov chain has a limit distribution $\rho = (\rho_0, \rho_1, \dots, \rho_N)$:

$$\lim_{n \rightarrow \infty} \mathbb{P}(X_n = k) = \rho_k > 0,$$

for any initial distribution for X_0 . The limit distribution ρ is also invariant (or stationary), in that if X_0 has distribution ρ then X_n has distribution ρ for every n . The distribution ρ satisfies the balance equations

$$\rho = \rho P,$$

where $\rho_0 + \dots + \rho_N = 1$.

Once more, the binomial conditional distributions make some aspects of the process simple to calculate. For example,

$$\mathbb{E}(X_n) = \mathbb{E}\mathbb{E}(X_n | X_{n-1}) = N\mu_B + (1 - \mu_A - \mu_B)\mathbb{E}(X_{n-1}).$$

At stationarity, $\mathbb{E}(X_n) = \mathbb{E}(X_{n-1}) \equiv \mathbb{E}(X)$, so

$$\mathbb{E}(X) = \frac{N\mu_B}{\mu_A + \mu_B}. \quad (3.1.2)$$

This is also the limiting value of $\mathbb{E}(X_n)$ as $n \rightarrow \infty$.

Now we investigate the stationary distribution ρ when N is large. To get a non-degenerate limit, we assume that the mutation probabilities μ_A and μ_B satisfy

$$\lim_{N \rightarrow \infty} 2N\mu_A = \theta_A > 0, \quad \lim_{N \rightarrow \infty} 2N\mu_B = \theta_B > 0, \quad (3.1.3)$$

so that mutation rates are of the order of the reciprocal of the population size. We define the total mutation rate θ by

$$\theta = \theta_A + \theta_B.$$

Given $X_n = i$, X_{n+1} is binomially distributed with parameters N and π_i given by (3.1.1). Exploiting simple properties of the binomial distribution shows that the diffusion approximation for the fraction of allele A in the population has

$$\mu(x) = -x\theta_A/2 + (1-x)\theta_B/2, \quad \sigma^2(x) = x(1-x), \quad 0 < x < 1. \quad (3.1.4)$$

The stationary density $\pi(y)$ of $Y(\cdot)$ satisfies the ordinary differential equation

$$-\mu(y)\pi(y) + \frac{1}{2} \frac{d\{\sigma^2(y)\pi(y)\}}{dy} = 0,$$

and it follows readily that

$$\pi(y) \propto \frac{1}{\sigma^2(y)} \exp\left(\int^y 2 \frac{\mu(u)}{\sigma^2(u)} du\right).$$

Hence $\pi(y) \propto y^{\theta_B-1}(1-y)^{\theta_A-1}$ and we see that at stationarity the fraction of A alleles has the beta distribution with parameters θ_B and θ_A . The density π is given by

$$\pi(y) = \frac{\Gamma(\theta)}{\Gamma(\theta_A)\Gamma(\theta_B)} y^{\theta_B-1}(1-y)^{\theta_A-1}, \quad 0 < y < 1.$$

In particular,

$$\mathbb{E}(Y) = \frac{\theta_B}{\theta}, \quad \text{Var}(Y) = \frac{\theta_A\theta_B}{\theta^2(\theta+1)}. \quad (3.1.5)$$

Remark. An alternative description of the mutation model in this case is as follows. Mutations occur at rate $\theta/2$, and when a mutation occurs the resulting allele is A with probability π_A and B with probability π_B . This model can be identified with the earlier one with $\theta_A = \theta\pi_A$, $\theta_B = \theta\pi_B$.

Remark. In the case of the K -allele model with mutation rate $\theta/2$ and mutations resulting in allele A_i with probability $\pi_i > 0$, $i = 1, 2, \dots, K$, the stationary density of the (now $(K-1)$ -dimensional) diffusion is

$$\pi(y_1, \dots, y_K) = \frac{\Gamma(\theta)}{\Gamma(\theta\pi_1) \dots \Gamma(\theta\pi_K)} y_1^{\theta\pi_1-1} \dots y_K^{\theta\pi_K-1},$$

for $y_i > 0$, $i = 1, \dots, K$, $y_1 + \dots + y_K = 1$.

3.2 Estimating the mutation rate

Modern molecular techniques have made it possible to sample genomic variability in natural populations. As a result, we need to develop the appropriate sampling theory to describe the statistical properties of such samples. For the models described in this section, this is easy to do. If a sample of n chromosomes is drawn with replacement from a stationary population, it is straightforward to calculate the distribution of the number N_A of A alleles in the sample. This distribution follows from the fact that given the population frequency Y of the A allele, the sample is distributed like a binomial random variable with parameters n and Y . Thus

$$\mathbb{P}(N_A = k) = \mathbb{E} \left(\binom{n}{k} Y^k (1 - Y)^{n-k} \right).$$

Since Y has the Beta(θ_B, θ_A) density, we see that N_A has the Beta-Binomial distribution:

$$\mathbb{P}(N_A = k) = \binom{n}{k} \frac{\Gamma(\theta) \Gamma(k + \theta_B) \Gamma(n - k + \theta_A)}{\Gamma(\theta_B) \Gamma(\theta_A) \Gamma(n + \theta)}, \quad k = 0, 1, \dots, n. \quad (3.2.1)$$

It follows from this that

$$\mathbb{E}(N_A) = n \frac{\theta_B}{\theta}, \quad \text{Var}(N_A) = \frac{n(n + \theta) \theta_A \theta_B}{\theta^2 (\theta + 1)}. \quad (3.2.2)$$

The probability that a sample of size one is an A allele is just $p \equiv \theta_B/\theta$. Had we ignored the dependence in the sample, we might have assumed that the genes in the sample were independently labelled A with probability p . The number N_A of A s in the sample then has a binomial distribution with parameters n and p . If we wanted to estimate the parameter p , the natural estimator is $\hat{p} = N_A/n$, and

$$\text{Var}(\hat{p}) = p(1 - p)/n.$$

As $n \rightarrow \infty$, this variance tends to 0, so that \hat{p} is a (weakly) consistent estimator of p . Of course, the sampled genes are *not* independent, and the true variance of N_A/n is, from (3.2.2),

$$\text{Var}(N_A/n) = \left(1 + \frac{\theta}{n} \right) \frac{\theta_A \theta_B}{\theta^2 (1 + \theta)}.$$

It follows that $\text{Var}(N_A/n)$ tends to the positive limit $\text{Var}(Y)$ as $n \rightarrow \infty$. Indeed, N_A/n is not a consistent estimator of $p = \theta_A/\theta$, because (by the strong law of large numbers) $N_A/n \rightarrow Y$, the population frequency of the A allele. This simple example shows how strong the dependence in the sample can be, and shows why consistent estimators of parameters in this subject are

the exception rather than the rule. Consistency typically has to be generated, at least in principle, by sampling variability at many independent loci.

The example in this section is our first glimpse of the difficulties caused by the relatedness of sequences in the sample. This relatedness has led to a number of interesting approaches to estimation and inference for population genetics data. In the next sections we describe the Ewens sampling formula (Ewens (1972)), the first systematic treatment of the statistical properties of estimators of the compound mutation parameter θ .

3.3 Allozyme frequency data

By the late 1960s, it was possible to sample, albeit indirectly, the molecular variation in the DNA of a population. These data came in the form of allozyme frequencies. A sample of size n resulted in a set of genes in which differences between genes could be observed, but the precise nature of the differences was irrelevant. Two *Drosophila* allozyme frequency data sets, each having 7 distinct alleles, are given below:

- *D. tropicalis* Esterase-2 locus [$n = 298$]
234, 52, 4, 4, 2, 1, 1
- *D. simulans* Esterase-C locus [$n = 308$]
91, 76, 70, 57, 12, 1, 1

It is clear that these data come from different distributions. Of the first set, Sewall Wright (1978, p303) argued that

... the observations do not agree at all with the equal frequencies expected for neutral alleles in enormously large populations.

This raises the question of what shape these distributions should have under a neutral model. The answer to this was given by Ewens (1972). Because the labels are irrelevant, a sample of genes can be broken down into a set of alleles that occurred just once in the sample, another collection that occurred twice, and so on. We denote by $C_j(n)$ the number of alleles represented j times in the sample of size n . Because the sample has size n , we must have

$$C_1(n) + 2C_2(n) + \cdots + nC_n(n) = n.$$

In this section we derive the distribution of $(C_1(n), \dots, C_n(n))$, known as the Ewens Sampling Formula (henceforth abbreviated to ESF). To do this, we need to study the effects of mutations in the history of a sample.

Mutations on a genealogy

In Section 3 we will give a detailed description of the ancestral relationships among a sample of individuals. For now, we recall from the last section that in a large population, the number of distinct ancestors at times t in the past

is described by the ancestral process $A_n(t)$. It is clear by symmetry that when the ancestral process moves from k to $k - 1$, the two ancestors chosen to join are randomly chosen from the k possibilities. Thus the ancestral relationships among a sample of individuals can be represented as a random rooted bifurcating tree that starts with n leaves (or tips), and joins random pairs of ancestors together at times $T_n, T_n + T_{n-1}, \dots, W_n = T_n + \dots + T_2$. All the individuals in the sample are traced back to their most recent common ancestor at time W_n .

Next we examine the effects of mutation in the coalescent tree of a sample. Suppose that a mutation occurs with probability u per gene per generation. The expected number of mutations along a lineage of g generations is therefore gu . With time measured in units of N generations, this is of the form tNu which is finite if u is of order $1/N$. Just as in (3.1.3), we take

$$\theta = 2Nu$$

to be fixed as $N \rightarrow \infty$. In the discrete process, mutations arise in the ancestral lines independently on different branches of the genealogical tree. In the limit, it is clear that they arise at the points of independent Poisson processes of rate $\theta/2$ on each branch.

We can now superimpose mutations on the genealogical tree of the sample. For allozyme frequency data, we suppose that every mutation produces a type that has not been seen in the population before. One concrete way to achieve this is to label types by uniform random variables; whenever a mutation occurs, the resulting individual has a type that is uniformly distributed on $(0,1)$, independently of other labels. This model is an example of an *infinitely-many alleles model*.

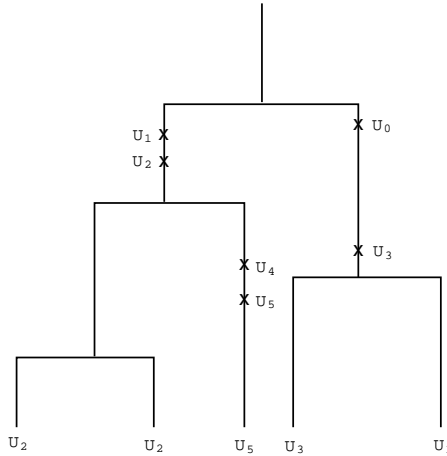
3.4 Simulating an infinitely-many alleles sample

As we will see, the reason that genealogical approaches have become so useful lies first in the fact that they provide a simple way to simulate samples from complex genetics models, and so to compare models with data. To simulate a sample, one need not simulate the whole population first and then sample from that – this makes these methods extremely appealing. Later in these notes we will see the same ideas applied in discrete settings as well, particularly for branching process models. This top down, or ‘goodness-of-fit’, approach has been used extensively since the introduction of the coalescent by Kingman (1982), Tajima (1983) and Hudson (1983) to simulate the behavior of test statistics which are intractable by analytical means.

To simulate samples of data following the infinitely-many-alleles model is, in principle, elementary. First simulate the genealogical tree of the sample by simulating observations from the waiting times T_n, T_{n-1}, \dots, T_2 and choosing pairs of nodes to join at random. Then we superimpose mutations according to a Poisson process of rate $\theta/2$, independently on each branch.

The effects of each mutation are determined by the mutation process. In the present case, the result of a mutation on a branch replaces the current label with an independently generated uniform random variable. An example is given in Figure 3.1, and the types represented in the sample are labelled U_5, U_2, U_2, U_3, U_3 respectively.

Fig. 3.1. A coalescent tree for $n = 5$ with mutations



3.5 A recursion for the ESF

To derive the ESF, we use a coalescent argument to find a recursion satisfied by the joint distribution of the sample configuration in an infinitely-many-alleles model. Under the infinitely-many-alleles mutation scheme, a sample of size n may be represented as a configuration $\mathbf{c} = (c_1, \dots, c_n)$, where

$$c_i = \text{number of alleles represented } i \text{ times}$$

and $|\mathbf{c}| \equiv c_1 + 2c_2 + \dots + nc_n = n$. It is convenient to think of the configuration \mathbf{b} of samples of size $j < n$ as being an n -vector with coordinates $(b_1, b_2, \dots, b_j, 0, \dots, 0)$, and we assume this in the remainder of this section. We define $\mathbf{e}_i = (0, 0, \dots, 0, 1, 0, \dots, 0)$, the i th unit vector.

We derive an equation satisfied by the sampling probabilities $q(\mathbf{c}), n = |\mathbf{c}| > 1$ defined by

$$q(\mathbf{c}) = \mathbb{P}(\text{sample of size } |\mathbf{c}| \text{ taken at stationarity has configuration } \mathbf{c}), \tag{3.5.1}$$

with $q(\mathbf{e}_1) = 1$. Suppose then that the configuration is \mathbf{c} . Looking at the history of the sample, we will either find a mutation or we will be able to

trace two individuals back to a common ancestor. The first event occurs with probability

$$\frac{n\theta/2}{n\theta/2 + n(n-1)/2} = \frac{\theta}{\theta + n - 1},$$

and results in the configuration \mathbf{c} if the configuration just before the mutation was \mathbf{b} , where

- (i) $\mathbf{b} = \mathbf{c}$, and mutation occurred to one of the c_1 singleton lines (probability c_1/n);
- (ii) $\mathbf{b} = \mathbf{c} - 2\mathbf{e}_1 + \mathbf{e}_2$, and a mutation occurred to an individual in the 2-class (probability $2(c_2 + 1)/n$);
- (iii) $\mathbf{b} = \mathbf{c} - \mathbf{e}_1 - \mathbf{e}_{j-1} + \mathbf{e}_j$ and the mutation occurred to an individual in a j -class, producing a singleton mutant and a new $(j-1)$ -class (probability $j(c_j + 1)/n$).

On the other hand, the ancestral join occurred with probability $(n-1)/(\theta + n - 1)$, and in that case the configuration $\mathbf{b} = \mathbf{c} + \mathbf{e}_j - \mathbf{e}_{j+1}$, and an individual in one of $c_j + 1$ allelic classes of size j had an offspring, reducing the number of j -classes to c_j , and increasing the number of $(j+1)$ -classes to c_{j+1} . This event has probability $j(c_j + 1)/(n-1)$, $j = 1, \dots, n-1$. Combining these possibilities, we get

$$q(\mathbf{c}) = \frac{\theta}{\theta + n - 1} \left[\frac{c_1}{n} q(\mathbf{c}) + \sum_{j=2}^n \frac{j(c_j + 1)}{n} q(\mathbf{c} - \mathbf{e}_1 - \mathbf{e}_{j-1} + \mathbf{e}_j) \right] + \frac{n-1}{\theta + n - 1} \left[\sum_{j=1}^{n-1} \frac{j(c_j + 1)}{n-1} q(\mathbf{c} + \mathbf{e}_j - \mathbf{e}_{j+1}) \right], \quad (3.5.2)$$

where we use the convention that $q(\mathbf{c}) = 0$ if any $c_i < 0$. Ewens (1972) established the following result.

Theorem 3.1 *In a stationary sample of size n , the probability of sample configuration \mathbf{c} is*

$$q(\mathbf{c}) = \mathbb{P}(C_1(n) = c_1, \dots, C_n(n) = c_n) = \mathbb{1}(|\mathbf{c}| = n) \frac{n!}{\theta_{(n)}} \prod_{j=1}^n \binom{\theta}{j}^{c_j} \frac{1}{c_j!}, \quad (3.5.3)$$

where (as earlier) we have written $x_{(j)} = x(x+1)\cdots(x+j-1)$, $j = 1, 2, \dots$, and $|\mathbf{c}| = c_1 + 2c_2 + \cdots + nc_n$.

Proof. This can be verified by induction on $n = |\mathbf{c}|$ and $k = \|\mathbf{c}\| := c_1 + \cdots + c_n$ in the equation (3.5.2) by noting that the right-hand side of the equation has terms with $|\mathbf{b}| = n-1$ and $\|\mathbf{b}\| \leq k$, or with $|\mathbf{b}| = n$ and $\|\mathbf{b}\| < k$.

Remark. Watterson (1974) noted that if Z_1, Z_2, \dots are independent Poisson random variables with $\mathbb{E}Z_j = \theta/j$, then

$$\mathcal{L}(C_1(n), C_2(n), \dots, C_n(n)) = \mathcal{L}\left(Z_1, Z_2, \dots, Z_n \mid \sum_{i=1}^n iZ_i = n\right), \quad (3.5.4)$$

where $\mathcal{L}(X)$ means ‘the distribution of X .’

The ESF typically has a very skewed distribution, assigning most mass to configurations with several alleles represented a few times. In particular, the distribution is far from ‘flat’; recall Wright’s observation cited in the introduction of this section. In the remainder of the section, we will explore some of the properties of the ESF.

Remark. The ESF arises in many other settings. See Tavaré and Ewens (1997) and Ewens and Tavaré (1998) for a flavor of this.

3.6 The number of alleles in a sample

The random variable $K_n = C_1(n) + \dots + C_n(n)$ is the number of distinct alleles observed in a sample. Its distribution can be found directly from (3.5.3):

$$\mathbb{P}(K_n = k) = \sum_{\mathbf{c}: |\mathbf{c}|=k} q(\mathbf{c}) = \frac{\theta^k}{\theta_{(n)}} n! \sum_{\mathbf{c}: |\mathbf{c}|=k} \left(\frac{1}{j}\right)^{c_j} \frac{1}{c_j!} = \frac{\theta^k |S_n^k|}{\theta_{(n)}}, \quad (3.6.1)$$

where $|S_n^k|$ is the Stirling number of the first kind,

$$|S_n^k| = \text{coefficient of } x^k \text{ in } x(x+1)\cdots(x+n-1),$$

and the last equality follows from Cauchy’s formula for the number of permutations of n symbols having k distinct cycles.

Another representation of the distribution of K_n can be found by noting that

$$\begin{aligned} \mathbb{E}_S^{K_n} &= \sum_{l=1}^n s^l \frac{\theta^l |S_n^l|}{\theta_{(n)}} = \frac{(\theta s)_{(n)}}{\theta_{(n)}} = \frac{\theta s(\theta s + 1)\cdots(\theta s + n - 1)}{\theta(\theta + 1)\cdots(\theta + n - 1)} \\ &= s \left(\frac{1}{\theta + 1} + \frac{\theta}{\theta + 1} s \right) \cdots \left(\frac{n-1}{\theta + n - 1} + \frac{\theta}{\theta + n - 1} s \right) = \prod_{j=1}^n \mathbb{E}_S^{\xi_j} \end{aligned}$$

where the ξ_j are independent Bernoulli random variables satisfying

$$\mathbb{P}(\xi_j = 1) = 1 - \mathbb{P}(\xi_j = 0) = \frac{\theta}{\theta + j - 1}, \quad j = 1, \dots, n. \quad (3.6.2)$$

It follows that we can write

$$K_n = \xi_1 + \cdots + \xi_n, \quad (3.6.3)$$

a sum of independent, but not identically distributed, Bernoulli random variables. Therefore

$$\mathbb{E}(K_n) = \sum_{j=1}^n \mathbb{E}\xi_j = \sum_{j=0}^{n-1} \frac{\theta}{\theta + j}, \quad (3.6.4)$$

and

$$\text{Var}(K_n) = \sum_{j=1}^n \text{Var}(\xi_j) = \sum_{j=0}^{n-1} \frac{\theta}{\theta + j} - \sum_{j=0}^{n-1} \frac{\theta^2}{(\theta + j)^2} = \sum_{j=0}^{n-1} \frac{\theta j}{(\theta + j)^2}. \quad (3.6.5)$$

For large n , we see that $\mathbb{E}K_n \sim \theta \log n$ and $\text{Var}(K_n) \sim \theta \log n$. It can be shown (cf. Barbour, Holst and Janson (1992)) that the total variation distance between a sum $W = \xi_1 + \cdots + \xi_n$ of independent Bernoulli random variables ξ_i with means p_i , and a Poisson random variable P with mean $p_1 + \cdots + p_n$ satisfies

$$d_{TV}(\mathcal{L}(W), \mathcal{L}(P)) \leq \frac{p_1^2 + \cdots + p_n^2}{p_1 + \cdots + p_n}.$$

It follows from the representation (3.6.3) that there is a constant c such that

$$d_{TV}(\mathcal{L}(K_n), \mathcal{L}(P_n)) \leq \frac{c}{\log n}, \quad (3.6.6)$$

where P_n is a Poisson random variable with mean $\mathbb{E}K_n$. As a consequence,

$$\frac{K_n - \mathbb{E}K_n}{\sqrt{\text{Var}K_n}} \Rightarrow N(0, 1), \quad (3.6.7)$$

and the same result holds if the mean and variance of K_n are replaced by $\theta \log n$.

3.7 Estimating θ

In this section, we return to the question of inference about θ from the sample. We begin with an approach used by population geneticists prior to the advent of the ESF.

The sample homozygosity

It is a simple consequence of the ESF (with $n = 2$) that

$$\mathbb{P}(\text{two randomly chosen genes are identical}) = \frac{1}{1 + \theta}.$$

In a sample of size n , define for $i \neq j$

$$\delta_{ij} = \begin{cases} 1 & \text{if genes } i \text{ and } j \text{ are identical} \\ 0 & \text{otherwise} \end{cases}$$

and set

$$F_n^* = \frac{2}{n(n-1)} \sum_{i < j} \delta_{ij}.$$

We call F_n^* the *homozygosity* of the sample; it is the probability that two randomly chosen distinct members of the sample of size n have identical types. It is elementary to show that

$$\mathbb{E}(F_n^*) = \frac{1}{1 + \theta}. \tag{3.7.1}$$

The variance of F_n^* is more difficult to calculate, but it can be shown that

$$\mathbb{E}(F_n^*)^2 = \frac{1}{n(n-1)} \left(\frac{2}{1 + \theta} + \frac{8(n-2)}{(1 + \theta)(2 + \theta)} + \frac{(n-2)(n-3)(6 + \theta)}{(1 + \theta)(2 + \theta)(3 + \theta)} \right). \tag{3.7.2}$$

The results in (3.7.1) and (3.7.2) can be combined to calculate $\text{Var}(F_n^*)$. We see that as $n \rightarrow \infty$,

$$\text{Var}(F_n^*) \rightarrow \frac{2\theta}{(1 + \theta)^2(2 + \theta)(3 + \theta)}, \tag{3.7.3}$$

as found by Stewart (1976). It turns out that F_n^* converges in distribution as $n \rightarrow \infty$ to a limiting random variable F^* having variance given in (3.7.3).

If there are l types in the sample, with μ_j of type $j, j = 1, \dots, l$, then

$$F_n^* = \sum_{j=1}^l \frac{\mu_j(\mu_j - 1)}{n(n-1)}. \tag{3.7.4}$$

We note that the homozygosity is often calculated as

$$F_n = \sum_{j=1}^l \left(\frac{\mu_j}{n} \right)^2. \tag{3.7.5}$$

The difference between F_n and F_n^* is of order n^{-1} : F_n is the probability that two genes taken *with* replacement are identical, F_n^* the probability that two genes sampled *without* replacement are identical.

We have seen that $\mathbb{E}(F_n^*) = 1/(1 + \theta)$. This suggests a method of moments estimator for θ obtained by equating the observed sample homozygosity to its expectation:

$$\tilde{\theta} = \frac{1}{F_n^*} - 1$$

The right hand side of (3.7.4) shows that $\tilde{\theta}$ depends largely on the partition of the data into types, and not on the number of types. We will see that the

latter is sufficient for θ , so standard statistical theory suggests that $\tilde{\theta}$ might not be a good estimator – it is based largely on those parts of the data which are uninformative for θ . To examine this issue further, we used a coalescent simulation to generate 10,000 samples of size 100 from the infinitely-many-alleles process for different values of the target θ , and computed the estimator $\tilde{\theta}$ for each of them. In Table 1 below are some summary statistics from these simulations.

Table 1. Simulated properties of $\tilde{\theta}$ in samples of size $n = 100$

	$\theta = 0.1$	$\theta = 1.0$	$\theta = 5.0$	$\theta = 10.0$
mean	0.15	1.38	6.00	11.38
std. dev.	0.32	1.03	2.60	4.15
RMSE†	0.32	1.10	2.79	4.37
median	0.00	1.19	5.73	11.01
5th %ile	0.00	0.09	2.21	5.25
95th %ile	0.94	3.36	10.73	18.80

†RMSE: root mean square error. 10,000 replicates used.

It can be seen that the estimator $\tilde{\theta}$ is biased upwards. This might be anticipated, because

$$\mathbb{E}(\tilde{\theta}) = \mathbb{E}(1/F_n^* - 1) \geq 1/\mathbb{E}(F_n^*) - 1 = \theta,$$

the inequality following from an application of Jensen's Inequality. We note that the estimator $\tilde{\theta}$ has a non-degenerate limit as $n \rightarrow \infty$, precisely because F_n^* does. Thus $\tilde{\theta}$ is *not* a consistent estimator of θ . However, a consistent estimator can be derived by using the number of types observed in the sample, as we now show.

Estimation using the number of types in the sample

Notice from (3.5.3) and (3.6.1) that the conditional distribution of \mathbf{c} , given that $K_n = k$, does not depend on θ :

$$\begin{aligned} \mathbb{P}(\mathbf{c}|K_n = k) &= q(\mathbf{c}) / \mathbb{P}(K_n = k) \\ &= \frac{n! \theta^k}{\theta_{(n)}} \prod_{j=1}^n \binom{1}{j}^{c_j} \frac{1}{c_j!} \bigg/ \frac{\theta^k |S_n^k|}{\theta_{(n)}} \\ &= \frac{n!}{|S_n^k|} \prod_{j=1}^n \binom{1}{j}^{c_j} \frac{1}{c_j!}. \end{aligned} \tag{3.7.6}$$

It follows that K_n is a sufficient statistic for θ ; it contains all the information useful for estimating θ . The maximum likelihood estimator of θ may be found from (3.6.1). If k alleles are observed in the sample, then the log-likelihood is

$$\log L(\theta) = \log(|S_n^k|) + k \log \theta - \sum_{j=0}^{n-1} \log(\theta + j).$$

Differentiating with respect to θ shows that the maximum likelihood estimator $\hat{\theta}$ of θ may be found by solving the equation

$$k = \sum_{j=0}^{n-1} \frac{\theta}{\theta + j}. \quad (3.7.7)$$

As can be seen from (3.6.4), this is also the moment estimator of θ . The Fisher information may be calculated readily from the log-likelihood, and we find that the asymptotic variance of $\hat{\theta}$ is

$$\text{Var}(\hat{\theta}) \approx \theta^2 / \text{Var}(K_n). \quad (3.7.8)$$

Therefore $\hat{\theta}$ is consistent for θ . Indeed, asymptotically $\hat{\theta}$ has a Normal distribution with mean θ and variance $\theta / \log n$. We used the simulated data described above to assess the properties of the estimator $\hat{\theta}$. Some results are given in Table 2. It can be seen that the distribution of $\hat{\theta}$ is much more concentrated than that of $\tilde{\theta}$, and $\hat{\theta}$ seems to be somewhat less biased than $\tilde{\theta}$. Histograms comparing the two estimators appear in Figure 3.2.

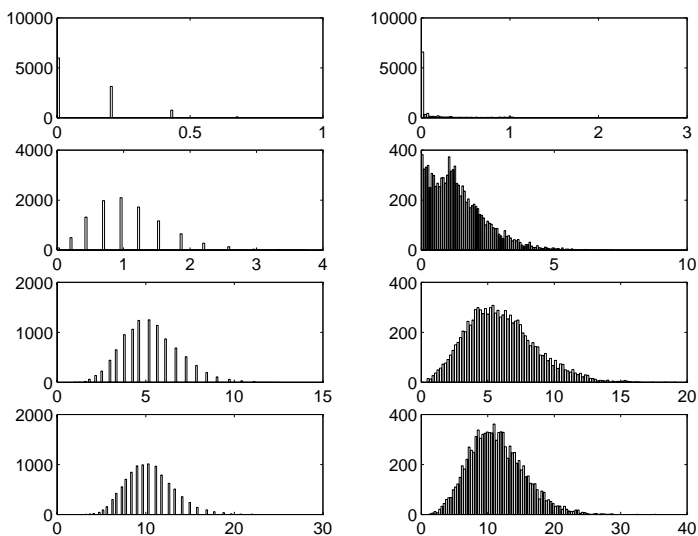
It is worth relating these two approaches to estimating θ . If we were given the values of each δ_{ij} , $1 \leq i < j \leq n$, then we would be able to calculate the value of K_n , and each of the allele frequencies. We can see that summarizing the δ_{ij} in the form of F_n^* throws away a lot of information – for example, the summary statistic results in an inconsistent estimator of θ . We shall meet this phenomenon again when we investigate estimation in the infinitely-many-sites model.

3.8 Testing for selective neutrality

One might try to perform a “goodness of fit” test on genetic data to see whether the Ewens sampling formula is appropriate. If the fit is rejected, it may be evidence of selection (or of geographical structure, variation in population sizes, other mutation mechanisms or other unnamed departures from the model). Watterson (1977) suggested using the sample homozygosity F_n defined in (3.7.5) as a test statistic. Under neutrality, the conditional distribution of the counts is given by (3.7.6), from which the null distribution of F_n follows. F_n will tend to have larger values when the allele frequencies are skewed, and smaller values when the allele frequencies are more equal. When testing for heterosis, small values of the test statistic lead to rejection

Table 2. Simulated properties of $\hat{\theta}$ in samples of size $n = 100$.

	$\theta = 0.1$	$\theta = 1.0$	$\theta = 5.0$	$\theta = 10.0$
mean	0.11	1.03	5.12	10.17
std. dev.	0.15	0.54	1.57	2.70
RMSE	0.15	1.03	1.57	2.71
median	0.00	0.95	5.14	9.70
5th %ile	0.00	0.20	2.95	6.15
95th %ile	0.43	1.87	7.82	15.00

Fig. 3.2. Histograms of 10,000 replicates of estimators of θ based on samples of size $n = 100$. Left hand column is MLE $\hat{\theta}$, right hand column is $\tilde{\theta}$. First row corresponds to $\theta = 0.1$, second to $\theta = 1.0$, third to $\theta = 5.0$, and fourth to $\theta = 10.0$.

of neutrality. For the *D. tropicalis* data in the introduction we have $F_{298} = 0.6475$, while for the *D. simulans* data we have $F_{308} = 0.2356$. Significance points of the distribution under neutrality were given in Watterson (1978), but they can be simulated rapidly. One approach, with ties to combinatorics, is outlined in the complements. Using this method, the P-value for the first set is 0.87, while for the second set it is 0.03. Thus, in contrast to Wright's expectation, the *D. simulans* do not fit neutral expectations. We will not focus further on tests of neutrality in these notes. An up-to-date discussion about detecting neutrality is given in Kreitman (2000).

4 The Coalescent

In the last two sections we studied the behavior of the genealogy of a sample from a Wright-Fisher model when the population size N is large. We introduced the ancestral process $A_n(t)$ that records the number of distinct ancestors of a sample of size n a time t earlier, and we studied some of its properties. In this section we describe in more detail the structure of Kingman's coalescent, a continuous time process whose state space is the set of equivalence relations on the set $[n] \equiv \{1, 2, \dots, n\}$. We also give an alternative representation as a bifurcating tree, and we discuss the robustness of these approximations to different models of reproduction.

4.1 Who is related to whom?

We record information not only about the number of ancestors at various times in the past, but also information about which individuals are descended from which ancestors. For some fixed time t , one way of doing this is by labelling the individuals in the sample from the set $\{1, \dots, n\}$ and defining a (random) equivalence relation \sim on $[n]$ by

$i \sim j$ if and only if individuals i and j share a common ancestor at time t .

It is often easiest to describe the equivalence relation by listing the equivalence classes. Note that each equivalence class corresponds to a particular ancestor of the sample at time t , and that the individuals in the equivalence class are exactly those who are descended from the ancestor of the class.

More formally, we could label the individuals in the sample from the set $[n]$. If at time t there are $A_n(t) = k$ ancestors of the sample, we could list the members of the sample descended from each particular ancestor. This would give us an unordered collection $E_1 \equiv \{i_{11}, \dots, i_{1l_1}\}, E_2 \equiv \{i_{21}, \dots, i_{2l_2}\}, \dots, E_k \equiv \{i_{k1}, \dots, i_{kl_k}\}$ of sets which would partition $[n]$, *i.e.* $E_i \cap E_j = \emptyset$ $i \neq j$ and $E_1 \cup \dots \cup E_k = [n]$. We often refer to the sets E_1, \dots, E_k as *classes*, or *equivalence classes*.

Denote by $C(t)$ the (random) partition (or equivalently, equivalence relation) which is obtained from the genealogy in this way. What are the dynamics of the process $\{C(t) : t \geq 0\}$? Suppose that $C(t) = \alpha$ for some partition α with k classes (we write $|\alpha| = k$). As t increases and we go further into the past, the process will remain constant until the first occasion that two of the k individuals who are the ancestors of the classes are involved in a coalescence. When this happens, those two individuals and hence all their descendants in the two equivalence classes will share a common ancestor. The effect is to merge or coalesce the two classes corresponding to these two individuals. The rate at which this happens to a particular pair of individuals (and hence to a particular pair of classes) is 1. Note that this argument and the fact that population events happen at the points of a Poisson process ensures that the process $C(\cdot)$ is Markovian.

In summary, denote by \mathcal{E}_n the set of equivalence relations on $[n]$. The process $\{C(t) : t \geq 0\}$ is a continuous-time Markov chain on \mathcal{E}_n with

$$\begin{aligned} C(0) = \Delta &\equiv \{(i, i), i = 1, 2, \dots, n\} \\ &\equiv \{\{1\}\{2\} \dots \{n\}\}, \end{aligned}$$

the state in which “nobody is related to anyone else”, and transition rates $\{q_{\alpha\beta}, \alpha, \beta \in \mathcal{E}_n\}$ given by

$$q_{\alpha\beta} = \begin{cases} -\binom{k}{2} & \text{if } \alpha = \beta, |\alpha| = k \\ 1 & \text{if } \alpha \prec \beta \\ 0 & \text{otherwise} \end{cases} \quad (4.1.1)$$

where the notation $\alpha \prec \beta$ means that the partition β may be obtained from α by merging two of the classes in α . The observation that the sample may eventually be traced back to a single common ancestor means that almost surely

$$\begin{aligned} \lim_{t \rightarrow \infty} C(t) = \Theta &\equiv \{(i, j), i, j = 1, 2, \dots, n\} \\ &= \{\{1, 2, \dots, n\}\} \end{aligned}$$

so that everybody is related to everybody else and there is just one class.

The process $\{C(t), t \geq 0\}$ is known as the n -*coalescent*, or *coalescent*. To calculate its distribution, it is convenient to study the discrete time (embedded) jump chain $\{\mathcal{C}_k; k = n, n - 1, \dots, 1\}$ obtained by watching the continuous-time process $C(\cdot)$ only at those times when it changes state. This chain starts from $\mathcal{C}_n = \Delta$ and has transition probabilities

$$\mathbb{P}(\mathcal{C}_{k-1} = \beta | \mathcal{C}_k = \alpha) = \begin{cases} \binom{k}{2}^{-1} & \text{if } \alpha \prec \beta, |\alpha| = k \\ 0 & \text{otherwise.} \end{cases}$$

Thus $C(\cdot)$ moves through a sequence $\Delta = \mathcal{C}_n \prec \mathcal{C}_{n-1} \prec \dots \prec \mathcal{C}_1 = \Theta$, spending (independent) exponential amounts of time in each state $\mathcal{C}_k \in \mathcal{E}_n$ with respective parameters $\binom{k}{2}$, $k = n, n - 1, \dots, 2$, before being absorbed in state Θ .

Notice that in $C(\cdot)$ transition rates from a state α (and hence the time spent in α) depend on α only through $|\alpha|$, and that

$$|C(t)| = A_n(t),$$

since classes in $C(t)$ correspond to ancestors of the sample. Thus the joint distributions of $\{A_n(t); t \geq 0\}$ conditional on the sequence $\mathcal{C}_n, \dots, \mathcal{C}_1$ are just the same as its unconditional distributions. Hence $\{\mathcal{C}_k\}$ and $\{A_n(t)\}$ are independent processes. Thus

$$C(t) = \mathcal{C}_{A_n(t)}, t \geq 0$$

and

$$\begin{aligned}
\mathbb{P}(C(t) = \alpha) &= \sum_{j=1}^n \mathbb{P}(C(t) = \alpha | A_n(t) = j) \mathbb{P}(A_n(t) = j) \\
&= \sum_{j=1}^n \mathbb{P}(\mathcal{C}_j = \alpha) \mathbb{P}(A_n(t) = j) \\
&= \mathbb{P}(A_n(t) = |\alpha|) \mathbb{P}(\mathcal{C}_{|\alpha|} = \alpha).
\end{aligned}$$

The distribution of $A_n(t)$ has been given earlier. That of \mathcal{C}_j is given in the following theorem of Kingman (1982a).

Theorem 4.1 *For the jump chain of the n -coalescent,*

$$\mathbb{P}(\mathcal{C}_j = \alpha) = \frac{(n-j)!j!(j-1)!}{n!(n-1)!} \lambda_1! \cdots \lambda_j!$$

where $|\alpha| = j$ and $\lambda_1, \dots, \lambda_j$ are the sizes of the equivalence classes of α .

Proof. Use backward induction. The result is clearly true when $j = n$. Then

$$\begin{aligned}
\mathbb{P}(\mathcal{C}_{j-1} = \beta) &\equiv p_{j-1}(\beta) = \sum_{\alpha \in \mathcal{E}_n} p_j(\alpha) \mathbb{P}(\mathcal{C}_{j-1} = \beta | \mathcal{C}_j = \alpha) \\
&= \sum_{\alpha < \beta} p_j(\alpha) \frac{2}{j(j-1)}
\end{aligned}$$

Write $\lambda_1, \dots, \lambda_{j-1}$ for the sizes of the equivalence classes of β . Then those of α are $\lambda_1, \dots, \lambda_{l-1}, m, \lambda_l - m, \lambda_{l+1}, \dots, \lambda_{j-1}$ for some l , $l = 1, \dots, j-1$ and some m , $m = 1, 2, \dots, \lambda_l - 1$. Using the inductive hypothesis, we have

$$\begin{aligned}
p_{j-1}(\beta) &= \sum_{l=1}^{j-1} \sum_{m=1}^{\lambda_l-1} \frac{2}{j(j-1)} \frac{(n-j)!j!(j-1)!}{n!(n-1)!} \\
&\quad \times \lambda_1! \cdots \lambda_{l-1}! m! (\lambda_l - m)! \lambda_{l+1}! \cdots \lambda_{j-1}! \frac{1}{2} \binom{\lambda_l}{m} \\
&= \frac{(n-j)!(j-1)!(j-2)!}{n!(n-1)!} \lambda_1! \cdots \lambda_{j-1}! \sum_{l=1}^{j-1} \sum_{m=1}^{\lambda_l-1} 1 \\
&= \frac{(n-j+1)(n-j)!(j-1)!(j-2)!}{n!(n-1)!} \lambda_1! \cdots \lambda_{j-1}!
\end{aligned}$$

as required. \square

Note that the distribution of \mathcal{C} and hence $C(\cdot)$ depends only on the sizes of the equivalence classes rather than on which individuals are in these classes.

4.2 Genealogical trees

Knowledge of a sample path of the n -coalescent, the value of $C(t)$ for all $t \geq 0$, specifies the time for which there are n distinct ancestors of the sample, which two individuals share an ancestor when the number of ancestors drops by 1, the time for which there are $n - 1$ distinct ancestors, which two ancestors share an ancestor when the number drops from $n - 1$ to $n - 2$, and so on. Eventually we have information about the times between coalescences and knowledge of which ancestors coalesce. Another, perhaps more natural, way of representing this information is as a genealogical tree. The lengths of the various branches are proportional to the times between the various events.

It is convenient to think of the n -coalescent as a random, rooted, binary tree, with lengths attached to the edges, instead of its original form as a stochastic process where values are partitions of $[n]$. The structure of the genealogical process translates easily to the random tree: the leaves of the tree represent the n sequences in the sample. The first join in the tree occurs at time T_n , and results in the joining of two randomly chosen sequences. There are now $n - 1$ nodes in the tree, and the next coalescence event occurs a time T_{n-1} later, and results in the joining of two nodes chosen at random from the $n - 1$. This structure is continued until the final two nodes are joined at the most common ancestor, at time W_n .

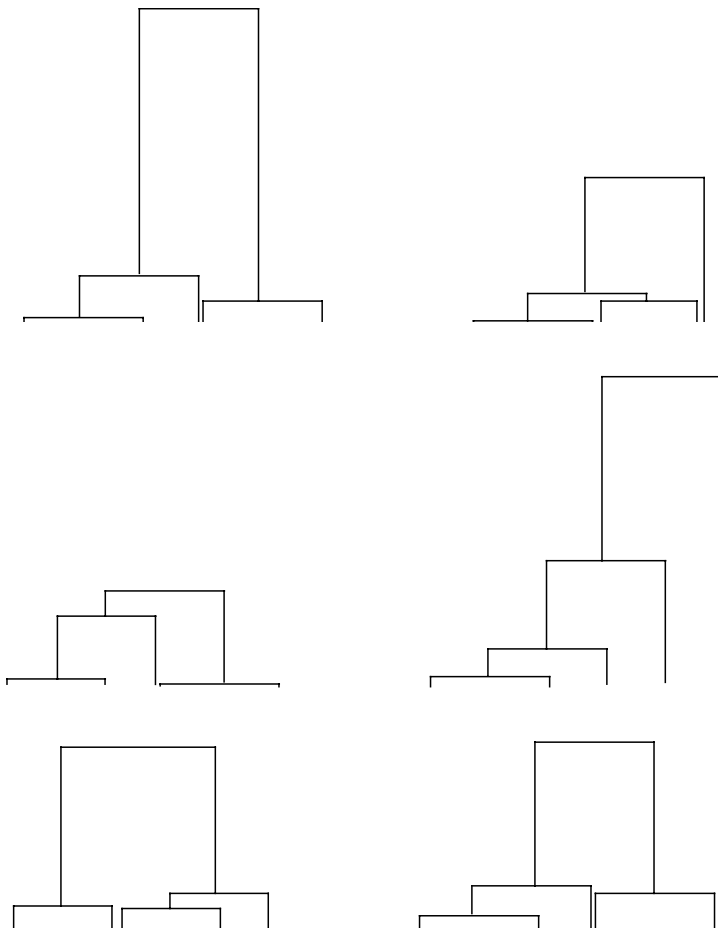
Some simulated genealogical trees for a sample of size 5 from a constant population are shown in Figure 4.1. It is instructive derive the values of the coalescent process from such a tree.

In Figure 4.2 coalescent trees for samples of size 6 and 32 from a constant size population are shown, and in Figure 4.3 trees for samples of size 6 in both constant and exponentially growing populations are shown. One of the most striking qualitative properties, which is evident in Figure 4.2, is the extent to which the tree is dominated by the last few branches. The mean time for which the tree has two branches is 1. The mean time for which the tree has more than two branches, namely $(1 - 2/n)$, is smaller: for much of the time since its common ancestor, the sample has only two ancestors. Further, for any sample size n , the variability in the time T_2 for which there are two branches in the tree accounts for most of the variability in the depth of the whole tree. These observations reinforce the theoretical results given earlier in Section 2.3. The simulated tree with exponential growth in Figure 4.3 clearly displays the star-like nature of the tree alluded to in Section 2.4.

4.3 Robustness in the coalescent

We have seen that the genealogy of the Wright-Fisher model can be described by the coalescent when the population size is large. In this section, we outline how the coalescent arises as an approximation for a wide variety of other reproduction models having constant population size.

Fig. 4.1. Six realizations, drawn on the same scale, of coalescent trees for a sample of $n = 5$. (In each tree the labels 1,2,3,4,5 should be assigned at random to the leaves.)



We noted earlier that in the Wright-Fisher model individuals have independent Poisson-distributed numbers of offspring, conditioned on the requirement that the total population size be fixed at N . Let ν_i be the number of offspring born to individual i , $i = 1, 2, \dots, N$. We saw in (2.2.1) that $\boldsymbol{\nu} = (\nu_1, \dots, \nu_N)$ has a multinomial distribution:

$$\mathbb{P}(\nu_1 = m_1, \dots, \nu_N = m_N) = \frac{N!}{m_1! \cdots m_N!} \left(\frac{1}{N} \right)^N$$

provided $m_1 + \cdots + m_N = N$. In particular the ν_i are identically distributed (but not of course independent), and

Fig. 4.2. Coalescent trees for samples of size 6 and 32 from a population of constant size

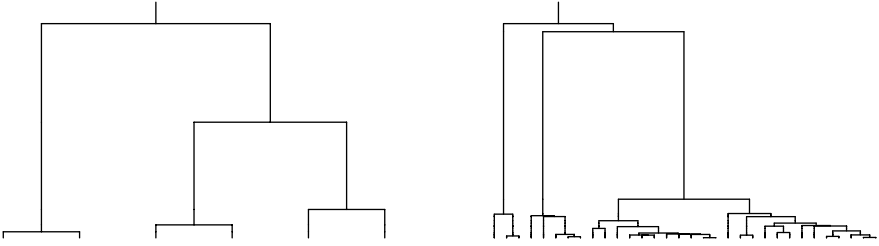
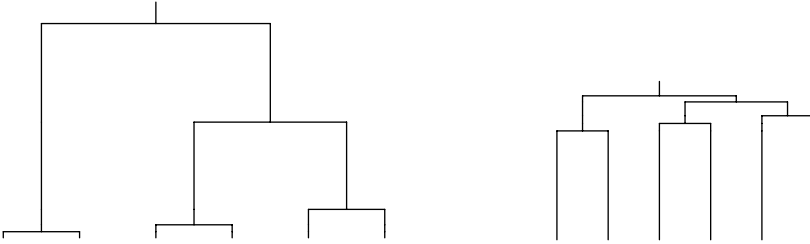


Fig. 4.3. The coalescent tree of a sample of size 6 (constant population size in left panel, exponentially growing population in right panel)



$$\mathbb{E}(\nu_1) = 1, \sigma_N^2 \equiv \text{Var}(\nu_1) = 1 - \frac{1}{N}. \tag{4.3.1}$$

Next we consider two other reproduction models that capture some of the features of the Wright-Fisher case. Suppose first that $\boldsymbol{\nu} \equiv (1, 1, \dots, 1)$, so that each individual has precisely one offspring. For this model,

$$\mathbb{E}(\nu_1) = 1, \sigma_N^2 = 0.$$

Now consider the opposite extreme, in which precisely one individual has all the offspring. Then $\boldsymbol{\nu} = N\mathbf{e}_i = N(0, \dots, 1, 0, \dots, 0)$ for some $i = 1, \dots, N$. For this case,

$$\mathbb{E}(\nu_1) = 1, \sigma_N^2 = N - 1. \tag{4.3.2}$$

Our interest focuses on the asymptotic behavior of the genealogy as $N \rightarrow \infty$. In the second model the individuals in the sample never share common ancestors, and in the third the sample can be traced back to a single individual in one generation. Clearly neither of these models has an interesting genealogy! We shall see that the way to distinguish the three models can be based on the behavior of σ_N^2 : for the Wright-Fisher model, $\sigma_N^2 \rightarrow 1$, for the second model $\sigma_N^2 = 0$, and for the third model $\sigma_N^2 \rightarrow \infty$. If time is to be rescaled in units proportional to N , then we get a non-degenerate genealogy if $\sigma_N^2 \rightarrow \sigma^2 \in (0, \infty)$.

General reproduction models with reproductive symmetry, introduced by Cannings (1974), can be formulated as follows.

- (i) Constant population size requires that $\nu_1 + \dots + \nu_N = N$.
- (ii) The collection of random variables ν_1, \dots, ν_N is exchangeable. That is, the distribution of offspring numbers does not depend on the way in which the individuals are labelled.
- (iii) The distribution of (ν_1, \dots, ν_N) is the same in each generation. This is time stationarity.
- (iv) The joint distribution of (ν_1, \dots, ν_N) is independent of family sizes in other generations. This is neutrality: offspring numbers for particular individuals do not depend on ancestral offspring numbers.

Some properties of this general model are elementary to obtain. For example, since

$$\nu_1 + \dots + \nu_N = N \tag{4.3.3}$$

and the ν_i have identical distributions it follows that

$$\mathbb{E}(\nu_1) = 1.$$

Squaring (4.3.3) and taking expectations shows that

$$\text{Cov}(\nu_1, \nu_2) = \frac{-\sigma_N^2}{N-1}. \tag{4.3.4}$$

Any particular distribution for (ν_1, \dots, ν_N) which satisfies the conditions above specifies a model for the reproduction of the population. The main result is that under minor additional conditions, the n -coalescent provides a good approximation of the genealogy of such a model when the population size is large, and time is measured in units proportional to N generations.

We begin by studying the ancestral process in a sample of size n from a population model of size N . The analog of (2.2.3) is given in the next lemma; cf. Cannings (1974) and Gladstien (1978).

Lemma 4.2 *For $1 \leq k \leq n$, we have*

$$g_{kj} = \binom{N}{k}^{-1} \binom{N}{j} \sum_{\mathbf{b} \in \Delta_j^k} \mathbb{E} \binom{\nu_1}{b_1} \dots \binom{\nu_j}{b_j} \tag{4.3.5}$$

where

$$\Delta_j^k = \{(l_1, \dots, l_j) : l_i \in \mathbb{N}, i = 1, \dots, j; l_1 + \dots + l_j = k\}.$$

Proof. Conditional on the offspring numbers $\nu = (\nu_1, \dots, \nu_N)$ we have

$$\mathbb{P}(k \text{ have } j \text{ distinct parents} | \nu) = \sum_{\substack{l_1, \dots, l_j \\ \text{distinct} \in [N]}} \sum_{\mathbf{b} \in \Delta_j^k} \binom{N}{k}^{-1} \prod_{m=1}^j \binom{\nu_{l_m}}{b_m}.$$

Taking expectations and using the exchangeability assumption completes the proof. \square

Kingman’s celebrated result gives conditions under which the genealogy of a sample is approximated by the coalescent. He showed (Kingman (1982b)) that if

- (i) $\sigma_N^2 \equiv \text{Var}(\nu_1) \rightarrow \sigma^2 \in (0, \infty)$ as $N \rightarrow \infty$;
- (ii) $\sup_N \mathbb{E}(\nu_1^k) < \infty \quad k = 3, 4, \dots$

and time is measured in units of $N\sigma^{-2}$ generations, then in the limit as $N \rightarrow \infty$, the genealogical structure of the sample is well approximated by the coalescent. Thus any result which follows from the fact that sample genealogies are described by an n -coalescent will be approximately true for any large population evolving according to an exchangeable model. The assumption of large population is reasonable in many genetics applications.

Note that the variance of the offspring distribution plays a role in the approximation of genealogy by the coalescent. If time is scaled in units of N generations, then the ancestral process appropriate for the sample is given by $A_n(\sigma^2 t), t \geq 0$. On this time scale, the waiting time T_j while the sample has j distinct ancestors has an exponential distribution with mean

$$\mathbb{E}T_j = \frac{2}{\sigma^2 j(j-1)}$$

in coalescent units, or

$$\frac{2N}{\sigma^2 j(j-1)}$$

in units of generations. It should be clear that when inferring properties of the ancestral tree from data, the parameter σ^2 has to be estimated.

Remark. As noted in Kingman (2000)), his attempt to understand the structure of the Ewens sampling formula led directly to his development of the coalescent. Kingman (1982c) derives the Ewens Sampling Formula in (3.5.3) directly from the effects of mutation in the coalescent. Define a relation $\mathcal{R} \in \mathcal{E}_n$ which contains (i, j) if, on watching the equivalence classes of $C(t)$ containing i and j until the time they coincide, we observe no mutations to either. Kingman gives the distribution of \mathcal{R} as

$$\mathbb{P}(\mathcal{R} = \xi) = \frac{\theta^k}{\theta_{(n)}} \prod_{j=1}^k (\lambda_j - 1)!, \tag{4.3.6}$$

where $\lambda_1, \dots, \lambda_k$ are the sizes of the equivalence classes of \mathcal{R} . If we multiply this by the number of ξ that have the given sizes, namely

$$\frac{n!}{\lambda_1! \cdots \lambda_k! c_1! \cdots c_n!},$$

where c_j is the number of the λ_i equal to j , we obtain the ESF. Thus the ESF is indeed a consequence of mutation in the coalescent.

4.4 Generalizations

Since the introduction of Kingman's coalescent several authors have studied related approximations. For populations of constant size, Möhle (1998) has phrased the approximations in terms of the parameter

$$c_N = \frac{\text{Var}(\nu_1)}{N-1},$$

which is the probability that two individuals chosen at random without replacement from the same generation have the same parent; cf. (4.3.4). The natural time scale is then in units of $\lfloor c_N^{-1} \rfloor$ generations.

We assume in what follows that $c_N > 0$ for sufficiently large N , and that, for integers $k_1 \geq \dots \geq k_j \geq 2$ the limits

$$\phi_j(k_1, \dots, k_j) = \lim_{N \rightarrow \infty} \frac{\mathbb{E}((\nu_1)_{[k_1]} \cdots (\nu_j)_{[k_j]})}{N^{k_1 + \dots + k_j - j} c_N} \quad (4.4.1)$$

exist, and that

$$c = \lim_{N \rightarrow \infty} c_N \quad (4.4.2)$$

exists.

A complete classification of the limiting behavior of the finite population coalescent process (run on the new time scale) is given by Möhle and Sagitov (2001). In the case

$$c = 0, \quad \phi_j(k_1, \dots, k_j) = 0 \text{ for } j \geq 2$$

the limiting process is Kingman's coalescent described earlier.

More generally, when $c = 0$ the limiting process is a continuous time Markov chain on the space of equivalence relations \mathcal{E}_n , with transition rates given by

$$q_{\alpha\beta} = \begin{cases} \phi_a(b_1, \dots, b_a) & \text{if } \alpha \subseteq \beta, \\ 0 & \text{otherwise} \end{cases} \quad (4.4.3)$$

In (4.4.3), a is the number of equivalence classes in α , $b_1 \geq b_2 \geq \dots \geq b_a$ are the ordered sizes of the groups of merging equivalence classes of β , and b is the number of equivalence classes of β . Note that $\phi_1(2) = 1$, so this does indeed reduce to the transition rates in (4.1.1) in the Kingman case. For rates

of convergence of such approximations see Möhle (2000), and for analogous results in the case of variable population size see Möhle (2002).

When $c > 0$, the limit process is a discrete time Markov chain on \mathcal{E}_n , with transition matrix P given by $P = I + cQ$, where Q has entries given in (4.4.3). This case obtains, for example, when some of the family sizes are of order N with positive probability. In these limits many groups of individuals can coalesce at the same time, and the resulting coalescent tree need not be bifurcating. Examples of this type arise when a small number of individuals has a high chance of producing most of the offspring, as is the case in some fish populations. For related material, see also Pitman (1999), Sagitov (1999) and Schweinsberg (2000).

4.5 Coalescent reviews

Coalescents have been devised for numerous other population genetics settings, most importantly to include recombination (Hudson (1983)), a subject we return to later in the notes. There have been numerous reviews of aspects of coalescent theory over the years, including Hudson (1991, 1992), Ewens (1990), Tavaré (1993), Donnelly and Tavaré (1995), Fu and Li (1999), Li and Fu (1999) and Neuhauser and Tavaré (2001). Nordborg (2001) has the most comprehensive review of the structure of the coalescent that includes selfing, substructure, migration, selection and much more.

5 The Infinitely-many-sites Model

We begin this section by introducing a data set that will motivate the developments that follow. The data are part of a more extensive mitochondrial data set obtained by Ward *et al.* (1991). Table 3 describes the segregating sites (those nucleotide positions that are not identical in all individuals in the sample) in a collection of sequences of length 360 base pairs sampled from the D-loop of 55 members of the Nuu Chah Nulth native American Indian tribe. The data exhibit a number of important features. First, each segregating site is either *purine* (A, G) or *pyrimidine* (C, T); no transversions are observed in the data. Thus at each segregating site one of two possible nucleotides is present. The segregating sites are divided into 5 purine sites and 13 pyrimidine sites. The right-most column in the table gives the multiplicity of each distinct allele (here we call each distinct sequence an allele). Notice that some alleles, such as e and j , appear frequently whereas others, such as c and n appear only once. We would like to explore the nature of the mutation process that gave rise to these data, to estimate relevant genetic parameters and to uncover any signal the data might contain concerning the demographic history of the sample. Along the way, we introduce several aspects of the theory of the infinitely-many-sites model.

The mutations represented on a tree

In our example, there are $n = 14$ distinct sequences, and each column consists of two possible characters, labelled 0 and 1 for simplicity. In order to summarize these data, we compute the numbers $II(i, j)$ giving the number of coordinates at which the i th and j th of the n sequences differ. $II(i, j)$ is the Hamming distance between sequences i and j . This results in a symmetric $n \times n$ matrix II with 0 down the diagonal. For our example, the off-diagonal elements of II are given in Table 4

It is known (cf. Buneman (1971), Waterman (1995) Chapter 14, Gusfield (1997) Chapter 17) that if an $n \times s$ data matrix representing n sequences each of k binary characters, satisfies the *four-point condition*

$$\begin{aligned} &\text{For every pair of columns, not more than three} \\ &\text{of the patterns } 00, 01, 10, 11 \text{ occur} \end{aligned} \tag{5.0.1}$$

then there is an unrooted tree linking the n sequences in such a way that the distance from sequence i to sequence j is given by the elements of the matrix D . Our example set does indeed satisfy this condition.

If the character state 0 corresponds to the ancestral base at each site, then we can check for the existence of a rooted tree by verifying the *three-point condition*

$$\begin{aligned} &\text{For every pair of columns, not more than two} \\ &\text{of the patterns } 01, 10, 11 \text{ occur} \end{aligned} \tag{5.0.2}$$

Table 3. Segregating sites in a sample of mitochondrial sequences

Position	1 1 2 2 3	1 1 1 1 1 2 2 2 2 3 3	allele freqs.
	0 9 5 9 4	8 9 2 4 6 6 9 3 6 7 7 1 3	
	6 0 1 6 4	8 1 4 9 2 6 4 3 7 1 5 9 9	
Site	1 2 3 4 5	6 7 8 9 10 11 12 13 14 15 16 17 18	
allele			
<i>a</i>	A G G A A	T C C T C T T C T C T T C	2
<i>b</i>	A G G A A	T C C T T T T C T C T T C	2
<i>c</i>	G A G G A	C C C T C T T C C C T T T	1
<i>d</i>	G G A G A	C C C C C T T C C C T T C	3
<i>e</i>	G G G A A	T C C T C T T C T C T T C	19
<i>f</i>	G G G A G	T C C T C T T C T C T T C	1
<i>g</i>	G G G G A	C C C T C C C C C C T T T	1
<i>h</i>	G G G G A	C C C T C C C T C C T T T	1
<i>i</i>	G G G G A	C C C T C T T C C C C C T	4
<i>j</i>	G G G G A	C C C T C T T C C C C T T	8
<i>k</i>	G G G G A	C C C T C T T C C C T T C	5
<i>l</i>	G G G G A	C C C T C T T C C C T T C	4
<i>m</i>	G G G G A	C C T T C T T C C C T T C	3
<i>n</i>	G G G G A	C T C T C T T C C T T T C	1

Mitochondrial data from Ward *et al.* (1991). Variable purine and pyrimidine positions in the control region. Position 69 corresponds to position 16,092 in the human reference sequence published by Anderson *et al.* (1981)

It is known that if the most frequent type at each site is labelled 0 (ancestral), then the unrooted tree exists if and only if the rooted tree exists. Gusfield (1991) gives a $O(ns)$ time algorithm for finding a rooted tree:

Algorithm 5.1 Algorithm to find rooted tree for binary data matrix

1. Remove duplicate columns in the data matrix.
2. Consider each column as a binary number. Sort the columns into decreasing order, with the largest in column 1.
3. Construct paths from the leaves to the root in the tree by labelling nodes by mutation column labels and reading vertices in paths from right to left where 1s occur in rows.

Table 4. Distance between sequences for the Ward data

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>	<i>i</i>	<i>j</i>	<i>k</i>	<i>l</i>	<i>m</i>	<i>n</i>
<i>a</i>														
<i>b</i>	1													
<i>c</i>	6	7												
<i>d</i>	6	7	4											
<i>e</i>	1	2	5	5										
<i>f</i>	2	3	6	6	1									
<i>g</i>	7	8	3	5	6	7								
<i>h</i>	8	9	4	6	7	8	1							
<i>i</i>	7	8	3	5	6	7	4	5						
<i>j</i>	6	7	2	4	5	6	3	4	1					
<i>k</i>	4	5	2	2	3	4	3	4	3	2				
<i>l</i>	5	6	1	3	4	5	2	3	2	1	1			
<i>m</i>	5	6	3	3	4	5	4	5	4	3	1	2		
<i>n</i>	6	7	4	4	5	6	5	6	5	4	2	3	3	

Figure 5.1 shows the resulting rooted tree for the Ward data, and Figure 5.2 shows corresponding unrooted tree. Note that the distances between any two sequences in the tree is indeed given by the appropriate entry of the matrix in Table 4. We emphasize that these trees are equivalent representations of the original data matrix.

In this section we develop a stochastic model for the evolution of such trees, beginning with summary statistics such as the number of segregating sites seen in the data.

5.1 Measures of diversity in a sample

We begin our study by describing some simple measures of the amount of diversity seen in a sample of DNA sequences. For a sample of n sequences of length s base pairs, write $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{is})$ for the sequence of bases from sequence i , $1 \leq i \leq n$, and define $\Pi(i, j)$ to be the number of sites at which sequences i and j differ:

$$\Pi(i, j) = \sum_{l=1}^s \mathbb{1}(y_{il} \neq y_{jl}), \quad i \neq j. \quad (5.1.1)$$

The *nucleotide diversity* Π_n in the sample is the mean pairwise difference defined by

$$\Pi_n = \frac{1}{n(n-1)} \sum_{i \neq j} \Pi(i, j), \quad (5.1.2)$$

and the per site nucleotide diversity is defined as

Fig. 5.1. Rooted tree for the Ward data found from Gusfield's algorithm

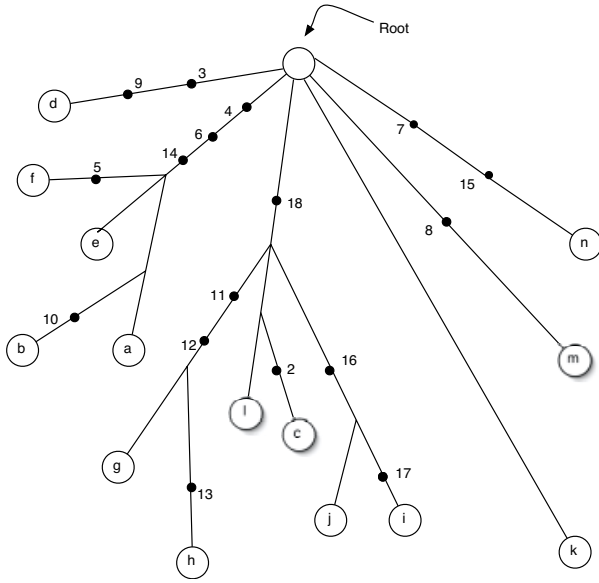
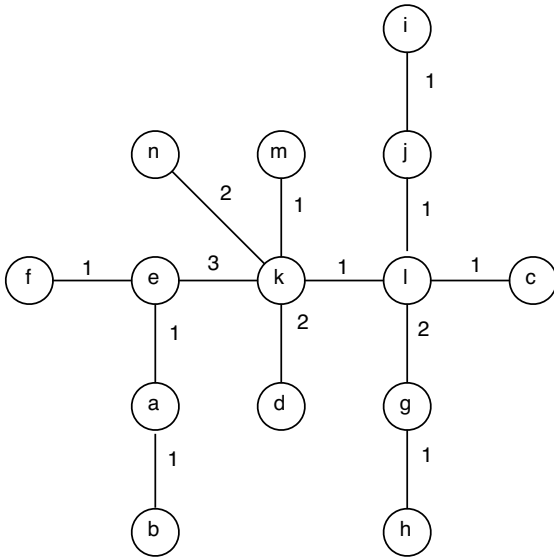


Fig. 5.2. Unrooted tree for the Ward data found from Figure 5.1. The numbers on the branches correspond to the number of sites on that branch.



$$\pi_n = \Pi_n/s.$$

Suppose that each position in the sequences being compared is from an alphabet \mathcal{A} having α different letters (so that $\alpha = 4$ in the usual nucleotide alphabet), and write n_{la} for the number of times the letter a appears in site l in the sample. Then it is straightforward to show that

$$\Pi_n = \frac{1}{n(n-1)} \sum_{l=1}^s \sum_{a \in \mathcal{A}} n_{la}(n - n_{la}) := \frac{n}{n-1} \sum_{l=1}^s H_l, \quad (5.1.3)$$

where H_l is the heterozygosity at site l , defined by

$$H_l = \sum_{a \in \mathcal{A}} \frac{n_{la}}{n} \left(1 - \frac{n_{la}}{n}\right).$$

Thus, but for the correction factor $n/(n-1)$, the per site nucleotide diversity is just the average heterozygosity across the region; that is,

$$\pi_n = \frac{n}{n-1} \frac{1}{s} \sum_{l=1}^s H_l.$$

The sampling distribution of Π_n depends of course on the mutation mechanism that operates in the region. In the case of the infinitely-many-sites mutation model, we have

$$\begin{aligned} \mathbb{E}\Pi_n &= \frac{1}{n(n-1)} \sum_{i \neq j} \Pi(i, j) \\ &= \mathbb{E}\Pi(1, 2) \quad (\text{by symmetry}) \\ &= \mathbb{E}(\# \text{ of segregating sites in sample of size } 2) \\ &= \theta \mathbb{E}(T_2), \end{aligned}$$

where T_2 is the time taken to find the MRCA of a sample of size two. In the case of constant population size, we have

$$\mathbb{E}\Pi_n = \theta. \quad (5.1.4)$$

The variance of Π_n was found by Tajima (1983), who showed that

$$\text{Var}(\Pi_n) = \frac{n+1}{3(n-1)}\theta + \frac{2(n^2+n+3)}{9n(n-1)}\theta^2. \quad (5.1.5)$$

The nucleotide diversity statistic is a rather crude summary of the variability in the data. In the next section, we study pairwise difference curves.

5.2 Pairwise difference curves

The random variables $\Pi(i, j)$ are identically distributed, but they are of course not independent. Their common distribution can be found from the observation, exploited several times already, that

$$\mathbb{P}(\Pi(1, 2) = k) = \mathbb{E}\mathbb{P}(\Pi(1, 2) = k | T_2),$$

Conditional on T_2 , $\Pi(1, 2)$ has a Poisson distribution with parameter $2T_2\theta/2 = \theta T_2$, so that for a population varying with rate function $\lambda(t)$,

$$\mathbb{P}(\Pi(1, 2) = k) = \int_0^\infty e^{-\theta t} \frac{(\theta t)^k}{k!} \lambda(t) e^{-\Lambda(t)} dt. \quad (5.2.1)$$

In the case of a constant size, when $\lambda(t) = 1$ and $\Lambda(t) = t$, the integral can be evaluated explicitly, giving

$$\mathbb{P}(\Pi(1, 2) = k) = \frac{1}{1 + \theta} \left(\frac{\theta}{1 + \theta} \right)^k, \quad k = 0, 1, \dots \quad (5.2.2)$$

Thus $\Pi(1, 2)$ has a geometric distribution with mean θ .

The *pairwise difference curve* is obtained by using the empirical distribution of the set $\Pi(i, j), 1 \leq i \neq j \leq n$ to estimate the probabilities in (5.2.1). Define

$$H_{nk} = \frac{1}{n(n-1)} \sum_{i \neq j} \mathbb{1}(\Pi(i, j) = k), \quad (5.2.3)$$

the fraction of pairs of sequences separated by k segregating sites. By symmetry, we have

$$\mathbb{E}(H_{nk}) = \mathbb{P}(\Pi(1, 2) = k), \quad k = 0, 1, \dots \quad (5.2.4)$$

5.3 The number of segregating sites

The basic properties of the infinitely-many-sites model were found by Watterson (1975). Because each mutation is assumed to produce a new segregating site, the number of segregating sites observed in a sample is just the total number of mutations S_n since the MRCA of the sample. Conditional on L_n , S_n has a Poisson distribution with mean $\theta L_n/2$. We say that S_n has a *mixed Poisson distribution*, written $S_n \sim \text{Po}(\theta L_n/2)$. It follows that

$$\begin{aligned} \mathbb{E}(S_n) &= \mathbb{E}(\mathbb{E}(S_n | L_n)) \\ &= \mathbb{E}(\theta L_n/2) \\ &= \frac{\theta}{2} \sum_{j=2}^n j \frac{2}{j(j-1)} = \theta \sum_{j=1}^{n-1} \frac{1}{j}. \end{aligned} \quad (5.3.1)$$

Notice that for large n , $\mathbb{E}(S_n) \sim \theta \log(n)$.

We can write $S_n = Y_2 + \dots + Y_n$ where Y_j is the number of mutations that arise while the sample has j ancestors. Since the T_j are independent, the Y_j are also independent. As above, Y_j has a mixed Poisson distribution, $\text{Po}(\theta j T_j / 2)$. It follows that

$$\begin{aligned} \mathbb{E}(s^{Y_j}) &= \mathbb{E}(\mathbb{E}(s^{Y_j} | T_j)) \\ &= \mathbb{E}(\exp(-[\theta j T_j / 2](1-s))) \\ &= \frac{j-1}{j-1 + \theta(1-s)}, \end{aligned} \quad (5.3.2)$$

showing (Watterson (1975)) that Y_j has a geometric distribution with parameter $(j-1)/(j-1 + \theta)$:

$$\mathbb{P}(Y_j = k) = \left(\frac{\theta}{\theta + j - 1} \right)^k \left(\frac{j-1}{\theta + j - 1} \right) \quad k = 0, 1, \dots \quad (5.3.3)$$

Since the Y_j are independent for different j , it follows that

$$\text{Var}(S_n) = \sum_{j=2}^n \text{Var}(Y_j) = \theta \sum_{j=1}^{n-1} \frac{1}{j} + \theta^2 \sum_{j=1}^{n-1} \frac{1}{j^2}. \quad (5.3.4)$$

The probability generating function of S_n satisfies

$$\mathbb{E}(s^{S_n}) = \prod_{j=2}^n \mathbb{E}(s^{Y_j}) = \prod_{j=2}^n \frac{j-1}{j-1 + \theta(1-s)} \quad (5.3.5)$$

from which further properties may be found. In particular, it follows from this that for $m = 0, 1, \dots$

$$\mathbb{P}(S_n = m) = \frac{n-1}{\theta} \sum_{l=1}^{n-1} (-1)^{l-1} \binom{n-2}{l-1} \left(\frac{\theta}{l+\theta} \right)^{m+1}. \quad (5.3.6)$$

Estimating θ

It follows from (5.3.1) that

$$\theta_W = S_n \bigg/ \sum_{j=1}^{n-1} \frac{1}{j} \quad (5.3.7)$$

is an unbiased estimator of θ . From (5.3.4) we see that the variance of θ_W is

$$\text{Var}(\theta_W) = \left[\theta \sum_{j=1}^{n-1} \frac{1}{j} + \theta^2 \sum_{j=1}^{n-1} \frac{1}{j^2} \right] \left[\sum_{j=1}^{n-1} \frac{1}{j} \right]^{-2}. \quad (5.3.8)$$

Notice that as $n \rightarrow \infty$, $\text{Var}(\theta_W) \rightarrow 0$, so that the estimator θ_W is weakly consistent for θ .

An alternative estimator of θ is the moment estimator derived from (5.1.4), namely

$$\theta_T = \Pi_n. \tag{5.3.9}$$

The variance of θ_T follows immediately from (5.1.5). In fact, Π_n has a non-degenerate limit distribution as $n \rightarrow \infty$, so that θ_T cannot be consistent. This parallels the discussion in Section 3 about estimating θ on the basis of the number K_n of alleles or via the sample homozygosity F_n . The inconsistency of the pairwise estimators arises because these summary statistics lose a lot of information available in the sample.

We used the coalescent simulation algorithm to assess the properties of the estimators θ_W and θ_T for samples of size $n = 100$. The results of 10,000 simulations are given in Tables 5 and 6 for a variety of values of θ . It can be seen that the distribution of θ_W is much more concentrated than that of θ_T . Histograms comparing the two estimators appear in Figure 5.3.

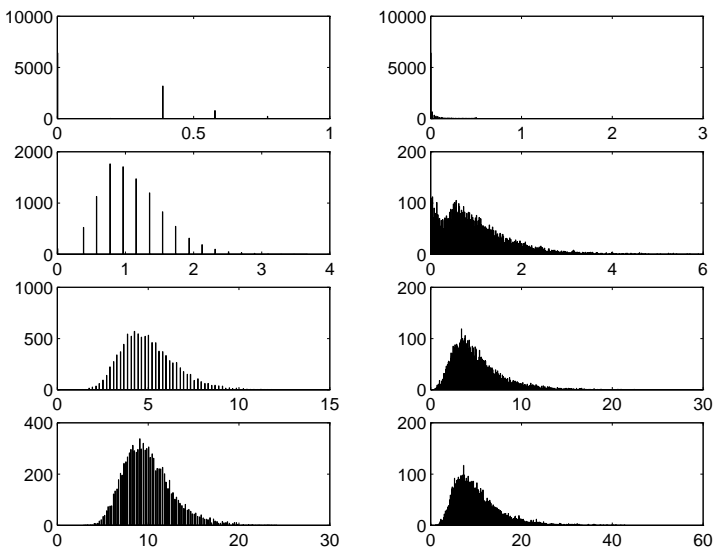
Table 5. Simulated properties of θ_W in samples of size $n = 100$.

	$\theta = 0.1$	$\theta = 1.0$	$\theta = 5.0$	$\theta = 10.0$
mean	0.18	1.10	5.03	9.99
std dev	0.23	0.48	1.53	2.75
median	0.00	0.97	4.83	9.66
5th %ile	0.00	0.39	2.90	6.18
95th %ile	0.39	1.93	7.73	15.07

Table 6. Simulated properties of θ_T in samples of size $n = 100$.

	$\theta = 0.1$	$\theta = 1.0$	$\theta = 5.0$	$\theta = 10.0$
mean	0.10	1.00	4.95	9.97
std dev	0.19	0.75	2.65	4.98
median	0.00	0.84	4.35	8.91
5th %ile	0.00	0.08	1.79	4.13
95th %ile	0.40	2.42	10.16	19.48

Fig. 5.3. Histograms of 10,000 replicates of estimators of θ based on samples of size $n = 100$. Left hand column is θ_W , right hand column is θ_T . First row corresponds to $\theta = 0.1$, second to $\theta = 1.0$, third to $\theta = 5.0$, and fourth to $\theta = 10.0$.



How well can we do?

The estimators θ_W and θ_T are based on summary statistics of the original sequence data. It is of interest to know how well these unbiased estimators might in principle behave. In this section, we examine this question in more detail for the case of constant population size.

If we knew how many mutations had occurred on each of the j branches of length T_j , $j = 2, \dots, n$ in the coalescent tree, then we could construct a simple estimator of θ using standard results for independent random variables. Let Y_{jk} , $k = 1, \dots, j$; $j = 2, \dots, n$ denote the number of mutations on the k^{th} branch of length T_j and set $Y_j = \sum_{k=1}^j Y_{jk}$. Y_j is the observed number of mutations that occur during the time the sample has j distinct ancestors. Since each mutation produces a new segregating site, this is just the number of segregating sites that arise during this time. Since the T_j are independent, so too are the Y_j . We have already met the distribution of Y_j in equation (5.3.3), and it follows that the likelihood for observations Y_j , $j = 2, \dots, n$ is

$$\begin{aligned} L_n(\theta) &= \prod_{j=2}^n \left(\frac{\theta}{j-1+\theta} \right)^{Y_j} \binom{j-1}{j-1+\theta} \\ &= \theta^{S_n} (n-1)! \prod_{j=2}^n (j-1+\theta)^{-(Y_j+1)}, \end{aligned}$$

where $S_n = \sum_{j=2}^n Y_j$ is the number of segregating sites. The maximum likelihood estimator based on this approach is therefore the solution of the equation

$$\theta = S_n \left/ \sum_{j=2}^n \frac{Y_j + 1}{j-1+\theta} \right. \tag{5.3.10}$$

Furthermore,

$$\frac{\partial^2 \log L_n}{\partial \theta^2} = -\frac{S_n}{\theta^2} + \sum_{j=2}^n \frac{(Y_j + 1)}{(j-1+\theta)^2},$$

so that

$$\begin{aligned} -\mathbb{E} \left(\frac{\partial^2 \log L_n}{\partial \theta^2} \right) &= \frac{\theta \sum_1^{n-1} \frac{1}{j}}{\theta^2} - \sum_{j=2}^n \left(\frac{\theta}{j-1} + 1 \right) \frac{1}{(j-1+\theta)^2} \\ &= \frac{1}{\theta} \sum_1^{n-1} \frac{1}{j} - \sum_1^{n-1} \frac{1}{j(j+\theta)} \\ &= \frac{1}{\theta} \sum_1^{n-1} \frac{1}{j+\theta} \end{aligned} \tag{5.3.11}$$

Hence the variance of unbiased estimators θ_U of θ satisfies

$$\text{Var}(\theta_U) \geq \theta \left/ \sum_1^{n-1} \frac{1}{j+\theta} \right.,$$

as shown by Fu and Li (1993). The right-hand side is also the large-sample variance of the estimator θ_F in (5.3.10).

How does this bound compare with that in (5.3.8)? Certainly

$$\text{Var}(\theta_F) \leq \text{Var}(\theta_W), \tag{5.3.12}$$

and we can see that if θ is fixed and $n \rightarrow \infty$ then

$$\frac{\text{Var}(\theta_F)}{\text{Var}(\theta_W)} \rightarrow 1.$$

If, on the other hand, n is fixed and θ is large, we see that

$$\frac{\text{Var}(\theta_F)}{\text{Var}(\theta_W)} \rightarrow \left(\sum_1^{n-1} \frac{1}{j} \right)^2 \left/ (n-1) \sum_1^{n-1} \frac{1}{j^2} \right.,$$

so that there can be a marked decrease in efficiency in using the estimator θ_W when θ is large. We cannot, of course, determine the numbers Y_j from data; this is more information than we have in practice. However, it does suggest that we explore the MLE of θ using the likelihoods formed from the full data rather than summary statistics. Addressing this issue leads us to study the underlying tree structure of infinitely-many-sites data in more detail, as well as to develop some computational algorithms for computing MLEs.

5.4 The infinitely-many-sites model and the coalescent

The infinitely-many-sites model is an early attempt to model the evolution of a completely linked sequence of sites in a DNA sequence. The term ‘completely linked’ means that no recombination is allowed. Each mutation on the coalescent tree of the sample introduces a mutant base at a site that has not previously experienced a mutation. One formal description treats the type of an individual as an element (x_1, x_2, \dots) of $E = \cup_{r \geq 1} [0, 1]^r$. If a mutation occurs in an offspring of an individual of type (x_1, x_2, \dots, x_r) , then the offspring has type $(x_1, x_2, \dots, x_r, U)$, where U is a uniformly distributed random variable independent of the past history of the process.

Figure 3.1 provides a trajectory of the process. It results in a sample of five sequences, their types being (U_1, U_2) , (U_1, U_2) , (U_1, U_2, U_4, U_5) , (U_0, U_3) , (U_0, U_3) respectively.

There are several other ways to represent such sequences, of which we mention just one. Consider the example above once more. Each sequence gives a mutational path from the individual back to the most recent common ancestor of the sample. We can think of these as labels of locations at which new mutant sites have been introduced. In this sample there are six such sites, each resulting in a new segregating site. We can therefore represent the sequences as strings of 0s and 1s, each of length six. At each location, a 1 denotes a mutant type and a 0 the original or ‘wild’ type. Arbitrarily labelling the sites 1, 2, \dots , 6 corresponding to the mutations at U_0, U_1, \dots, U_5 , we can write the five sample sequences as

$$\begin{aligned} (U_1, U_2, U_4, U_5) &= 011011 \\ (U_1, U_2) &= 011000 \\ (U_1, U_2) &= 011000 \\ (U_0, U_3) &= 100100 \\ (U_0, U_3) &= 100100 \end{aligned}$$

These now look more like aligned DNA sequences! Of course, in reality we do not know which type at a given segregating site is ancestral and which is mutant, and the ordering of sites by time of mutation is also unknown.

5.5 The tree structure of the infinitely-many-sites model

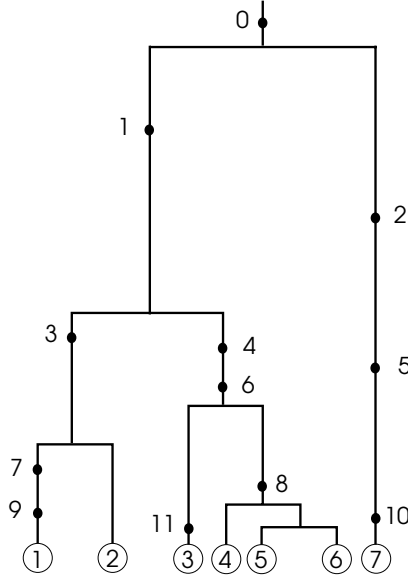
We have just seen that in the infinitely-many-sites model, each gene can be thought of as an infinite sequence of completely linked sites, each labelled 0 or 1. A 0 denotes the ancestral (original) type, and a 1 the mutant type. The mutation mechanism is such that a mutant offspring gets a mutation at a single new site that has never before seen a mutation. This changes the 0 to a 1 at that site, and introduces another segregating site into the sample. By way of example, a sample of 7 sequences might have the following structure:

```

gene 1 ... 1 0 1 0 0 0 1 0 1 0 0 ...
gene 2 ... 1 0 1 0 0 0 0 0 0 0 0 ...
gene 3 ... 1 0 0 1 0 1 0 0 0 0 1 ...
gene 4 ... 1 0 0 1 0 1 0 1 0 0 0 ...
gene 5 ... 1 0 0 1 0 1 0 1 0 0 0 ...
gene 6 ... 1 0 0 1 0 1 0 1 0 0 0 ...
gene 7 ... 0 1 0 0 1 0 0 0 0 1 0 ...
    
```

the dots indicating non-segregating sites. Many different coalescent trees can give rise to a given set of sequences. Figure 5.4 shows one of them.

Fig. 5.4. Coalescent tree with mutations



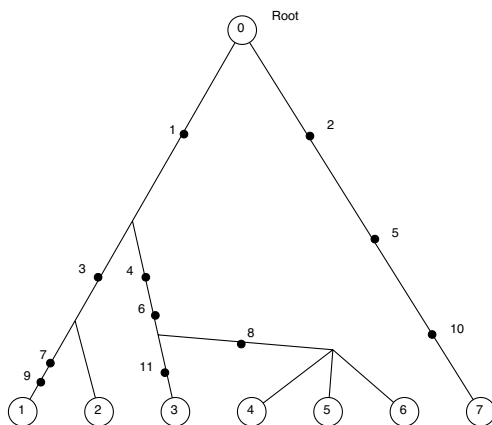
The coalescent tree with mutations can be condensed into a genealogical tree with no time scale by labelling each sequence by a list of mutations up

to the common ancestor. For the example in Figure 5.4, the sequences may be represented as follows:

gene 1 (9,7,3,1,0)
 gene 2 (3,1,0)
 gene 3 (11,6,4,1,0)
 gene 4 (8,6,4,1,0)
 gene 5 (8,6,4,1,0)
 gene 6 (8,6,4,1,0)
 gene 7 (10,5,2,0)

The condensed genealogical tree is shown in Figure 5.5. The leaves in the tree

Fig. 5.5. Genealogical tree corresponding to Figure 5.4



are the tips, corresponding to the sequences in the sample. The branches in the tree are the internal links between different mutations. The 0s in each sequence are used to indicate that the sequences can be traced back to a common ancestor.

Thus we have three ways to represent the sequences in the sample: (i) as a list of paths from the sequence to the root; (ii) as a *rooted* genealogical tree; and (iii) as a matrix with entries in $\{0, 1\}$ where a 0 corresponds to the ancestral type at a site, and a 1 the mutant type. In our example, the 0-1 matrix given above is equivalent to the representations in Figures 5.4 and 5.5. Finally, the number of segregating sites is precisely the number of mutations in the tree. In the next section, we discuss the structure of these tree representations in more detail.

5.6 Rooted genealogical trees

Following Ethier and Griffiths (1987), we think of the i th gene in the sample as a sequence $\mathbf{x}_i = (x_{i0}, x_{i1}, \dots)$ where each $x_{ij} \in \mathbb{Z}_+$. (In our earlier parlance, the type space E of a gene is the space \mathbb{Z}_+^∞ .) It is convenient to think of x_{i0}, x_{i1}, \dots as representing the most recently mutated site, the next most recently, and so on. A sample of n genes may therefore be represented as n sequences $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. The assumption that members of the sample have an ancestral tree and that mutations never occur at sites that have previously mutated imply that the sequences $\mathbf{x}_1, \dots, \mathbf{x}_n$ satisfy:

- (1) Coordinates within each sequence are distinct
- (2) If for some $i, i' \in \{1, \dots, n\}$ and $j, j' \in \mathbb{Z}_+$ we have $x_{ij} = x_{i'j'}$, then $x_{i,j+k} = x_{i',j'+k}$, $k = 1, 2, \dots$
- (3) there is a coordinate common to all n sequences.

Rules (2) and (3) above say that the part of the sequences inherited from the common ancestor appears at the right-hand end of the sequences. In practice we can discard from each \mathbf{x} sequence those entries that are common to all of the sequences in the sample; these are the coordinates after the value common to all the sequences. It is the segregating sites, and not the non-segregating sites, that are important to us. In what follows, we use these representations interchangeably.

Trees are called *labelled* if the sequences (tips) are labelled. Two labelled trees are identical if there is a renumbering of the sites that makes the labelled trees the same. More formally, let $\mathcal{T}_n = \{(\mathbf{x}_1, \dots, \mathbf{x}_n) \text{ is a tree}\}$. Define an equivalence relation \sim by writing $(\mathbf{x}_1, \dots, \mathbf{x}_n) \sim (\mathbf{y}_1, \dots, \mathbf{y}_n)$ if there is a bijection $\xi : \mathbb{Z}_+ \rightarrow \mathbb{Z}_+$ with $y_{ij} = \xi(x_{ij})$, $i = 1, \dots, n$, $j = 0, 1, \dots$. Then \mathcal{T}_n / \sim corresponds to labelled trees. Usually, we do not distinguish between an equivalence class and a typical member.

An *ordered labelled* tree is one where the sequences are labelled, and considered to be in a particular order. Visually this corresponds to a tree diagram with ordered leaves. An *unlabelled* (and so unordered) tree is a tree where the sequences are not labelled. Visually two unlabelled trees are identical if they can be drawn identically by rearranging the leaves and corresponding paths in one of the trees. Define a second equivalence relation \approx by $(\mathbf{x}_1, \dots, \mathbf{x}_n) \approx (\mathbf{y}_1, \dots, \mathbf{y}_n)$ if there is a bijection $\xi : \mathbb{Z}_+ \rightarrow \mathbb{Z}_+$ and a permutation σ of $1, 2, \dots, n$ such that $y_{\sigma(i),j} = \xi(x_{ij})$, $i = 1, \dots, n$, $j = 0, 1, \dots$. Then \mathcal{T}_n / \approx corresponds to unlabelled trees.

Usually trees are unlabelled, with sequences and sites then labelled for convenience. However it is easiest to deal with ordered labelled trees in a combinatorial and probabilistic sense, then deduce results about unlabelled trees from the labelled variety. Define

$$(\mathcal{T}_d / \sim)_0 = \{T \in \mathcal{T}_d / \sim : \mathbf{x}_1, \dots, \mathbf{x}_d \text{ all distinct}\}$$

and similarly for $(\mathcal{J}_d/\approx)_0$. $T \in \cup_{d \geq 1} (\mathcal{J}_d/\sim)_0$ corresponds to the conventional graph theoretic tree, with multiple tips removed. There is a one-to-one correspondence between trees formed from the sequences and binary sequences of sites. Let $\mathbf{x}_1, \dots, \mathbf{x}_d$ be distinct sequences of sites satisfying (1), (2) and (3), and let \mathcal{J} be the incidence matrix of segregating sites. If u_1, \dots, u_k are the segregating sites (arranged in an arbitrary order) then

$$\mathcal{J}_{ij} = 1 \text{ if } u_j \in \mathbf{x}_i, \quad i = 1, \dots, d, \quad j = 1, \dots, k.$$

The sites which are not segregating do not contain information about the tree.

Deducing the tree from a set of d binary sequences is not a priori simple, because sites where mutations occur are unordered with respect to time and any permutation of the columns of \mathcal{J} produces the same tree. In addition, binary data often have unknown ancestral labelling, adding a further complication to the picture. However, these trees are equivalent to the rooted trees discussed in the introduction. It follows that we can use the three-point condition in (5.0.2) to check whether a matrix of segregating sites is consistent with this model, and if it is, we can reconstruct the tree using Gusfield's algorithm 5.1. We turn now to computing the distribution of such a rooted tree.

5.7 Rooted genealogical tree probabilities

Let $p(T, \mathbf{n})$ be the probability of obtaining the alleles $T \in (\mathcal{J}_d/\sim)_0$ with multiplicities $\mathbf{n} = (n_1, \dots, n_d)$ and let $n = \sum_1^d n_i$. This is the probability of getting a particular *ordered* sample of distinct sequences with the indicated multiplicities. Ethier and Griffiths (1987) and Griffiths (1989) established the following:

Theorem 5.1 $p(T, \mathbf{n})$ satisfies the equation

$$\begin{aligned} n(n-1+\theta)p(T, \mathbf{n}) &= \sum_{k:n_k \geq 2} n_k(n_k-1)p(T, \mathbf{n} - \mathbf{e}_k) \\ &+ \theta \sum_{\substack{k:n_k=1, \mathbf{x}_{k0} \text{ distinct,} \\ \mathbf{x}_k \neq \mathbf{x}_j \quad \forall j}} p(\mathcal{S}_k T, \mathbf{n}) \\ &+ \theta \sum_{\substack{k:n_k=1, \\ \mathbf{x}_{k0} \text{ distinct.}}} \sum_{j:\mathbf{x}_k=\mathbf{x}_j} p(\mathcal{R}_k T, \mathcal{R}_k(\mathbf{n} + \mathbf{e}_j)). \end{aligned} \tag{5.7.1}$$

In equation (5.7.1), \mathbf{e}_j is the j th unit vector, \mathcal{S} is a shift operator which deletes the first coordinate of a sequence, $\mathcal{S}_k T$ deletes the first coordinate of the k^{th} sequence of T , $\mathcal{R}_k T$ removes the k^{th} sequence of T , and ' \mathbf{x}_{k0} distinct' means that $x_{k0} \neq x_{ij}$ for all $(\mathbf{x}_1, \dots, \mathbf{x}_d)$ and $(i, j) \neq (k, 0)$. The boundary condition is $p(T_1, (1)) = 1$.

Remark. The system (5.7.1) is recursive in the quantity $\{n-1 + \text{number of vertices in } T\}$.

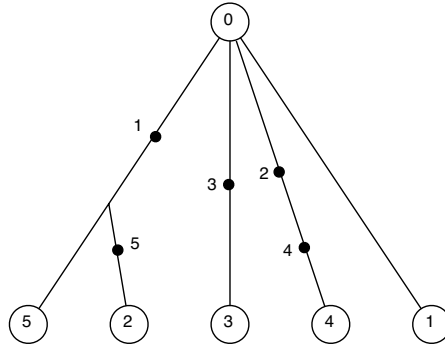


Proof. Equation (5.7.1) can be validated by a simple coalescent argument, by looking backwards in time for the first event in the ancestry of the sample. The first term on the right of (5.7.1) corresponds to a coalescence occurring first. This event has probability $(n - 1)/(\theta + n - 1)$. For any k with $n_k \geq 2$, the two individuals who coalesce may come from an allele with n_k copies, and the tree after the coalescence would be $(T, \mathbf{n} - \mathbf{e}_k)$. The contribution to (T, \mathbf{n}) form events of this sort is therefore

$$\frac{n - 1}{\theta + n - 1} \sum_{k: n_k \geq 2} \frac{n_k}{n} \left(\frac{n_k - 1}{n - 1} \right) p(T, \mathbf{n} - \mathbf{e}_k).$$

The second terms on the right of (5.7.1) correspond to events where a mutation occurs first. Suppose then that the mutation gave rise to sequence \mathbf{x}_k . There are two different cases to consider, these being determined by whether or not the sequence $\mathfrak{S}\mathbf{x}_k$ that resulted in \mathbf{x}_k is already in the sample, or not. These two cases are illustrated in the tree in Figure 5.6. The sequences are

Fig. 5.6. Representative tree



- $\mathbf{x}_1 = (0)$
- $\mathbf{x}_2 = (5 \ 1 \ 0)$
- $\mathbf{x}_3 = (3 \ 0)$
- $\mathbf{x}_4 = (2 \ 4 \ 0)$
- $\mathbf{x}_5 = (1 \ 0)$

Note that $\mathfrak{S}\mathbf{x}_2 = (1 \ 0) = \mathbf{x}_5$, so the ancestral type of \mathbf{x}_2 is in the sample. This corresponds to the third term on the right of (5.7.1). On the other hand, $\mathfrak{S}\mathbf{x}_4 = (4 \ 0)$, a type not now in the sample. This corresponds to second term on the right of (5.7.1). The phrase ‘ x_{k0} distinct’ that occurs in these two sums is

required because not all leaves with $n_k = 1$ can be removed; some cannot have arisen in the evolution of the process. The sequence \mathbf{x}_5 provides an example.

Combining these probabilities gives a contribution to $p(T, \mathbf{n})$ of

$$\frac{\theta}{\theta + n - 1} \left\{ \sum \frac{1}{n} p(\mathcal{S}_k T, \mathbf{n}) + \sum \sum \frac{1}{n} p(\mathcal{R}_k T, \mathcal{R}_k(\mathbf{n} + \mathbf{e}_j)) \right\},$$

and completes the proof. \square

It is sometimes more convenient to consider the recursion satisfied by the quantities $p^0(T, \mathbf{n})$ defined by

$$p^0(T, \mathbf{n}) = \frac{n!}{n_1! \dots n_d!} p(T, \mathbf{n}). \quad (5.7.2)$$

$p^0(T, \mathbf{n})$ is the probability of the labelled tree T , without regard to the order of the sequences in the sample. Using (5.7.1), this may be written in the form

$$\begin{aligned} n(n-1+\theta)p^0(T, \mathbf{n}) &= \sum_{k:n_k \geq 2} n(n_k-1)p^0(T, \mathbf{n} - \mathbf{e}_k) \\ &\quad + \theta \sum_{\substack{k:n_k=1, x_{k0} \text{ distinct,} \\ \mathcal{S}\mathbf{x}_k \neq \mathbf{x}_j \forall j}} p^0(\mathcal{S}_k T, \mathbf{n}) \\ &\quad + \theta \sum_{\substack{k:n_k=1, \\ x_{k0} \text{ distinct.}}} \sum_{j:\mathcal{S}\mathbf{x}_k = \mathbf{x}_j} (n_j+1)p^0(\mathcal{R}_k T, \mathcal{R}_k(\mathbf{n} + \mathbf{e}_j)). \end{aligned} \quad (5.7.3)$$

Let $p^*(T, \mathbf{n})$ be the probability of a corresponding unlabelled tree with multiplicity of the sequences given by \mathbf{n} . p^* is related to p^0 by a combinatorial factor, as follows. Let S_d denote the set of permutations of $(1, \dots, d)$. Given a tree T and $\sigma \in S_d$, define $T_\sigma = (\mathbf{x}_{\sigma(1)}, \dots, \mathbf{x}_{\sigma(d)})$ and $\mathbf{n}_\sigma = (n_{\sigma(1)}, \dots, n_{\sigma(d)})$. Letting

$$a(T, \mathbf{n}) = |\{\sigma \in S_d : T_\sigma = T, \mathbf{n}_\sigma = \mathbf{n}\}|, \quad (5.7.4)$$

we have

$$p^*(T, \mathbf{n}) = \frac{1}{a(T, \mathbf{n})} p^0(T, \mathbf{n}). \quad (5.7.5)$$

Informally, the number of distinct ordered labelled trees corresponding to the unlabelled tree is

$$\frac{n!}{n_1! \dots n_d! a(T, \mathbf{n})}.$$

In the tree shown in Figure 5.5, $a(T, \mathbf{n}) = 1$. A subsample of three genes $(9, 7, 3, 1, 0)$, $(11, 6, 4, 1, 0)$, $(10, 5, 2, 0)$, forming a tree T' with frequencies $\mathbf{n}' = (1, 1, 1)$, has $a(T', \mathbf{n}') = 2$, because the first two sequences are equivalent in an unlabelled tree.

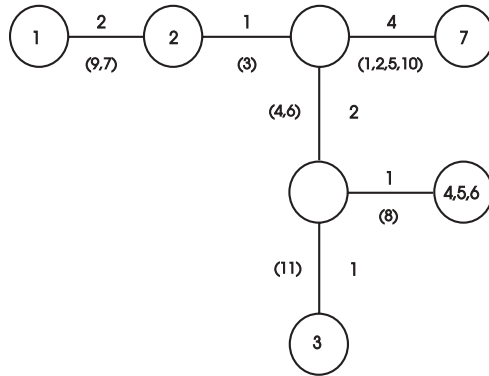
These recursions may be solved for small trees, and the resulting genealogical tree probabilities used to estimate θ by true maximum likelihood methods.

One drawback is that the method depends on knowing the ancestral type at each site, an assumption rarely met in practice. We turn now to the tree structure that underlies the process when the ancestral labelling is unknown.

5.8 Unrooted genealogical trees

When the ancestral base at each site is unknown there is an *unrooted* genealogical tree that corresponds to the sequences. In these unrooted trees, the vertices represent sequences and the number of mutations between sequences are represented by numbers along the edges; see Griffiths and Tavaré (1995). It is convenient to label the vertices to show the sequences they represent. The unrooted tree for the example sequences is shown in Figure 5.7.

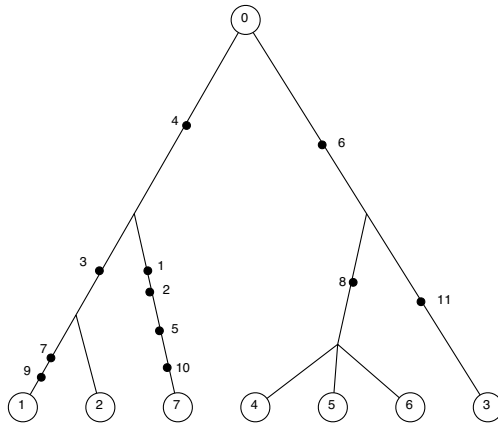
Fig. 5.7. Unrooted genealogical tree corresponding to Figure 5.4



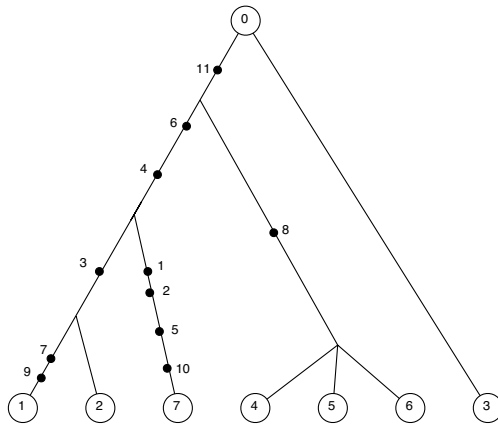
Given a single rooted tree, the unrooted genealogy can be found. The constructive way to do this is to put potential ancestral sequences at the nodes in the rooted tree (ignoring the root). There are three such nodes in the example in Figure 5.5. The ancestral sequence might be represented in the sample (as with sequence 2 in that figure), or it may be an inferred sequence not represented in the sample.

Given a rooted genealogy, we have seen how the corresponding unrooted tree can be found. Conversely, the class of rooted trees produced from an unrooted genealogy may be constructed by placing the root at one of the sequences, or between mutations along an edge. This corresponds to picking up the unrooted tree at that point and shaking it. Two examples are given in Figure 5.8. In the first, the root corresponds to the third sequence, and in the second it is between the two mutations between the two inferred sequences. The unrooted tree constructed from any of these rooted trees is of course unique.

Fig. 5.8. Moving the root



Tree with root between mutations



Tree with root the third sequence

If there are α sequences (including the inferred sequences), with m_1, m_2, \dots mutations along the edges, and s segregating sites, then there are

$$\alpha + \sum_j (m_j - 1) = s + 1 \tag{5.8.1}$$

rooted trees when the sequences are labelled. There may be fewer unlabelled rooted trees, as some can be identical after unlabelling the sequences. In the example there are 11 segregating sites, and so 12 labelled rooted trees, which correspond to distinct unlabelled rooted trees as well.

The class of rooted trees corresponds to those constructed from toggling the ancestor labels 0 and 1 at sites. The number of the 2^s possible relabellings that are consistent with the sequences having come from a tree is

$$\alpha + \sum_j \sum_{k=1}^{m_j-1} \binom{m_j}{k} = \alpha + \sum_j (2^{m_j} - 2). \tag{5.8.2}$$

This follows from the observation that if there is a collection of m segregating sites which correspond to mutations between sequences, then the corresponding data columns of the 0-1 sequences (with 0 the ancestral state) are identical or complementary. Any of the $\binom{m}{k}$ configurations of k identical and $m-k$ complementary columns correspond to the same labelled tree with a root placed after the k th mutation. The correspondence between different rooted labelled trees and the matrix of segregating sites can be described as follows: in order to move the root from one position to another, toggle those sites that occur on the branches between the two roots.

The upper tree in Figure 5.8 has incidence matrix

```
gene 1 0 0 1 1 0 0 1 0 1 0 0
gene 2 0 0 1 1 0 0 0 0 0 0 0
gene 3 0 0 0 0 0 1 0 0 0 0 1
gene 4 0 0 0 0 0 1 0 1 0 0 0
gene 5 0 0 0 0 0 1 0 1 0 0 0
gene 6 0 0 0 0 0 1 0 1 0 0 0
gene 7 1 1 0 1 1 0 0 0 0 1 0
```

whereas the lower tree in Figure 5.8 has incidence matrix

```
gene 1 0 0 1 1 0 1 1 0 1 0 1
gene 2 0 0 1 1 0 1 0 0 0 0 1
gene 3 0 0 0 0 0 0 0 0 0 0 0
gene 4 0 0 0 0 0 0 0 1 0 0 1
gene 5 0 0 0 0 0 0 0 1 0 0 1
gene 6 0 0 0 0 0 0 0 1 0 0 1
gene 7 1 1 0 1 1 1 0 0 0 1 1
```

It can readily be checked that the sites between the two roots are those numbered 6 and 11, and if these are toggled then one tree is converted into the other.

5.9 Unrooted genealogical tree probabilities

A labelled unrooted genealogical tree of a sample of sequences has a vertex set V which corresponds to the labels of the sample sequences and any inferred sequences in the tree. Let \mathcal{Q} be the edges of the tree, described by $(m_{ij}, i, j \in V)$, where m_{ij} is the number of mutations between vertices i and

j . Let \mathbf{n} denote the multiplicities of the sequences. It is convenient to include the inferred sequences $\ell \in V$ with $n_\ell = 0$. Then the unrooted genealogy is described by (\mathbf{Q}, \mathbf{n}) .

Define $p(\mathbf{Q}, \mathbf{n})$, $p^0(\mathbf{Q}, \mathbf{n})$, $p^*(\mathbf{Q}, \mathbf{n})$ analogously to the probabilities for T . The combinatorial factor relating $p^*(\mathbf{Q}, \mathbf{n})$ and $p^0(\mathbf{Q}, \mathbf{n})$ is

$$a(\mathbf{Q}, \mathbf{n}) = |\{\sigma \in S_{|V|} : \mathbf{Q}_\sigma = \mathbf{Q}, \mathbf{n}_\sigma = \mathbf{n}\}|. \quad (5.9.1)$$

The quantities $p(\mathbf{Q}, \mathbf{n})$ and $p^0(\mathbf{Q}, \mathbf{n})$ satisfy recursions similar to (5.7.1) and (5.7.3), which can be derived by considering whether the last event back in time was a coalescence or a mutation. The recursion for $p(\mathbf{Q}, \mathbf{n})$ is

$$\begin{aligned} n(n-1+\theta)p(\mathbf{Q}, \mathbf{n}) &= \sum_{k:n_k \geq 2} n_k(n_k-1)p(\mathbf{Q}, \mathbf{n} - \mathbf{e}_k) \\ &+ \theta \sum_{\substack{k:n_k=1, |k|=1, \\ k \rightarrow j, m_{kj} > 1}} p(\mathbf{Q} - \mathbf{e}_{kj}, \mathbf{n}) \\ &+ \theta \sum_{\substack{k:n_k=1, |k|=1, \\ k \rightarrow j, m_{kj}=1}} p(\mathbf{Q} - \mathbf{e}_{kj}, \mathbf{n} + \mathbf{e}_j - \mathbf{e}_k), \end{aligned} \quad (5.9.2)$$

where $|k| = 1$ means that the degree of the vertex k is 1 (that is, k is a leaf), and $k \rightarrow j$ means that vertex k is joined to vertex j . In the last term on the right of (5.9.2), vertex k is removed from \mathbf{Q} . The boundary conditions in (5.9.2) for $n = 2$ are

$$p((0), 2\mathbf{e}_1) = \frac{1}{1+\theta},$$

and

$$p((m), \mathbf{e}_1 + \mathbf{e}_2) = \left(\frac{\theta}{1+\theta}\right)^m \frac{1}{1+\theta}, \quad m = 1, 2, \dots$$

The probability of a labelled unrooted genealogical tree \mathbf{Q} is

$$p(\mathbf{Q}, \mathbf{n}) = \sum_{T \in C(\mathbf{Q})} p(T, \mathbf{n}), \quad (5.9.3)$$

where $C(\mathbf{Q})$ is the class of distinct labelled rooted trees constructed from \mathbf{Q} . The same relationship holds in (5.9.3) if p is replaced by p^0 .

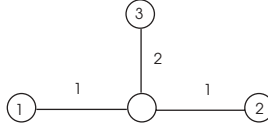
5.10 A numerical example

In this example we suppose that the ancestral states are unknown, and that the sequences, each with multiplicity unity, are:

$$\begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{array}$$

For convenience, label the segregating sites 1, 2, 3, and 4 from the left. When 0 is the ancestral state, a possible rooted tree for these sequences has paths to the root of (1, 0), (2, 3, 0), and (4, 0). It is then straightforward to construct the corresponding unrooted genealogy, which is shown in Figure 5.9. The central sequence is inferred. There are five possible labelled rooted trees

Fig. 5.9. Unrooted Genealogy



constructed from the unrooted genealogy, corresponding to the root being at one of the sequences, or between the two mutations on the edge. These five trees are shown in Figure 5.10, together with their probabilities $p(T, \mathbf{n})$, computed exactly from the recursion (5.7.1) when $\theta = 2.0$. $p(\mathbf{Q}, \mathbf{n})$ is the sum of these probabilities, 0.004973. The factor in (5.9.1) is 2, and the multinomial coefficient $3!/1!1!1! = 6$ so $p^*(\mathbf{Q}, \mathbf{n}) = 3 \times 0.00497256 = 0.014919$. Note that the trees (b) and (e) are identical unlabelled rooted trees, but are distinct labelled rooted trees, so are both counted in calculating $p^*(\mathbf{Q}, \mathbf{n})$.

In this small genealogy, the coalescent trees with four mutations can be enumerated to find the probability of the genealogy. The trees which produce the tree in Figure 5.9 are shown in Figure 5.11, with the correspondence to the trees in Figure 5.10 highlighted.

Let T_3 be the time during which the sample has three ancestors, and T_2 the time during which it has two. T_3 and T_2 are independent exponential random variables with respective rates 3 and 1. By considering the Poisson nature of the mutations along the edges of the coalescent tree, the probability of each type of tree can be calculated. For example, the probability $p_{(a1)}$ of the first tree labelled (a1) is

$$\begin{aligned}
 p_{(a1)} &= \mathbb{E} \left[\left(e^{-\theta T_3/2} \frac{\theta T_3}{2} \right)^2 e^{-\theta T_2/2} e^{-\theta(T_2+T_3)/2} \frac{1}{2!} (\theta(T_2 + T_3)/2)^2 \right] \\
 &= \frac{\theta^4}{32} \mathbb{E} \left[e^{-\theta(3T_3/2+T_2)} T_3^2 (T_2 + T_3)^2 \right] \\
 &= \frac{\theta^4(17\theta^2 + 46\theta + 32)}{27(\theta + 1)^3(\theta + 2)^5}.
 \end{aligned}$$

In a similar way the other tree probabilities may be calculated. We obtain

Fig. 5.10. Labelled rooted tree probabilities

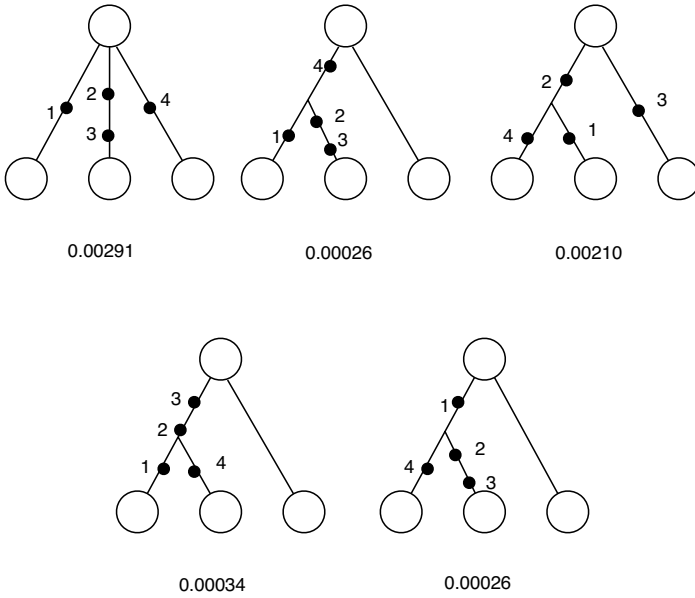
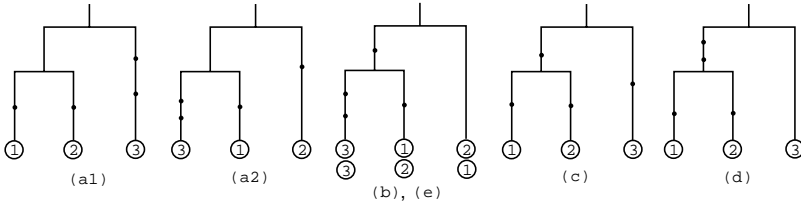


Fig. 5.11. Possible coalescent trees leading to the trees in Figure 5.10



$$\begin{aligned}
 p_{(a2)} &= \frac{\theta^4}{16} \mathbb{E} \left[2e^{-\theta(3T_3/2+T_2)} T_3^3 (T_2 + T_3) / 2 \right] \\
 &= \frac{2\theta^4(11\theta + 14)}{27(\theta + 1)^2(\theta + 2)^5}, \\
 p_{(b)} = p_{(e)} &= \frac{\theta^4}{16} \mathbb{E} \left[e^{-\theta(3T_3/2+T_2)} T_3^3 T_2 / 2 \right] \\
 &= \frac{\theta^4}{9(\theta + 1)^2(\theta + 2)^4}, \\
 p_{(c)} &= \frac{\theta^4}{16} \mathbb{E} \left[e^{-\theta(3T_3/2+T_2)} (T_2 + T_3) T_3^2 T_2 \right] \\
 &= \frac{\theta^4(2\theta + 3)}{9(\theta + 1)^3(\theta + 2)^4}, \\
 p_{(d)} &= \frac{\theta^4}{16} \mathbb{E} \left[e^{-\theta(3T_3/2+T_2)} T_3^2 T_2^2 / 2 \right] \\
 &= \frac{2\theta^4}{9(\theta + 1)^3(\theta + 2)^3}.
 \end{aligned}$$

Note that there are two coalescent trees that correspond to case (a2), depending on whether 1 coalesced with 3 first, or 2 did. When $\theta = 2$, these probabilities reduce to $p_{(a1)} = 0.004115, p_{(a2)} = 0.004630, p_{(b),(e)} = 0.000772, p_{(c)} = 0.003601, p_{(d)} = 0.001029$. From these we deduce that $p(T(a), \mathbf{n}) = (0.004115 + 0.004630)/3 = 0.002915, p(T(b), \mathbf{n}) = p(T(e), \mathbf{n}) = 0.000772/3 = 0.000257, p(T(c), \mathbf{n}) = 0.003601/3 = 0.001203, \text{ and } p(T(d), \mathbf{n}) = 0.001029/3 = 0.000343$, so that $p(\mathbf{Q}, \mathbf{n}) = 0.004973$, in agreement with the recursive solution.

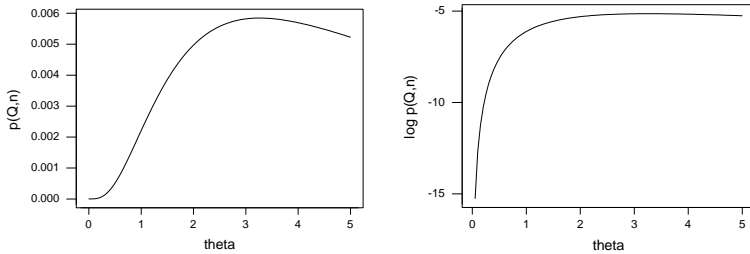
5.11 Maximum likelihood estimation

For the example in the previous section, it can be shown that the likelihood is

$$p(\mathbf{Q}, \mathbf{n}) = \frac{4\theta^4(5\theta^2 + 14\theta + 10)}{27(\theta + 1)^3(\theta + 2)^5}.$$

This has the value 0.004973 when $\theta = 2$, as we found above. The maximum likelihood estimator of θ is $\hat{\theta} = 3.265$, and the approximate variance (found from the second derivative of the log-likelihood) is 8.24. The likelihood curves are plotted in Figure 5.12.

Fig. 5.12. Likelihood $p(\mathbf{Q}, \mathbf{n})$ plotted as a function of θ , together with log-likelihood.



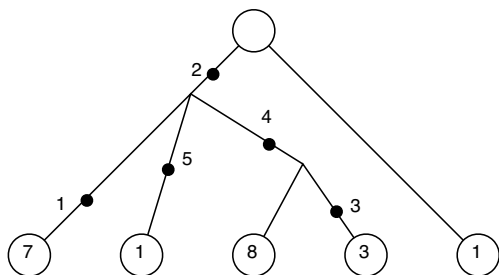
As might be expected, there is little information in such a small sample. Now consider a data set with 20 sequences, 5 segregating sites and multiplicities given below. The reduced genealogical tree is given in Figure 5.13.

```

0 1 0 1 0 : 8
0 1 1 1 0 : 3
0 0 0 0 0 : 1
0 1 0 0 1 : 1
1 1 0 0 0 : 7
    
```

Assuming that the ancestral labels are known, the probabilities $p^*(T, \mathbf{n})$ may be found using the recursion in (5.7.1), and they give a value of the MLE as $\hat{\theta} = 1.40$.

Fig. 5.13. Rooted genealogical tree for example data set. [Here, leaf labels refer to multiplicities of sequences]



To develop a practical method of maximum likelihood we need to be able to solve the recursions for p^0 for large sample sizes and large numbers of segregating sites. A general method for doing this is discussed in the next section.

6 Estimation in the Infinitely-many-sites Model

In this section we describe some likelihood methods for the infinitely-many-sites model, with a view to estimation of the compound mutation parameter θ . The method described here originated with Griffiths and Tavaré (1994), and has since been revisited by Felsenstein *et al.* (1999) and Stephens and Donnelly (2000). As we saw at the end of the previous section, exact calculation using the recursion approach is possible for relatively small sample sizes. For larger samples a different approach is required. We begin this section with Monte Carlo-based method for approximating these sampling probabilities by simulation backwards along the sample paths of the coalescent. Later in the section we relate this approach to importance sampling and show how to improve the original approach.

6.1 Computing likelihoods

Griffiths and Tavaré's approach is based on an elementary result about Markov chains given below.

Lemma 6.1 *Let $\{X_k; k \geq 0\}$ be a Markov chain with state space S and transition matrix P . Let A be a set of states for which the hitting time*

$$\eta = \inf\{k \geq 0 : X_k \in A\}$$

is finite with probability one starting from any state $x \in T \equiv S \setminus A$. Let $f \geq 0$ be a function on S , and define

$$u_x(f) = \mathbb{E}_x \prod_{k=0}^{\eta} f(X_k) \tag{6.1.1}$$

for all $X_0 = x \in S$, so that

$$u_x(f) = f(x), x \in A$$

Then for all $x \in T$

$$u_x(f) = f(x) \sum_{y \in S} p_{xy} u_y(f). \tag{6.1.2}$$

Proof.

$$\begin{aligned}
u_x(f) &= \mathbb{E}_x \left(\prod_{k=0}^{\eta} f(X_k) \right) \\
&= f(x) \mathbb{E}_x \left(\prod_{k=1}^{\eta} f(X_k) \right) \\
&= f(x) \mathbb{E}_x \left(\mathbb{E}_x \left(\prod_{k=1}^{\eta} f(X_k) \right) \middle| X_1 \right) \\
&= f(x) \mathbb{E}_x \left(\mathbb{E}_{X_1} \left(\prod_{k=0}^{\eta} f(X_k) \right) \right) \text{ (by the Markov property)} \\
&= f(x) \mathbb{E}_x u(X_1) \\
&= f(x) \sum_{y \in S} p_{xy} u_y(f).
\end{aligned}$$

□

This result immediately suggests a simulation method for solving equations like that on the right of (6.1.2): simulate a trajectory of the chain X starting at x until it hits A at time η , compute the value of the product $\prod_{k=0}^{\eta} f(X_k)$, and repeat this several times. Averaging these values provides an estimate of $u_x(f)$.

One application of this method is calculation of the sample tree probabilities $p^0(T, \mathbf{n})$ for the infinitely-many-sites model using the recursion in (5.7.3). In this case the appropriate Markov chain $\{X_k, k \geq 0\}$ has a tree state space, and makes transitions as follows:

$$(T, \mathbf{n}) \rightarrow (T, \mathbf{n} - \mathbf{e}_k) \text{ with probability } \frac{(n_k - 1)}{f(T, \mathbf{n})(n + \theta - 1)} \quad (6.1.3)$$

$$\rightarrow (\mathcal{S}_k T, \mathbf{n}) \text{ with probability } \frac{\theta}{f(T, \mathbf{n})n(n + \theta - 1)} \quad (6.1.4)$$

$$\rightarrow (\mathcal{R}_k T, \mathcal{R}_k(\mathbf{n} + \mathbf{e}_j)) \text{ with prob. } \frac{\theta(n_j + 1)}{f(T, \mathbf{n})n(n + \theta - 1)} \quad (6.1.5)$$

The first type of transition is only possible if $n_k > 1$, and the second or third if $n_k = 1$. In the last two transitions a distinct singleton first coordinate in a sequence is removed. The resulting sequence is still distinct from the others in (6.1.4), but in (6.1.5) the shifted k th sequence is equal to the j th sequence. The scaling factor is

$$f(T, \mathbf{n}) \equiv f_{\theta}(T, \mathbf{n}) = \sum_{k=1}^d \frac{(n_k - 1)}{(n + \theta - 1)} + \frac{\theta m}{n(n + \theta - 1)},$$

where m is given by

$$m = |\{k : n_k = 1, x_{k,0} \text{ distinct, } \mathcal{S}\mathbf{x}_k \neq \mathbf{x}_j \forall j\}| + \sum_{k:n_k=1, x_{k,0} \text{ distinct}} \sum_{j:\mathcal{S}\mathbf{x}_k=\mathbf{x}_j} (n_j + 1).$$

The idea is to run the process starting from an initial tree (T, \mathbf{n}) until the time τ at which there are two sequences (x_{10}, \dots, x_{1i}) and (x_{20}, \dots, x_{2j}) with $x_{1i} = x_{2j}$ (corresponding to the root of the tree) representing a tree T_2 . The probability of such a tree is

$$p^0(T_2) = (2 - \delta_{i+j,0}) \binom{i+j}{j} \left[\frac{\theta}{2(1+\theta)} \right]^{i+j} \frac{1}{1+\theta}.$$

The representation of $p^0(T, \mathbf{n})$ is now

$$p^0(T, \mathbf{n}) = \mathbb{E}_{(T, \mathbf{n})} \left[\prod_{l=0}^{\tau-1} f(T(l), \mathbf{n}(l)) \right] p^0(T_2), \tag{6.1.6}$$

where $X(l) \equiv (T(l), \mathbf{n}(l))$ is the tree at time l . Equation (6.1.6) may be used to produce an estimate of $p^0(T, \mathbf{n})$ by simulating independent copies of the tree process $\{X(l), l = 0, 1, \dots\}$, and computing $\left[\prod_{l=0}^{\tau-1} f(T(l), \mathbf{n}(l)) \right] p^0(T_2)$ for each run. The average over all runs is then an unbiased estimator of $p^0(T, \mathbf{n})$. An estimate of $p^*(T, \mathbf{n})$ can then be found by dividing by $a(T, \mathbf{n})$.

6.2 Simulating likelihood surfaces

The distribution $p^0(T, \mathbf{n})$ provides the likelihood of the data (T, \mathbf{n}) , and so can be exploited for maximum likelihood approaches. One way to do this is to simulate the likelihood *independently* at a grid of points, and examine the shape of the resulting curve. In practice, this can be a very time consuming approach. In this section we describe another approach, based on importance sampling, for approximating a likelihood surface at a grid of points using just one run of the simulation algorithm.

The method uses the following lemma, a generalization of Lemma 6.1. The proof is essentially the same, and is omitted.

Lemma 6.2 *Let $\{X_k; k \geq 0\}$ be a Markov chain with state space S and transition matrix P . Let A be a set of states for which the hitting time*

$$\eta \equiv \eta_A = \inf\{k \geq 0 : X_k \in A\}$$

is finite with probability one starting from any state $x \in T \equiv S \setminus A$. Let $h \geq 0$ be a given function on A , let $f \geq 0$ be a function on $S \times S$ and define

$$u_x(f) = \mathbb{E}_x h(X_\eta) \prod_{k=0}^{\eta-1} f(X_k, X_{k+1}) \tag{6.2.1}$$

for all $X_0 = x \in S$, so that

$$u_x(f) = h(x), x \in A.$$

Then for all $x \in T$

$$u_x(f) = \sum_{y \in S} f(x, y) p_{xy} u_y(f). \quad (6.2.2)$$

It is convenient to recast the required equations in a more generic form, corresponding to the notation in Lemma 6.2. We denote by $q_\theta(x)$ the probability of the data x when the unknown parameters have value θ , which might be vector-valued. Equations such as (5.7.3) can then be recast in the form

$$q_\theta(x) = \sum_y f_\theta(x, y) p_\theta(x, y) q_\theta(y) \quad (6.2.3)$$

for some appropriate transition matrix $p_\theta(x, y)$. Now suppose that θ_0 is a particular set of parameters satisfying

$$f_\theta(x) p_\theta(x, y) > 0 \Rightarrow p_{\theta_0}(x, y) > 0.$$

We can recast the equations (6.2.3) in the form

$$q_\theta(x) = \sum_y f_\theta(x, y) \frac{p_\theta(x, y)}{p_{\theta_0}(x, y)} p_{\theta_0}(x, y) q_\theta(y) \quad (6.2.4)$$

so that from Lemma 6.2

$$q_\theta(x) = \mathbb{E}_x q_\theta(X(\eta)) \prod_{j=0}^{\eta-1} f_{\theta, \theta_0}(X(j), X(j+1)) \quad (6.2.5)$$

where $\{X(k), k \geq 0\}$ is the Markov chain with parameters θ_0 and

$$f_{\theta, \theta_0}(x, y) = f_\theta(x) \frac{p_\theta(x, y)}{p_{\theta_0}(x, y)}. \quad (6.2.6)$$

It follows that $q_\theta(x)$ can be calculated from the realizations of a single Markov chain, by choosing a value of θ_0 to drive the simulations, and evaluating the functional $q(X(\eta)) \prod_{j=0}^{\eta-1} f_{\theta, \theta_0}(X(j), X(j+1))$ along the sample path for each of the different values of θ of interest.

6.3 Combining likelihoods

It is useful to use independent runs for several values of θ_0 to estimate $q_\theta(x)$ on a grid of θ -values. For each such θ , the estimates for different θ_0 have the required mean $q_\theta(x)$, but they have different variances for different θ_0 . This

raises the question about how estimated likelihoods from different runs might be combined. Suppose then that we are approximating the likelihood on a set of g grid points, $\theta_1, \dots, \theta_g$, using r values of θ_0 and t runs of each simulation. Let \hat{q}_{ij} be the sample average of the t runs at the j th grid point for the i th value of θ_0 . For large t , the vectors $\hat{\mathbf{q}}_i \equiv (\hat{q}_{i1}, \dots, \hat{q}_{ig}), i = 1, \dots, r$ have independent and approximately multivariate Normal distributions with common mean vector $(q_{\theta_1}(x), \dots, q_{\theta_g}(x))$ and variance matrices $t^{-1}\Sigma_1, \dots, t^{-1}\Sigma_r$ respectively. The matrices $\Sigma_1, \dots, \Sigma_r$ are unknown but may be estimated in the conventional way from the simulations. Define the log-likelihood estimates $\hat{\mathbf{l}}_i \equiv (\hat{l}_{ij}, j = 1, 2, \dots, g)$ by

$$\hat{l}_{ij} = \log \hat{q}_{ij}, \quad j = 1, \dots, g, \quad i = 1, \dots, r.$$

By the delta method, the vectors $\hat{\mathbf{l}}_i, i = 1, \dots, r$ are independent, asymptotically Normal random vectors with common mean vector $\mathbf{l} \equiv (l_1, \dots, l_g)$ given by

$$l_i = \log q_{\theta_i}(x),$$

and covariance matrices $t^{-1}\Sigma_i^*$ determined by

$$(\Sigma_i^*)_{lm} = \frac{(\Sigma_i)_{lm}}{q_{\theta_i}(x) q_{\theta_m}(x)}. \tag{6.3.1}$$

If the Σ_j^* were assumed known, the minimum variance unbiased estimator of \mathbf{l} would be

$$\hat{\mathbf{l}} = \left(\sum_{j=1}^r (\Sigma_j^*)^{-1} \right)^{-1} \sum_{j=1}^r (\Sigma_j^*)^{-1} \hat{\mathbf{q}}'_j. \tag{6.3.2}$$

If the observations for different θ_j are not too correlated, it is useful to consider the simpler estimator with $\Sigma'_j \equiv \text{diag } \Sigma_j^*$ replacing Σ_j^* in (6.3.2). This estimator requires a lot less computing than that in (6.3.2). In practice, we use the estimated values \hat{q}_{il} and \hat{q}_{im} from the i th run to estimate the terms in the denominator of (6.3.1).

6.4 Unrooted tree probabilities

The importance sampling approach can be used to find the likelihood of an unrooted genealogy. However it seems best to proceed by finding all the possible rooted labelled trees corresponding to an unrooted genealogy, and their individual likelihoods. Simulate the chain $\{(T(l), \mathbf{n}(l)), l = 0, 1, \dots\}$ with a particular value θ_0 as parameter, and obtain the likelihood surface for other values of θ using the representation

$$p_\theta^0(T, \mathbf{n}) = \mathbb{E}_{(T, \mathbf{n})}^{\theta_0} \left[\prod_{l=0}^{\tau-1} h((T(l), \mathbf{n}(l)), (T(l+1), \mathbf{n}(l+1))) \right] p_\theta^0(T_2), \tag{6.4.1}$$

where $(T(l), \mathbf{n}(l))$ is the tree at time l , and h is determined by

$$h((T, \mathbf{n}), (T, \mathbf{n} - \mathbf{e}_k)) = f_{\theta_0}(T, \mathbf{n}) \frac{n + \theta_0 - 1}{n + \theta - 1},$$

and

$$h((T, \mathbf{n}), (T', \mathbf{n}')) = f_{\theta_0}(T', \mathbf{n}') \frac{\theta(n + \theta_0 - 1)}{\theta_0(n + \theta - 1)}.$$

where the last form holds for both transitions (6.1.4), when $(T', \mathbf{n}') = (\mathcal{S}_k T, \mathbf{n})$, and (6.1.5), when $(T', \mathbf{n}') = (\mathcal{R}_k T, \mathcal{R}_k(\mathbf{n} + \mathbf{e}_j))$.

Example

To illustrate the method we consider the following set of 30 sequences, with multiplicities given in parentheses:

0 0 1 0 0 0 1 (3)
 0 0 0 0 0 0 1 (4)
 0 0 0 0 0 0 0 (4)
 1 0 0 1 0 0 0 (11)
 1 0 0 0 0 0 0 (1)
 0 1 0 0 0 0 0 (2)
 0 0 0 0 1 0 1 (2)
 0 0 0 0 1 1 1 (3)

Simulations of the process on a grid of θ -values $\theta = 0.6(0.2)3.0$ for $\theta_0 = 1.0, 1.8$, and 2.6 were run for 30,000 replicates each. The curves of $\log p^0$ were combined as described earlier. This composite curve is compared with the true curve, obtained by direct numerical solution of the recursion, in Figure 6.1.

6.5 Methods for variable population size models

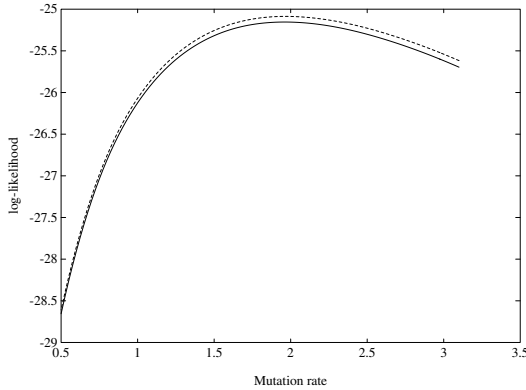
The present approach can also be used when the population size varies, as shown by Griffiths and Tavaré (1996, 1997). The appropriate recursions have a common form that may be written

$$q(t, x) = \int_t^\infty \sum_y r(s; x, y) q(s, y) g(t, x; s) ds \quad (6.5.1)$$

where $r(s; x, y) \geq 0$ and $g(t, x; s)$ is the density of the time to the first event in the ancestry of the sample after time t :

$$g(t, x; s) = \gamma(s, x) \exp\left(-\int_t^s \gamma(u, x) du\right). \quad (6.5.2)$$

Fig. 6.1. Log-likelihood curves. Dashed line: exact values. Solid line: Monte Carlo approximant.



Define

$$\begin{aligned}
 f(s; x) &= \sum_y r(s; x, y) \\
 P(s; x, y) &= \frac{r(s; x, y)}{f(s; x)},
 \end{aligned}
 \tag{6.5.3}$$

and rewrite (6.5.1) as

$$q(t, x) = \int_t^\infty f(s; x) \sum_y P(s; x, y) q(s, y) g(t, x; s) ds.
 \tag{6.5.4}$$

We associate a non-homogeneous Markov chain $\{X(t), t \geq 0\}$ with (6.5.4) as follows: Given that $X(t) = x$, the time spent in state x has density $g(t, x; s)$, and given that a change of state occurs at time s , the probability that the next state is y is $P(s; x, y)$. The process $X(\cdot)$ has a set of absorbing states, corresponding to those x for which $q(\cdot, x)$ is known. $X(\cdot)$ may be used to give a probabilistic representation of $q(t, x)$ analogous to the result in Lemma 6.1 in the following way: Let $\tau_1 < \tau_2 \cdots < \tau_k = \tau$ be the jump times of $X(\cdot)$, satisfying $\tau_0 \equiv t < \tau_1$, where τ is the time to hit the absorbing states. Then

$$q(t, x) = \mathbb{E}_{(t,x)} q(\tau, X(\tau)) \prod_{j=1}^k f(\tau_j; X(\tau_{j-1})),
 \tag{6.5.5}$$

where $\mathbb{E}_{(t,x)}$ denotes expectation with respect to $X(t) = x$.

Once more, the representation in (6.5.5) provides a means to approximate $q(x) \equiv q(0, x)$: Simulate many independent copies of the process $\{X(t), t \geq 0\}$

starting from $X(0) = x$, and compute the observed value of the functional under the expectation sign in (6.5.5) for each of them. The average of these functionals is an unbiased estimate of $q(x)$, and we may then use standard theory to see how accurately $q(x)$ has been estimated.

We have seen that it is important, particularly in the context of variance reduction, to have some flexibility in choosing the stopping time τ . Even in the varying environment setting, there are cases in which $q(\cdot, x)$ can be computed (for example by numerical integration) for a larger collection of states x , and then it is useful to choose τ to be the hitting time of this larger set.

The probability $q(t, x)$ is usually a function of some unknown parameters, which we denote once more by θ ; we write $q_\theta(t, x)$ to emphasize this dependence on θ . Importance sampling may be used as earlier to construct a single process $X(\cdot)$ with parameters θ_0 , from which estimates of $q_\theta(t, x)$ may be found for other values of θ . We have

$$q_\theta(t, x) = \int_t^\infty \sum_y f_{\theta, \theta_0}(t, x; s, y) P_{\theta_0}(s; x, y) q_\theta(s, y) g_{\theta_0}(t, x; s) ds \quad (6.5.6)$$

where

$$f_{\theta, \theta_0}(t, x; s, y) = \frac{f_\theta(s; x) g_\theta(t, x; s) P_\theta(s; x, y)}{g_{\theta_0}(t, x; s) P_{\theta_0}(s; x, y)}$$

and $f_\theta(s; x)$ and $P_\theta(s; x, y)$ are defined in (6.5.3). The representation analogous to (6.5.5) is

$$q_\theta(t, x) = \mathbb{E}_{(t, x)} q(\tau, X(\tau)) \prod_{j=1}^k f_{\theta, \theta_0}(\tau_{j-1}, X(\tau_{j-1}); \tau_j, X(\tau_j)), \quad (6.5.7)$$

and estimates of $q_\theta(t, x)$ may be simulated as described earlier in the Section.

6.6 More on simulating mutation models

The genetic variability we observe in samples of individuals is the consequence of mutation in the ancestry of these individuals. In this section, we continue the description of how mutation processes may be superimposed on the coalescent. We suppose that genetic types are labelled by elements of a set E , the ‘type space’. As mutations occur, the labels of individuals move around according to a mutation process on E .

We model mutation by supposing that a particular offspring of an individual of type $x \in E$ has a type in the set $B \subseteq E$ with probability $\Gamma(x, B)$. The mutation probabilities satisfy

$$\int_E \Gamma(x, dy) = 1, \quad \text{for all } x \in E.$$

When E is discrete, it is more usual to specify a transition matrix $\Gamma = (\gamma_{ij})$, where γ_{ij} is the probability that an offspring of an individual of type i is of type j . Such a mutation matrix Γ satisfies

$$\gamma_{ij} \geq 0, \quad \sum_{j \in E} \gamma_{ij} = 1 \text{ for each } i.$$

We assume that conditional its parent's type, the type of a particular offspring is independent of the types of other offspring, and of the demography of the population. In particular, the offspring of different individuals mutate independently.

In Section 3.4 we described a way to simulate samples from an infinitely-many-alleles model. This method can be generalized easily to any mutation mechanism. Generate the coalescent tree of the sample, sprinkle Poisson numbers of mutations on the branches at rate $\theta/2$ per branch, and superimpose the effects of the mutation process at each mutation. For discrete state spaces, this amounts to changing from type $i \in E$ to $j \in E$ with probability γ_{ij} at each mutation. This method works for variable population size, by running from the bottom up to generate the ancestral history, then from top down to add mutations.

When the population size is constant, it is possible to perform the simulation from the top down in one sweep.

Algorithm 6.1 To generate a stationary random sample of size n .

1. Choose a type at random according to the stationary distribution π of Γ . Copy this type, resulting in 2 lines.
2. If there are currently k lines, wait a random amount of time having exponential distribution with parameter $k(k + \theta - 1)/2$ and choose one of the lines at random. Split this line into 2 (each with same type as parent line) with probability $(k - 1)/(k + \theta - 1)$, and otherwise mutate the line according to Γ .
3. If there are fewer than $n + 1$ lines, return to step 2. Otherwise go back to the last time at which there were n lines and stop.

This algorithm is due to Ethier and Griffiths (1987); See also Donnelly and Kurtz (1996). Its nature comes from the 'competing exponentials' world, and it only works in the case of constant population size. For the infinitely-many-alleles and infinitely-many-sites models, the first step has to be modified so that the MRCA starts from an arbitrary label.

6.7 Importance sampling

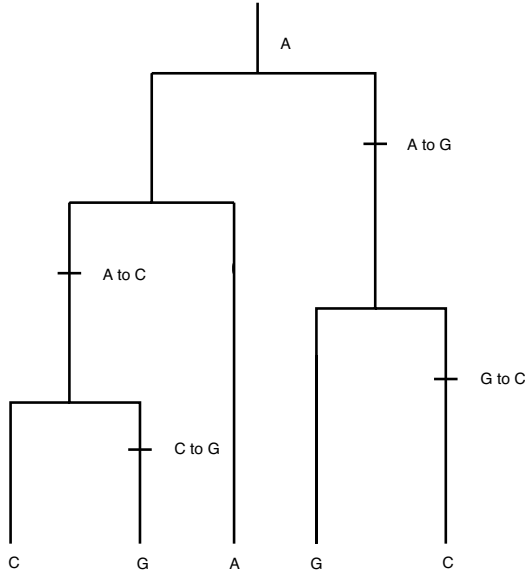
The next two sections are based on the papers of Felsenstein *et al.* (1999), and Stephens and Donnelly (2000). The review article of Stephens (2001) is also useful. In what follows, we assume a constant size population.

The *typed ancestry* \mathcal{A} of the sample is its genealogical tree G , together with the genetic type of the most recent common ancestor (MRCA) and the details and positions of the mutation events that occur along the branches of

G . An example is given in Figure 6.2. Algorithm 6.1 can be used to simulate observations having the distribution of \mathcal{A} .

The *history* \mathcal{H} is the typed ancestry \mathcal{A} with time and topology information removed. So \mathcal{H} is the type of the MRCA together with an ordered list of the split and mutation events which occur in \mathcal{A} (including the details of the types involved in each event, but not including which line is involved in each event). The history \mathcal{H} contains a record of the states $(H_{-m}, H_{-m+1}, \dots, H_{-1}, H_0)$ visited by the process beginning with the type $H_{-m} \in E$ of the MRCA and ending with genetic types $H_0 \in E^n$ of the sample. Here m is random, and the H_i are unordered lists of genetic types. Think of \mathcal{H} as $(H_{-m}, H_{-m+1}, \dots, H_{-1}, H_0)$, although it actually contains the details of which transitions occur between these states. In Figure 6.2, we have $\mathcal{H} = (\{A\}, \{A, A\}, \{A, G\}, \{A, A, G\}, \{A, C, G\}, \{A, C, G, G\}, \{A, C, C, G\}, \{A, C, C, C, G\}, \{A, C, C, G, G\})$.

Fig. 6.2. Genealogical tree G , typed ancestry \mathcal{A} and history \mathcal{H}



If H_i is obtained from H_{i-1} by a mutation from α to β , write $H_i = H_{i-1} - \alpha + \beta$, whereas if H_i is obtained from H_{i-1} by the split of a line of type α , write $H_i = H_{i-1} + \alpha$. The distribution $P_\theta(\mathcal{H})$ of \mathcal{H} is determined by the distribution π of the type of the MRCA, by the stopping rule in Algorithm 6.1, and by the Markov transition probabilities

$$\tilde{p}_\theta(H_i | H_{i-1}) = \begin{cases} \frac{n_\alpha}{n} \frac{\theta}{n-1+\theta} \Gamma_{\alpha\beta} & \text{if } H_i = H_{i-1} - \alpha + \beta \\ \frac{n_\alpha}{n} \frac{n-1}{n-1+\theta} & \text{if } H_i = H_{i-1} + \alpha \\ 0 & \text{otherwise} \end{cases} \quad (6.7.1)$$

where n_α is the number of chromosomes of type α in H_{i-1} and $n = \sum n_\alpha$.

We want to compute the distribution $q_\theta(\cdot)$ of the genetic types $\mathcal{D}_n = (a_1, \dots, a_n)$ in a random ordered sample. A sample from \mathcal{H} provides, through H_0 , a sample from q_θ . To get the ordered sample, we have to label the elements of H_0 , so that

$$q_\theta(\mathcal{D}_n | \mathcal{H}) = \begin{cases} (\prod_{\alpha \in E} n_\alpha!)/n! & \text{if } H_0 \text{ is consistent with } \mathcal{D}_n \\ 0 & \text{otherwise.} \end{cases} \quad (6.7.2)$$

We regard $L(\theta) \equiv q_\theta(\mathcal{D}_n)$ as the likelihood of the data \mathcal{D}_n . The Griffiths-Tavaré method uses the representation

$$L(\theta) = \mathbb{E} \left(\prod_{j=0}^{\tau} F(B_j) \mid B_0 = \mathcal{D}_n \right), \quad (6.7.3)$$

where B_0, B_1, \dots is a particular Markov chain and τ a stopping time for the chain; recall (6.1.6). Using (6.7.2), we can calculate

$$L(\theta) = \int q_\theta(\mathcal{D}_n | \mathcal{H}) P_\theta(\mathcal{H}) d\mathcal{H} \quad (6.7.4)$$

This immediately suggests a naive estimator of $L(\theta)$:

$$L(\theta) \approx \frac{1}{R} \sum_{i=1}^R q_\theta(\mathcal{D}_n | \mathcal{H}_i) \quad (6.7.5)$$

where $\mathcal{H}_i, i = 1, \dots, R$ are independent samples from $P_\theta(\mathcal{H})$. Unfortunately each term in the sum is with high probability equal to 0, so reliable estimation of $L(\theta)$ will require *enormous* values of R .

The importance sampling approach tries to circumvent this difficulty. Suppose that $Q_\theta(\cdot)$ is a distribution on histories that satisfies $\{\mathcal{H} : Q_\theta(\mathcal{H}) > 0\} \subset \{\mathcal{H} : P_\theta(\mathcal{H}) > 0\}$. Then we can write

$$L(\theta) = \int q_\theta(\mathcal{D}_n | \mathcal{H}) \frac{P_\theta(\mathcal{H})}{Q_\theta(\mathcal{H})} Q_\theta(\mathcal{H}) d\mathcal{H} \quad (6.7.6)$$

$$\approx \frac{1}{R} \sum_{i=1}^R q_\theta(\mathcal{D}_n | \mathcal{H}_i) \frac{P_\theta(\mathcal{H}_i)}{Q_\theta(\mathcal{H}_i)} := \frac{1}{R} \sum_{i=1}^R w_i, \quad (6.7.7)$$

where $\mathcal{H}_1, \dots, \mathcal{H}_R$ are independent samples from $Q_\theta(\cdot)$.

We call the distribution Q_θ the IS proposal distribution, and the w_i are called the IS weights. The idea of course is to choose the proposal distribution in such a way that the variance of the estimator in (6.7.7) is much smaller than that of the estimator in (6.7.5). The optimal choice Q_θ^* of Q_θ is

$$Q_\theta^*(\mathcal{H}) = P_\theta(\mathcal{H} \mid \mathcal{D}_n); \quad (6.7.8)$$

in this case

$$q_\theta(\mathcal{D}_n \mid \mathcal{H}) \frac{P_\theta(\mathcal{H})}{Q_\theta^*(\mathcal{H})} = L(\theta),$$

so the variance of the estimator is 0. Unfortunately, the required conditional distribution of histories is not known, so something else has to be tried.

In Section 6.2 we mentioned that estimating $L(\theta)$ on a grid of points can be done independently at each grid point, or perhaps by importance sampling, which in the present setting reduces to choosing the driving value θ_0 , and calculating

$$L(\theta) \approx \frac{1}{R} \sum_{i=1}^R q_\theta(\mathcal{D}_n \mid \mathcal{H}_i) \frac{P_\theta(\mathcal{H}_i)}{Q_{\theta_0}(\mathcal{H}_i)} \quad (6.7.9)$$

where $\mathcal{H}_1, \dots, \mathcal{H}_R$ are independent samples from $Q_{\theta_0}(\cdot)$.

6.8 Choosing the weights

A natural class of proposal distributions on histories arises by considering randomly reconstructing histories backward in time in a Markovian way, from the sample \mathcal{D}_n back to an MRCA. So a random history $\mathcal{H} = (H_{-m}, \dots, H_{-1}, H_0)$ may be sampled by choosing $H_0 = \mathcal{D}_n$, and successively generating H_{-1}, \dots, H_{-m} according to prespecified backward transition probabilities $p_\theta(H_{i-1} \mid H_i)$. The process stops at the first time that the configuration H_{-m} consists of a single chromosome.

In order for (6.7.6) to hold, we need to look at the subclass \mathcal{M} of these distributions for which, for each i , the support of $p_\theta(\cdot \mid H_i)$ is the set

$$\{H_{i-1} : \tilde{p}_\theta(H_i \mid H_{i-1}) > 0\}$$

where \tilde{p}_θ is given in (6.7.1). Such a p_θ then specifies a distribution Q_θ whose support is the set of histories consistent with the data \mathcal{D}_n .

Felsenstein *et al.* (1999) showed that the Griffiths-Tavaré scheme in (6.7.3) is a special case of this strategy, with

$$p_\theta(H_{i-1} \mid H_i) \propto \tilde{p}_\theta(H_{i-1} \mid H_i). \quad (6.8.1)$$

The optimal choice of Q_θ^* turns out to be from the class \mathcal{M} . Stephens and Donnelly (2000) prove the following result:

Theorem 6.3 Define $\pi(\cdot | \mathcal{D})$ to be the conditional distribution of the type of an $(n + 1)$ th sampled chromosome, given the types \mathcal{D} of the first n sampled chromosomes. Thus

$$\pi(\alpha | \mathcal{D}) = \frac{q_\theta(\{\mathcal{D}, \alpha\})}{q_\theta(\mathcal{D})}.$$

The optimal proposal distribution Q_θ^* is in the class \mathcal{M} , with

$$p_\theta^*(H_{i-1} | H_i) = \begin{cases} C^{-1} \frac{\theta}{2} n_\alpha \frac{\pi(\beta | H_i - \alpha)}{\pi(\alpha | H_i - \alpha)} \Gamma_{\beta\alpha} & \text{if } H_{i-1} = H_i - \alpha + \beta, \\ C^{-1} \binom{n_\alpha}{2} \frac{1}{\pi(\alpha | H_i - \alpha)} & \text{if } H_{i-1} = H_i - \alpha, \\ 0 & \text{otherwise,} \end{cases} \quad (6.8.2)$$

where n_α is the number of chromosomes of type α in H_i , and $C = n(n-1+\theta)/2$ where n is the number of chromosomes in H_i .

It is clear that knowing p_θ^* is equivalent to knowing Q_θ^* , which in turn is equivalent to knowing $L(\theta)$. So it should come as no surprise that the conditional probabilities are unknown for most cases of interest. The only case that is known explicitly is that in which $\Gamma_{\alpha\beta} = \Gamma_\beta$ for all α, β . In this case

$$\pi(\beta | \mathcal{D}) = \frac{n_\beta + \theta \Gamma_\beta}{n + \theta}. \quad (6.8.3)$$

Donnelly and Stephens argue that under the optimal proposal distribution there will be a tendency for mutations to occur towards the rest of the sample, and that coalescences of unlikely types are more likely than those of likely types. This motivated their choice of approximation $\hat{\pi}(\cdot | \mathcal{D})$ to the sampling probabilities $\pi(\cdot | \mathcal{D})$. They define $\hat{\pi}(\cdot | \mathcal{D})$ by choosing an individual from \mathcal{D} at random, and mutating it a geometric number of times according to the mutation matrix Γ . So

$$\hat{\pi}(\beta | \mathcal{D}) = \sum_{\alpha \in E} \frac{n_\alpha}{n} \sum_{m=0}^{\infty} \left(\frac{\theta}{\theta + n} \right)^m \frac{n}{\theta + n} \Gamma_{\alpha\beta}^m \quad (6.8.4)$$

$$\equiv \sum_{\alpha \in E} \frac{n_\alpha}{n} M_{\alpha\beta}^{(n)}. \quad (6.8.5)$$

$\hat{\pi}$ has a number of interesting properties, among them the fact that when $\Gamma_{\alpha\beta} = \Gamma_\beta$ for all α, β we have $\hat{\pi}(\cdot | \mathcal{D}) = \pi(\cdot | \mathcal{D})$ and the fact that $\hat{\pi}(\cdot | \mathcal{D}) = \pi(\cdot | \mathcal{D})$ when $n = 1$ and Γ is reversible.

The proposal distribution \hat{Q}_θ^* , an approximation to Q_θ^* , is defined by substituting $\hat{\pi}(\cdot | \mathcal{D})$ into (6.8.2):

$$\hat{p}_\theta(H_{i-1} | H_i) = \begin{cases} C^{-1} \frac{\theta}{2} n_\alpha \frac{\hat{\pi}(\beta | H_i - \alpha)}{\hat{\pi}(\alpha | H_i - \alpha)} \Gamma_{\beta\alpha} & \text{if } H_{i-1} = H_i - \alpha + \beta, \\ C^{-1} \binom{n_\alpha}{2} \frac{1}{\hat{\pi}(\alpha | H_i - \alpha)} & \text{if } H_{i-1} = H_i - \alpha, \\ 0 & \text{otherwise,} \end{cases} \quad (6.8.6)$$

In order to sample from \hat{p}_θ efficiently, one can use the following algorithm.

Algorithm 6.2

1. Choose a chromosome uniformly at random from those in H_i , and denote its type by α .
2. For each type $\beta \in E$ for which $\Gamma_{\beta\alpha} > 0$, calculate $\hat{\pi}(\beta | H_i - \alpha)$ from equation (6.8.5).
3. Sample H_i by setting

$$H_{i-1} = \begin{cases} H_i - \alpha + \beta & \text{w.p. } \propto \theta \hat{\pi}(\beta | H_i - \alpha) \Gamma_{\beta\alpha} \\ H_i - \alpha & \text{w.p. } \propto n_\alpha - 1. \end{cases}$$

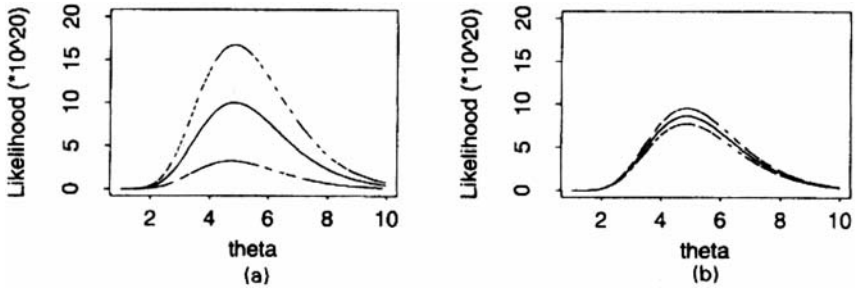
Example

Stephens and Donnelly give a number of examples of the use of their proposal distribution, including for the infinitely-many-sites model. In this case, the foregoing discussion has to be modified, because the type space E is uncountably infinite. However the principles behind the derivation of the proposal distribution \hat{Q}_θ can be used here too. Namely, we choose a chromosome uniformly at random from those present, and assume this chromosome is involved in the most recent event back in time. As we have seen (recall Theorem 5.1), the configuration of types H_i is equivalent to an unrooted genealogical tree, and the nature of mutations on that tree means that the chromosomes that can be involved in the most recent event backwards in time from H_i are limited:

- (a) any chromosome which is not the only one of its type may coalesce with another of that type;
- (b) any chromosome which is the only one of its type and has only one neighbor on the unrooted tree corresponding to H_i may have arisen from a mutation to that neighbor.

So their proposal distribution chooses the most recent event back in time by drawing a chromosome uniformly at random from those satisfying (a) or (b). Notice that this distribution does not depend on θ . In Figure 6.3 are shown a comparison of the Griffiths-Tavaré method with this new proposal distribution.

Fig. 6.3. (a) Likelihood surface estimate with ± 2 standard deviations from 100,000 runs of GT method, with $\theta_0 = 4$. (b) the same using 100,000 runs of the SD IS function. This is Fig. 7 from Stephens and Donnelly (2000).



It is an open problem to develop other, perhaps better, IS distributions for rooted and unrooted trees as well. The method presented here is also not appropriate for variable population size models, where the simple Markov structure of the process is lost. The representation of the Griffiths-Tavaré method as importance sampling, together with the results for the constant population size model, suggest that the development of much more efficient likelihood algorithms in that case. See Chapter 2 of Liu (2001) for an introduction to sequential importance sampling in this setting. The paper of Stephens and Donnelly has extensive remarks from a number of discussants on the general theme of computational estimation of likelihood surfaces.

7 Ancestral Inference in the Infinitely-many-sites Model

The methods in this section are motivated by the problem of inferring properties of the time to the most recent common ancestor of a sample given the data from that sample. For example, Dorit *et al.* (1996) sequenced a 729 bp region of the ZFY gene in a sample of $n = 38$ males and observed no variability; the number of segregating sites in the data is then $S_{38} = 0$. What can be said about the time to the MRCA (TMRCA) given the observation that $S_{38} = 0$?

Note that the time to the MRCA is an unobservable random variable in the coalescent setting, and so the natural quantity to report is the conditional distribution of W_n given the data \mathcal{D} , which in this case is just the event $\{S_n = 0\}$. In this section we derive some of properties of such conditional distributions. In later sections we consider much richer problems concerning inference about the structure of the coalescent tree conditional on a sample. The main reference for the material in this section is Tavaré *et al.* (1997).

7.1 Samples of size two

Under the infinitely-many-sites assumption, all of the information in the two sequences is captured in S_2 , the number of segregating sites. Our goal, then, is to describe T_2 , the time to the most recent common ancestor of the sample in the light of the data, which is the observed value of S_2 .

One approach is to treat the realized value of T_2 as an unknown parameter which is then naturally estimated by $\tilde{T}_2 = S_2/\theta$, since $\mathbb{E}(S_2|T_2) = \theta T_2$. Such an approach, however, does not use all of the available information. In particular, the information available about T_2 due to the effects of genealogy and demography are ignored.

Under the coalescent model, when $n = 2$ the coalescence time T_2 has an exponential distribution with mean 1 before the data are observed. As Tajima (1983) noted, it follows from Bayes Theorem that after observing $S_2 = k$, the distribution of T_2 is gamma with parameters $1 + k$ and $1 + \theta$, which has probability density function

$$f_{T_2}(t|S_2=k) = \frac{(1 + \theta)^{1+k}}{k!} t^k e^{-(1+\theta)t}, \quad t \geq 0. \quad (7.1.1)$$

In particular,

$$\mathbb{E}(T_2|S_2=k) = \frac{1 + k}{1 + \theta}, \quad (7.1.2)$$

$$\text{var}(T_2|S_2=k) = \frac{1 + k}{(1 + \theta)^2}. \quad (7.1.3)$$

The pdf (7.1.1) conveys all of the information available about T_2 in the light of both the data and the coalescent model.

If a point estimate were required, equation (7.1.2) suggests the choice $\hat{T}_2 = (1+S_2)/(1+\theta)$. Perhaps not surprisingly, the estimator \hat{T}_2 , which is based on all of the available information, is superior to \tilde{T}_2 which ignores the pre-data information. For example, writing MSE for the mean square error of an estimator, straightforward calculations show that

$$\text{MSE}(\hat{T}_2) = \frac{1}{1+\theta} < \frac{1}{\theta} = \text{MSE}(\tilde{T}_2).$$

The difference in mean square errors could be substantial for small θ . In addition, the estimator \tilde{T}_2 is clearly inappropriate when $S_2 = 0$.

7.2 No variability observed in the sample

We continue to assume the infinitely-many-sites mutation model with parameter θ , and derive the distribution of $W_n := T_n + \dots + T_2$ given $S_n = 0$ for the case of constant population size. Several authors have been motivated to study this particular problem, among them Fu and Li (1996), Donnelly *et al.* (1996) and Weiss and von Haeseler (1996). Because mutations occur according to independent Poisson processes on the branches of the coalescent tree, we see that

$$\begin{aligned} \mathbb{E}(\exp(-uW_n)\mathbb{1}(S_n = 0)) &= \mathbb{E}[\mathbb{E}(\exp(-uW_n)\mathbb{1}(S_n = 0) \mid T_n, \dots, T_2)] \\ &= \mathbb{E}[\exp(-uW_n)\mathbb{E}(\mathbb{1}(S_n = 0) \mid T_n, \dots, T_2)] \\ &= \mathbb{E}[\exp(-uW_n)\exp(-\theta L_n/2)] \\ &= \prod_{j=2}^n \mathbb{E} \exp(-(u + \theta j/2)T_j) \\ &= \prod_{j=2}^n \frac{\binom{j}{2}}{\binom{j}{2} + u + \frac{\theta j}{2}} \end{aligned}$$

Since

$$\mathbb{P}(S_n = 0) = \prod_{j=1}^{n-1} \frac{j}{j + \theta},$$

we see that

$$\mathbb{E}(\exp(-uW_n) \mid S_n = 0) = \prod_{j=2}^n \frac{j(j + \theta - 1)/2}{u + j(j + \theta - 1)/2}. \tag{7.2.1}$$

Let \tilde{W}_n denote a random variable with the same distribution as the conditional distribution of W_n given $S_n = 0$. Equation (7.2.1) shows that we can write

$$\tilde{W}_n = \tilde{T}_n + \dots + \tilde{T}_2 \tag{7.2.2}$$

where the \tilde{T}_i are independent exponential random variables with parameters $\binom{i}{2} + \frac{i\theta}{2}$ respectively. Many properties of \tilde{W}_n follow from this. In particular

$$\mathbb{E}(W_n | S_n = 0) = \sum_{j=2}^n \frac{2}{j(j + \theta - 1)}. \tag{7.2.3}$$

The conditional density function of W_n may be calculated from a partial fraction expansion, resulting in the expression

$$f_{W_n}(t | S_n = 0) = \sum_{j=2}^n (-1)^j \frac{(2j + \theta - 1)n_{[j]}(\theta + 1)_{(j)}}{2(j - 2)!(\theta + n)_{(j)}} e^{-j(\theta + j - 1)t/2}. \tag{7.2.4}$$

The corresponding distribution function follows from

$$\mathbb{P}(W_n > t | S_n = 0) = \sum_{j=2}^n (-1)^{j-2} \frac{(2j + \theta - 1)n_{[j]}(\theta + 1)_{(j)}}{(j - 2)!j(j + \theta - 1)(\theta + n)_{(j)}} e^{-j(\theta + j - 1)t/2}.$$

Intuition suggests that given the sample has no variability, the post-data TMRCA of the sample should be stochastically smaller than the pre-data TMRCA. This can be verified by the following simple coupling argument. Let E_2, \dots, E_n be independent exponential random variables with parameters $\theta, \dots, n\theta/2$ respectively, and let T_2, \dots, T_n be independent exponential random variables with parameters $\binom{2}{2}, \dots, \binom{n}{2}$ respectively, independent of the E_i . Noting that $\tilde{T}_i = \min(T_i, E_i)$, we see that

$$\begin{aligned} \tilde{W}_n &= \tilde{T}_n + \dots + \tilde{T}_2 \\ &= \min(T_n, E_n) + \dots + \min(T_2, E_2) \\ &\leq T_n + \dots + T_2 \\ &= W_n, \end{aligned}$$

establishing the claim.

7.3 The rejection method

The main purpose of this section is to develop the machinery that allows us to find the joint distribution of the coalescent tree \mathcal{T} conditional on the sample of size n having configuration \mathcal{D} . Here \mathcal{D} is determined by the mutation process acting on the genealogical tree \mathcal{T} of the sample. Such conditional distributions lead directly to the conditional distribution of the height W_n of the tree.

The basic result we exploit to study such quantities is contained in

Lemma 7.1 *For any real-valued function g for which $\mathbb{E}|g(\mathcal{T})| < \infty$, we have*

$$\mathbb{E}(g(\mathcal{T}) | \mathcal{D}) = \frac{\mathbb{E}(g(\mathcal{T})\mathbb{P}(\mathcal{D} | \mathcal{T}))}{\mathbb{P}(\mathcal{D})}. \tag{7.3.1}$$

Proof. We have

$$\begin{aligned}\mathbb{E}(g(\mathcal{T})\mathbb{1}(\mathcal{D})) &= \mathbb{E}(\mathbb{E}(g(\mathcal{T})\mathbb{1}(\mathcal{D}|\mathcal{T}))) \\ &= \mathbb{E}(g(\mathcal{T})\mathbb{E}(\mathbb{1}(\mathcal{D})|\mathcal{T})) \\ &= \mathbb{E}(g(\mathcal{T})\mathbb{P}(\mathcal{D}|\mathcal{T})).\end{aligned}$$

Dividing this by $\mathbb{P}(\mathcal{D})$ completes the proof. \square

For most mutation mechanisms, explicit results are not available for these expectations, but we can develop a simple simulation algorithm. The expectation in (7.3.1) has the form

$$\mathbb{E}(g(\mathcal{T})|\mathcal{D}) = \int g(t) \frac{\mathbb{P}(\mathcal{D}|t)}{\mathbb{P}(\mathcal{D})} f_n(t) dt, \quad (7.3.2)$$

where $f_n(t)$ denotes the density of \mathcal{T} . The expression in (7.3.2) is a classical set-up for the rejection method:

Algorithm 7.1 To simulate from the distribution of \mathcal{T} given \mathcal{D} .

1. Simulate an observation t from the coalescent distribution of \mathcal{T} .
2. Calculate $u = \mathbb{P}(\mathcal{D}|t)$.
3. Keep t with probability u , else go to Step 1.

The joint distribution of the *accepted* trees t is precisely the conditional distribution of \mathcal{T} given \mathcal{D} .

The average number of times the rejection step is repeated per output observation is $1/\mathbb{P}(\mathcal{D})$, so that for small values of $\mathbb{P}(\mathcal{D})$ the method is likely to be inefficient. It can be improved in several ways. If, for example, there is a constant c such that

$$\mathbb{P}(\mathcal{D}|t) \leq c \text{ for all values of } t,$$

then u in Step 2 of the algorithm can be replaced by u/c .

Note that if properties of W_n are of most interest, observations having the conditional distribution of W_n given \mathcal{D} can be found from the trees generated in algorithm 7.1. When the data are summarized by the number S_n of segregating sites, these methods become somewhat more explicit, as is shown in the next section.

7.4 Conditioning on the number of segregating sites

In this section we consider events of the form

$$\mathcal{D} \equiv \mathcal{D}_k = \{S_n = k\},$$

corresponding to the sample of size n having k segregating sites. Since each mutation in the coalescent tree corresponds to a segregating site, it follows that

$$\mathbb{P}(\mathcal{D}|\mathcal{T}) = \mathbb{P}(\mathcal{D}_k|L_n) = \text{Po}(\theta L_n/2)\{k\},$$

where $L_n = 2T_2 + \dots + nT_n$ is the total length of the ancestral tree of the sample and $\text{Po}(\lambda)\{k\}$ denotes the Poisson point probability

$$\text{Po}(\lambda)\{k\} = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k = 0, 1, \dots$$

Therefore

$$\mathbb{E}(g(W_n)|\mathcal{D}_k) = \frac{\mathbb{E}(g(W_n)\text{Po}(\theta L_n/2)\{k\})}{\mathbb{E}(\text{Po}(\theta L_n/2)\{k\})} \quad (7.4.1)$$

The simulation algorithm 7.1 then becomes

Algorithm 7.2 To simulate from the joint density of T_2, \dots, T_n given \mathcal{D}_k .

1. Simulate an observation $\mathbf{t} = (t_n, \dots, t_2)$ from the joint distribution of $\mathbf{T}_n = (T_n, \dots, T_2)$. Calculate $l = 2t_2 + \dots + nt_n$.
2. Calculate $u = \mathbb{P}(\mathcal{D}_k|\mathbf{t}) = \text{Po}(\theta l/2)\{k\}$.
3. Keep \mathbf{t} with probability u , else go to Step 1.

The joint distribution of the *accepted* vectors \mathbf{t} is precisely the conditional distribution of \mathbf{T}_n given \mathcal{D}_k .

Since

$$\mathbb{P}(S_n = k|\mathbf{t}) = \text{Po}(\theta l_n/2)\{k\} \leq \text{Po}(k)\{k\},$$

where we define $\text{Po}(0,0) = 1$, the modified algorithm becomes:

Algorithm 7.3 To simulate from the joint density of T_2, \dots, T_n given $S_n = k$.

1. Simulate an observation $\mathbf{t} = (t_n, \dots, t_2)$ from the joint distribution of $\mathbf{T}_n = (T_n, \dots, T_2)$.
2. Calculate $l = 2t_2 + \dots + nt_n$, and set

$$u = \frac{\text{Po}(l\theta/2)\{k\}}{\text{Po}(k)\{k\}}$$

3. Keep \mathbf{t} with probability u , else go to Step 1.

Values of $w_n = t_2 + \dots + t_n$ calculated from accepted vectors \mathbf{t} have the conditional distribution of W_n given $S_n = k$.

Notice that nowhere have we assumed a particular form for the distribution of \mathbf{T}_n . In particular, the method works when the population size is variable so long as \mathbf{T}_n has the distribution specified by (2.4.8). For an analytical approach to the constant population size case, see Fu (1996).

Remark. In these examples, we have simulated the ancestral process back to the common ancestor. It is clear, however, that the same approach can be used to simulate observations for any fixed time t into the past. All that is required is to simulate coalescence times back into the past until time t , and then the effects of mutation (together with the genetic types of the ancestors at time t) can be superimposed on the coalescent forest.

Example

We use this technique to generate observations from the model with variable population size when the conditioning event is \mathcal{D}_0 . The particular population size function we use for illustration is

$$f(x) = \alpha^{\min(t/v, 1)}, \quad (7.4.2)$$

corresponding to a population of constant relative size α more than (coalescent) time v ago, and exponential growth from time v until the present relative size of 1.

In the illustration, we chose $V = 50,000$ years, $N = 10^8$, a generation time of 20 years and $\alpha = 10^{-4}$. Thus $v = 2.5 \times 10^{-5}$. We compare the conditional distribution of W_n given \mathcal{D}_0 to that in the constant population size case with $N = 10^4$. Histograms of 5000 simulated observations are given in Figures 7.1 and 7.2. The mean of the conditional distribution in the constant population size case is 313,200 years, compared to 358,200 years in the variable case. Examination of other summary statistics of the simulated data (Table 7) shows that the distribution in the variable case is approximately that in the constant size case, plus about V years. This observation is supported by the plot of the empirical distribution functions of the two sets in Figure 7.3.

The intuition behind this is clear. Because of the small sample size relative to the initial population size N , the sample of size n will typically have about n distinct ancestors at the time of the expansion, V . These ancestors themselves form a random sample from a population of size αN .

Table 7. Summary statistics from 5000 simulation runs

	constant variable	
mean	313,170	358,200
std dev	156,490	158,360
median	279,590	323,210
5%	129,980	176,510
95%	611,550	660,260

Fig. 7.1. Histogram of 5000 replicates for constant population size, $N = 10^4$

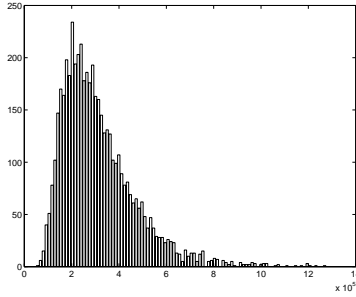


Fig. 7.2. Histogram of 5000 replicates for variable population size, $N = 10^8, T = 50,000, \alpha = 10^{-4}$

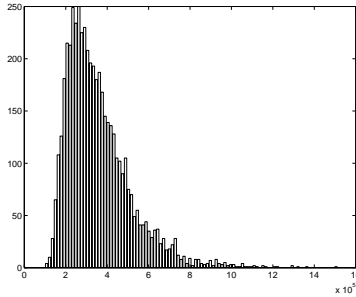
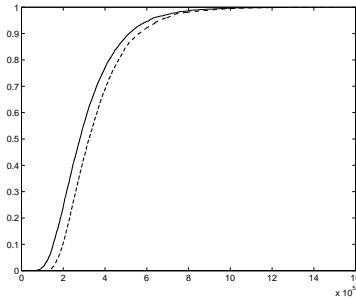


Fig. 7.3. Empirical distribution function. Solid line is constant population size case.



7.5 An importance sampling method

If moments of the post-data distribution of W_n , say, are required, then they can be found in the usual way from observations generated by Algorithm 7.2. As an alternative, an importance sampling scheme can be used. This is best illustrated by an example. Consider then the expression in (7.4.1). We have

$$\mathbb{E}(g(W_n)|\mathcal{D}_k) = \frac{\mathbb{E}(g(W_n)\text{Po}(\theta L_n/2)\{k\})}{\mathbb{E}(\text{Po}(\theta L_n/2)\{k\})}.$$

Point estimates of this quantity can be found by simulating independent copies $(W_n^{(j)}, L_n^{(j)})$, $j = 1, 2, \dots, R$ of the height and length of the ancestral tree and computing the ratio estimator

$$r_R = \frac{\sum_{j=1}^R g(W_n^{(j)})\text{Po}(\theta L_n^{(j)}/2)\{k\}}{\sum_{j=1}^R \text{Po}(\theta L_n^{(j)}/2)\{k\}}. \tag{7.5.1}$$

One application provides an estimate of the conditional distribution function of W_n given \mathcal{D}_k : Suppose that we have ordered the points $W_n^{(j)}$ and listed them as $W_n^{[1]} < W_n^{[2]} < \dots < W_n^{[R]}$. Let $L_n^{[1]}, \dots, L_n^{[R]}$ be the corresponding L -values. The empirical distribution function then has jumps of height

$$\frac{e^{-\theta L_n^{[l]}/2}}{\sum_{j=1}^R e^{-\theta L_n^{[j]}/2}}$$

at the points $W_n^{[l]}$, $l = 1, 2, \dots, R$.

This approach uses all the simulated observations, but requires either knowing which g are of interest, or storing a lot of observations. Asymptotic properties of the ratio estimator can be found from standard theory.

7.6 Modeling uncertainty in N and μ

In this section, we use prior information about the distribution of μ , as well as information that captures our uncertainty about the population size N . We begin by describing some methods for generating observations from the posterior distribution of the vector (W_n, N, μ) given the data \mathcal{D} . We use this to study the posterior distribution of the time W_n to a common ancestor, measured in years:

$$W_n^y = N \times G \times W_n.$$

The rejection method is based on the analog of (7.3.1):

$$\mathbb{E}(g(\mathbf{T}_n, N, \mu)|\mathcal{D}) = \frac{\mathbb{E}(g(\mathbf{T}_n, N, \mu)\mathbb{P}(\mathcal{D}|\mathbf{T}_n, N, \mu))}{\mathbb{P}(\mathcal{D})}. \tag{7.6.1}$$

This converts once more into a simulation algorithm; for definiteness we suppose once more that $\mathcal{D} = \{S_n = k\}$.

Algorithm 7.4 To simulate from conditional distribution of \mathbf{T}_n, N, μ given $S_n = k$.

1. Generate an observation \mathbf{t}, N, μ from the joint distribution of \mathbf{T}_n, N, μ .
2. calculate $l = 2t_2 + \dots + nt_n$, and

$$u = \frac{\text{Po}(lN\mu)\{k\}}{\text{Po}(k)\{k\}}$$

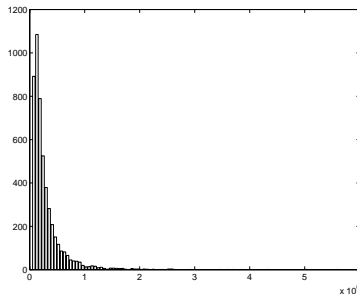
3. accept \mathbf{t}, N, μ with probability u , else go to Step 1.

Usually we assume that N and μ are independent of \mathbf{T}_n , and that N and μ are themselves independent.

Examples

Suppose that no variation is observed in the data, so that \mathcal{D}_0 . Suppose that N has a lognormal distribution with parameters $(10, 1)$, and that μ has a Gamma distribution with mean μ_0 and standard deviation $C\mu_0$. A constant size population is assumed. In the example, we took $\mu_0 = 2 \times 10^{-5}$ and $C = 1/20$ and $C = 1.0$. Histograms appear in Figures 7.4 and 7.5, and some summary statistics are given in Table 8.

Fig. 7.4. Histogram of 5000 replicates $C = 1/20$



Here we illustrate for the exponential growth model described earlier, with initial population size $N = 10^8$, and $\alpha = 10^{-4}$. We took N lognormally distributed with parameters 17.92, 1. (The choice of 17.92 makes the mean of $N = 10^8$.) For μ we took the Gamma prior with mean $= \mu_0$, and standard deviation $C\mu_0$. In the simulations, we used $C = 1$ and $C = 1/20$. Histograms of 5000 simulated observations are given in Figures 7.6 and 7.7. Some summary statistics are given in Table 9.

The importance sampling method also readily adapts to this Bayesian setting: apply the approach outlined in (7.5.1) to the expectation formula in (7.6.1).

Fig. 7.5. Histogram of 5000 replicates $C = 1$

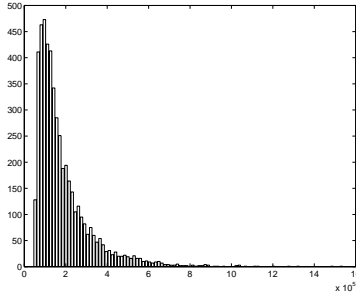


Table 8. Summary statistics from 5000 simulation runs. Prior mean $\mu_0 = 2 \times 10^{-5}$, $\mathcal{D} = \mathcal{D}_0$

	$C = 1.0$	$C = 1/20$
mean	647,821	262,590
median	369,850	204,020
5%	68,100	52,372
95%	2,100,000	676,890

Fig. 7.6. Histogram of 5000 replicates. Variable size model. $C = 1/20$

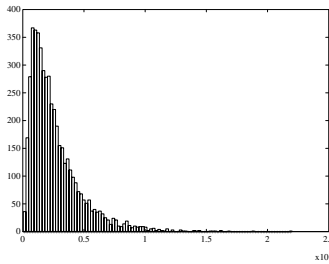


Fig. 7.7. Histogram of 5000 replicates. Variable size model. $C = 1$

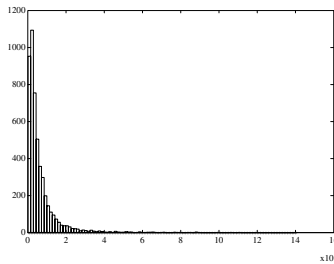


Table 9. Summary statistics from 5000 simulation runs. Prior mean $\mu_0 = 2 \times 10^{-5}$, $\mathcal{D} = \mathcal{D}_0$

	$C = 1$	$C = 1/20$
mean	292,000	186,000
median	194,000	141,490
5%	70,600	65,200
95%	829,400	462,000

7.7 Varying mutation rates

These rejection methods can be employed directly to study the behavior of the infinitely-many-sites model that allows for several regions with different mutation rates. Suppose then that there are r regions, with mutation rates μ_1, \dots, μ_r . The analysis also applies, for example, to r different types of mutations within a given region. We sample n individuals, and observe k_1 segregating sites in the first region, k_2 in the second, \dots , and k_r in the r^{th} . The problem is to find the conditional distribution of \mathcal{T} , given the vector (k_1, \dots, k_r) .

When N and the μ_i are assumed known, this can be handled by a modification of Algorithm 7.2. Conditional on L_n , the probability of (k_1, \dots, k_r) is

$$h(L_n) = \text{Po}(k_1, L_n\theta_1/2) \times \dots \times \text{Po}(k_r, L_n\theta_r/2),$$

where $\theta_i = 2N\mu_i$, $i = 1, 2, \dots, r$. It is easy to check that $h(L_n) \leq h(k/\theta)$, where

$$k = k_1 + \dots + k_r, \quad \theta = \theta_1 + \dots + \theta_r.$$

Therefore in the rejection algorithm we may take $u = h(L_n)/h(k/\theta)$ which simplifies to

$$u = h(L_n)/h(k/\theta) = \frac{\text{Po}(L_n\theta/2)\{k\}}{\text{Po}(k)\{k\}}. \quad (7.7.1)$$

Equation (7.7.1) establishes the perhaps surprising fact that the conditional distribution of W_n given (k_1, \dots, k_r) and $(\theta_1, \dots, \theta_r)$ depends on these values only through their respective totals: the total number of segregating sites k and the total mutation rate θ . Thus Algorithm 7.2 can be employed directly with the appropriate values of k and θ . This result justifies the common practice of analyzing segregating sites data through the total number of segregating sites, even though these sites may occur in regions of differing mutation rate.

If allowance is to be made for uncertainty about the μ_i , then this simplification no longer holds. However, Algorithm 7.3 can be employed with the rejection step replaced by (7.7.2):

$$u = \frac{\text{Po}(L_n\theta_1/2)\{k_1\}}{\text{Po}(k_1)\{k_1\}} \dots \frac{\text{Po}(L_n\theta_r/2)\{k_r\}}{\text{Po}(k_r)\{k_r\}}. \quad (7.7.2)$$

In this case, Step 2 requires generation of a vector of rates $\mu = (\mu_1, \dots, \mu_r)$ from the joint prior π_μ . Furthermore, the algorithm immediately extends to the case of variable population size.

7.8 The time to the MRCA of a population given data from a sample

In this section, we show how the rejection technique can be used to study the time T_m to the MRCA of a sample of m individuals, conditional on the number of segregating sites in a subsample of size n . In many applications of ancestral inference, the real interest is on the time to the MRCA of the *population*, given data on a *sample*. This can be obtained by setting $m = N$ below. See Tavaré (1997) and Tavaré *et al.* (1997) for further details and examples.

The quantities of interest here are A_m (the number of distinct ancestors of the sample), A_n (the number of distinct ancestors of the subsample), and W_n (the time to the MRCA of the subsample). The results of Saunders *et al.* (1984) justify the following algorithm:

Algorithm 7.5 Rejection algorithm for $f_{W_m}(t|S_n=k)$.

1. Set $A_m = m, A_n = n, W_n = 0, L_n = 0$
2. Generate E , exponential of rate $A_m(A_m - 1)/2$. Set $W_n = W_n + W, L_n = L_n + A_n \cdot E$.
3. Set $p = \frac{A_n(A_n-1)}{A_m(A_m-1)}$. Set $A_m = A_m - 1$. With probability p set $A_n = A_n - 1$. If $A_n > 1$ go to 2.
4. Set $u = \text{Po}(\theta L_n/2)\{k\}/\text{Po}(k)\{k\}$. Accept (A_m, W_n) with probability u , else go to 1.
5. If $A_m = 1$, set $T_{nm} = 0$, and return $W_m = W_n$. Else, generate independent exponentials E_j with parameter $j(j - 1)/2$, for $j = 2, 3, \dots, A_m$, and set $T_{nm} = E_2 + \dots + E_{A_m}$. Return $W_m = W_n + T_{nm}$.

Many aspects of the joint behavior of the sample and a subsample can be studied using this method. In particular, values of (A_m, W_n) accepted at step 5 have the joint conditional distribution of the number of ancestors of the sample at the time the subsample reaches its common ancestor and the time of the MRCA of the subsample, conditional on the number of segregating sites in the subsample. In addition, values of T_{nm} produced at step 5 have the conditional distribution of the time between the two most recent common ancestors. It is straightforward to modify the method to cover the case of variable population size, and the case where uncertainty in N and μ is modeled. With high probability, the sample and the subsample share a common ancestor and therefore a common time to the MRCA. However, if the two common ancestors differ then the times to the MRCA can differ substantially. This is explored further in the examples below.

Examples

Whitfield *et al.* (1995) describe another Y chromosome data set that includes a sample of $n = 5$ humans. The 15,680 bp region has three polymorphic nucleotides that once again are consistent with the infinitely-many-sites model. They estimated the coalescence time of the sample to be between 37,000 and 49,000 years. Again, we present several reanalyses, each of which is based on the number of segregating sites in the data. The results are summarized in Table 10 and illustrated in Figure 7.8.

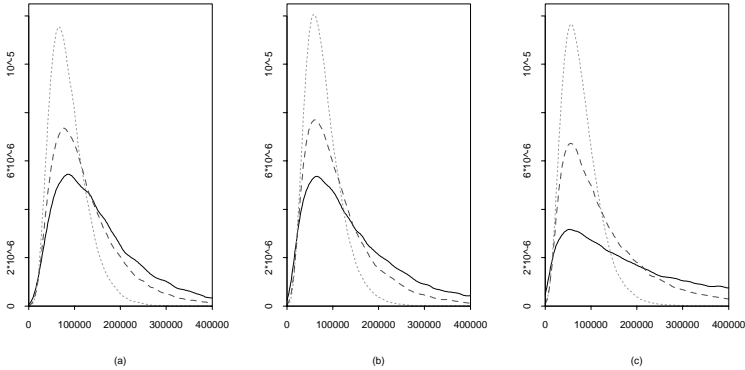
Table 10. Results of re-analyses of the data of Whitfield *et al.* In each case the data are $S_5 = 3$. Line (a) gives the interval reported by the authors (but note that they assigned no probability to their interval). Mean and 95% interval are estimated from samples of size 10,000. Details of the gamma and lognormal distributions are given in the text.

	Model	Mean of W_5 ($\times 10^3$)		95% Interval ($\times 10^3$)	
		pre-data	post-data	pre-data	post-data
(a)	Whitfield <i>et al.</i>				37 – 49
(b)	$N = 4,900$ $\mu_S = 3 \cdot 52 \times 10^{-4}$	157	87	31 – 429	30 – 184
(c)	$N = 4,900$ μ_S gamma	157	125	31– 429	32 – 321
(d)	N gamma $\mu_S = 3 \cdot 52 \times 10^{-4}$	159	80	21 – 517	26 – 175
(e)	N gamma μ_S gamma	159	117	21– 517	25 – 344
(f)	N lognormal μ_S gamma	428	149	19 – 2,200	22 – 543

In estimating the coalescence time, Whitfield *et al.* adopt a method which does not use population genetics modeling. While the method is not systematically biased, it may be inefficient to ignore pre-data information about plausible values of the coalescence time. In addition, the method substantially underrepresents the uncertainty associated with the estimates presented. Here, we contrast the results of such a method with those of one which does incorporate background information.

To determine the mutation rate, we use the average figure of $1 \cdot 123 \times 10^{-9}$ substitutions per nucleotide position per year given in Whitfield *et al.*, and a

Fig. 7.8. Probability density curves for W_5 . In each panel the three curves correspond to: solid, pre-data; dashed, post-data, assuming μ_S gamma; dotted, post-data assuming $\mu_S = 3 \cdot 52 \times 10^{-4}$. The three panels correspond to (a) $N = 4,900$; (b) N gamma; (c) N lognormal.



generation time of 20 years, to give $\mu = 15,680 \times 1.123 \times 10^{-9} \times 20 = 3.52 \times 10^{-4}$ substitutions per generation. For these parameter values, the post-data mean of W_5 is 87,000 years.

As noted in the previous section, the appropriate values of the parameters are not known. Analysis (c) incorporates uncertainty about μ , in the form of a gamma distribution with shape parameter 2 and mean $3 \cdot 52 \times 10^{-4}$, while continuing to assume that N is known to be 4,900. The effect is to greatly increase the post-data mean of W_5 . Allowing N to be uncertain while μ_S is known has, on the other hand, the effect of slightly reducing the post-data estimates of W_5 , compared with the case that N and μ_S are both known. This may be attributed to the data favoring values of N smaller than 4,900.

Analyses (e) and (f) incorporate uncertainty about both N and μ_S . They use the same prior distributions as analyses (g) and (i) respectively of the previous section. Note that, as should be expected, the uncertainty about T is larger than when one or both of N and μ_S are assumed known exactly.

Whitfield *et al.* (1995) point to their estimated coalescence time as being substantially shorter than those published for the human mitochondrial genome. In contrast, the ranges in each of our analyses (b) – (e) overlap with recent interval estimates for the time since mitochondrial Eve. In addition, recall that the quantity W_5 being estimated in Table 10 is the coalescence time of the sample of 5 males sequenced in the study. This time may be different from, and substantially shorter than, the coalescence time of *all* existing Y chromosomes. Under the assumption that $N = 4,900$ and $\mu = 3.52 \times 10^{-4}$, Algorithm 7.5 can be used to show that the mean time to the common ancestor

of the male *population*, given $S_5 = 3$, is 157,300 years, with a corresponding 95% interval of (58,900 – 409,800) years. These figures differ markedly from the corresponding values for the sample, given at line (b) of Table 10. It is the population values which are likely to be of primary interest.

7.9 Using the full data

The approach that conditions on the number of segregating sites in the data is convenient primarily because the rejection methods are quick and easy to program. However, it does not make full use of the data. In this section, we discuss how we can approximate the conditional distribution of TMRCA given the infinitely-many-sites rooted tree (T, \mathbf{n}) that corresponds to the data, or the corresponding unrooted tree (Q, \mathbf{n}) . See Griffiths and Tavaré (1994, 1999) for further details.

Consider first the rooted case. The probability $q(t, x)$ that a sample taken at time t has configuration x satisfies an equation of the form

$$q(t, x) = \int_t^\infty \sum_y r(s; x, y) q(s, y) g(t, x; s) ds$$

for a positive kernel r . For the case of an unrooted tree, we have $x = (T, \mathbf{n})$. Now define

$$q(t, x, w) = \mathbb{P}(\text{sample taken at time } t \text{ has configuration } x \\ \text{and TMRCA} \leq t + w)$$

By considering the time of the first event in the history of the sample, it can be seen that $q(t, x, w)$ satisfies the equation

$$q(t, x, w) = \int_t^\infty \sum_y r(s; x, y) q(s, y, t + w - s) g(t, x; s) ds \quad (7.9.1)$$

where we assume that $q(t, x, y) = 0$ if $y < t$. Recursions of this type can be solved using the Markov chain simulation technique described in Section 6. The simplest method is given in (6.5.3): we define

$$f(s; x) = \sum_y r(s; x, y) \\ P(s; x, y) = \frac{r(s; x, y)}{f(s; x)},$$

and rewrite (7.9.1) in the form

$$q(t, x, w) = \int_t^\infty f(s; x) \sum_y P(s; x, y) q(s, y, t + w - s) g(t, x; s) ds. \quad (7.9.2)$$

The Markov chain associated with the density g and the jump matrix P is once again denoted by $X(\cdot)$. The representation we use is then

$$q(t, x, w) = \mathbb{E}_{(t,x)} q(\tau, X(\tau), t + w - \tau) \prod_{j=1}^k f(\tau_j; X(\tau_{j-1})), \quad (7.9.3)$$

where $t = \tau_0 < \tau_1 < \dots < \tau_k = \tau$ are the jump times of $X(\cdot)$, and τ is the time taken to reach the set A that corresponds to a sample configuration x for a single individual. For the infinitely-many-sites tree, this corresponds to a tree of the form (T, \mathbf{e}_1) .

The natural initial condition is

$$q(t, x, w) = \mathbb{1}(w \geq 0), \quad x \in A,$$

so that

$$q(\tau, X(\tau), t + w - \tau) = \mathbb{1}(\tau < t + w).$$

The Monte Carlo method generates R independent copies of the X process, and for the i th copy calculates the observed value

$$F_i = \prod_{j=1}^{k_i} f(\tau_j^i; X^i(\tau_{j-1}^i)).$$

and estimates $q(t, x, w)$ by

$$\hat{q}(t, x, w) = \frac{\sum_{i=1}^R F_i \mathbb{1}(\tau^i \leq t + w)}{\sum_{i=1}^R F_i}.$$

The distribution function of TMRCA given the data (t, x) can be therefore be approximated by a step function that jumps a height $F_{(l)}/\sum F_i$ at the point $\tau_{(l)}$, where the $\tau_{(l)}$ are the increasing rearrangement of the times τ^i , and the $F_{(l)}$ are the corresponding values of the F_i .

This method can be used immediately when the data correspond to a rooted tree (T, \mathbf{n}) . When the data correspond to an unrooted tree (\mathbf{Q}, \mathbf{n}) we proceed slightly differently. Corresponding to the unrooted tree (\mathbf{Q}, \mathbf{n}) are rooted trees (T, \mathbf{n}) . An estimator of $\mathbb{P}(TMRCA \leq t + w, (T, \mathbf{n}))$ is given by

$$\frac{1}{R} \sum_{i=1}^R F_i(T) \mathbb{1}(\tau_i(T) \leq t + w),$$

the T denoting a particular rooted tree. Recalling (5.9.3), an estimator of $q(t, (\mathbf{Q}, \mathbf{n}), w)$ is therefore given by

$$\sum_T \frac{1}{R} \sum_{i=1}^R F_i(T) \mathbb{1}(\tau_i(T) \leq t + w),$$

and the conditional probability $q(t, (\mathbf{Q}, \mathbf{n}), w)/q(t, (\mathbf{Q}, \mathbf{n}))$ is estimated by

$$\frac{\sum_T \sum_{i=1}^R F_i(T) \mathbb{1}(\tau_i(T) \leq t + w)}{\sum_T \sum_{i=1}^R F_i(T)}.$$

The distribution of TMRCAs given data (\mathbf{Q}, \mathbf{n}) taken at time t is found by ranking all the times $\tau_j(T)$ over different T to get the increasing sequence $\tau_{(j)}$, together with the corresponding values $F_{(j)}$, and then approximating the distribution function by jumps of height $F_{(j)}/\sum F_{(j)}$ at the point $\tau_{(j)}$. Usually we take $t = 0$ in the previous results.

8 The Age of a Unique Event Polymorphism

In this section we study the age of an allele observed in a sample of chromosomes. Suppose then that a particular mutation Δ has arisen just once in the history of the population of interest. This mutation has an age (the time into the past at which it arose), and we want to infer its distribution given data \mathcal{D} . These data can take many forms:

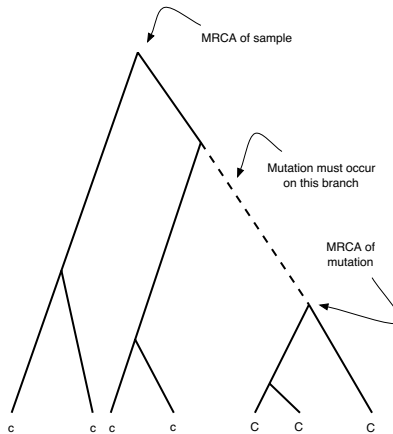
- the number of copies, b , of Δ observed in a sample of size n . Here we assume that $1 \leq b < n$, so that the mutation is segregating in the sample.
- the number of copies of Δ together with other molecular information about the region around Δ . For example, we might have an estimate of the number of mutations that have occurred in a linked region containing Δ .
- in addition, we might also have molecular information about the individuals in the sample who do not carry Δ .

The unique event polymorphism (UEP) assumption leads to an interesting class of coalescent trees that we study in the next section.

8.1 UEP trees

Suppose that the mutation Δ is represented b times in the sample. The UEP property means that the b sequences must coalesce together before any of the non- Δ sequences share any common ancestors with them. This situation is illustrated in Figure 8.1 for $n = 7$ and $b = 3$.

Fig. 8.1. Tree with UEP. The individuals carrying the special mutation Δ are labeled C , those not carrying the mutation are labeled c .



To understand the structure of these trees, we begin by studying the properties of trees that have the property \mathcal{E} that a particular b sequences coalesce together before any of the other $n - b$ join their subtree. To this end, let $n > J_{b-1} > \dots > J_1$ be the total number of distinct ancestors of the sample at the time the b first have $b - 1, \dots, 1$ distinct ancestors, and let J_0 ($1 \leq J_0 < J_1$) be the number of ancestors in the sample at the time the first of the other $n - b$ sequences shares a common ancestor with an ancestor of the b . In Figure 8.1, we have $J_2 = 5, J_1 = 4, J_0 = 2$.

It is elementary to find the distribution of J_{b-1}, \dots, J_0 . Recalling that in a coalescent tree joins are made at random, we find that

$$\begin{aligned} \mathbb{P}(J_r = j_r, r = b - 1, \dots, 0) &= \prod_{r=2}^b \left\{ \frac{\binom{j_r-r}{2}}{\binom{j_r}{2}} \dots \frac{\binom{j_{r-1}+2-r}{2}}{\binom{j_{r-1}}{2}} \frac{\binom{r}{2}}{\binom{j_{r-1}+1}{2}} \right\} \\ &\quad \times \frac{\binom{j_1-1}{2}}{\binom{j_1}{2}} \dots \frac{\binom{j_0+2-1}{2}}{\binom{j_0+2}{2}} \frac{j_0}{\binom{j_0+1}{2}} \end{aligned}$$

where we have defined $j_b = n$, and where $1 \leq j_0 < j_1 < \dots < j_{b-1} < n$. This expression can be simplified to give

$$\mathbb{P}(J_r = j_r, r = b - 1, \dots, 0) = \frac{2b!(b-1)!(n-b)!(n-b-1)!j_0}{n!(n-1)!}. \tag{8.1.1}$$

We can find $\mathbb{P}(\mathcal{E})$ by summing $1 \leq j_0 < j_1 < \dots < j_{b-1} < n$. Note that

$$\begin{aligned} \sum_{j_0=1}^{n-b} \sum_{j_0 < j_1 < \dots < j_{b-1} < n} 1 &= \sum_{j_0=1}^{n-b} j_0 \binom{n-j_0-1}{b-1} \\ &= \sum_{l=0}^{n-b-1} (l+1) \binom{n-1-l-1}{n-b-1-l} \\ &= \binom{n}{n-b-1}, \end{aligned}$$

the last equality coming from the identity

$$\sum_{k=1}^c \binom{c}{k} \binom{d+k}{d+1} = \binom{d}{c-1},$$

valid for integral c, d with $c = b, d = 2$. It follows that

$$\begin{aligned} \mathbb{P}(\mathcal{E}) &= \frac{2b!(b-1)!(n-b)!(n-b-1)!}{n!(n-1)!} \binom{n}{n-b-1} \\ &= \frac{2}{b+1} \binom{n-1}{b-1}^{-1}, \end{aligned} \tag{8.1.2}$$

as found by Wiuf and Donnelly (1999).

Now we can compute the conditional distribution of ‘everything’ given \mathcal{E} . For example it follows that for $1 \leq j_0 < j_1 < \dots < j_{b-1} < n$

$$\mathbb{P}(J_r = j_r, r = b - 1, \dots, 0 \mid \mathcal{E}) = j_0 \binom{n}{b+1}^{-1}, \quad (8.1.3)$$

while for $1 \leq j_0 < j_1 < n$,

$$\mathbb{P}(J_1 = j_1, J_0 = j_0 \mid \mathcal{E}) = j_0 \binom{n - j_1 - 1}{b - 2} \binom{n}{b+1}^{-1} \quad (8.1.4)$$

and for $1 < j_1 < j_2 \dots < j_{b-1} < n$,

$$\mathbb{P}(J_r = j_r, r = b - 1, \dots, 2 \mid J_1 = j_1, J_0 = j_0, \mathcal{E}) = \binom{n - j_1 - 1}{b - 2}^{-1}. \quad (8.1.5)$$

Having discussed the topological properties of UEP coalescent trees, we move on to the age of the mutation itself.

The distribution of J_Δ

Suppose that Δ mutations occur at rate $\mu/2$ on the branches of the coalescent tree. The random variable J_Δ gives the number of ancestors of the sample of size n when the mutation Δ occurs. Clearly, $J_0 < j_\Delta \leq J_1$. Its distribution can be found as follows. To get $J_\Delta = k$, a single mutation must arise on the branch of length T_k , and no other mutations must occur in the remainder of the coalescent tree. It follows from (8.1.4) that for $1 \leq j_0 < k \leq j_1 \leq n - b + 1$,

$$\mathbb{P}(J_1 = j_1, J_\Delta = k, J_0 = j_0 \mid \mathbf{T}, \mathcal{E}) = \frac{\mu}{2} T_k e^{-L_n \mu/2} j_0 \binom{n - j_1 - 1}{b - 2} \binom{n}{b+1}^{-1},$$

where $\mathbf{T} = (T_n, \dots, T_2)$ and $L_n = nT_n + \dots + 2T_2$ is the total length of the tree. Using the fact that for integral k ,

$$\sum_{j=0}^k \binom{c + k - j - 1}{k - j} \binom{d + j - 1}{j} = \binom{c + d + k - 1}{k}$$

we see that

$$\sum_{j_1=k}^{n-b+1} \binom{n - j_1 - 1}{b - 2} = \binom{n - k}{b - 1},$$

so that

$$\sum_{j_0=1}^{k-1} \sum_{j_1=k}^{n-b+1} j_0 \binom{n - j_1 - 1}{b - 2} = \frac{k(k - 1)}{2} \binom{n - k}{b - 1}.$$

Hence

$$\mathbb{P}(J_\Delta = k \mid \mathbf{T}, \mathcal{E}) = \frac{\mu}{2} T_k e^{-L_n \mu / 2} \frac{k(k-1)}{2} \binom{n-k}{b-1} \binom{n}{b+1}^{-1}, \quad (8.1.6)$$

and

$$\mathbb{P}(J_\Delta = k \mid \mathcal{E}) = \mathbb{E} \left(\frac{\mu}{2} T_k e^{-L_n \mu / 2} \right) \frac{k(k-1)}{2} \binom{n-k}{b-1} \binom{n}{b+1}^{-1}.$$

Letting \mathcal{U} denote the event that there is indeed a UEP, we have

$$\mathbb{P}(\mathcal{U} \mid \mathcal{E}) = \sum_{k=2}^{n-b+1} \mathbb{P}(J_\Delta = k \mid \mathcal{E}),$$

so that for $k = 2, \dots, n - b + 1$,

$$\mathbb{P}(J_\Delta = k \mid \mathcal{U} \cap \mathcal{E}) = \frac{k(k-1) \binom{n-k}{b-1} \mathbb{E} [T_k e^{-L_n \mu / 2}]}{\sum_{l=2}^{n-b+1} l(l-1) \binom{n-l}{b-1} \mathbb{E} [T_l e^{-L_n \mu / 2}]}. \quad (8.1.7)$$

Remark. In the constant population size case, this gives

$$\mathbb{P}(J_\Delta = k \mid \mathcal{U} \cap \mathcal{E}) = \frac{(k-1) \binom{n-k}{b-1} \frac{1}{k-1+\mu}}{\sum_{l=2}^{n-b+1} (l-1) \binom{n-l}{b-1} \frac{1}{l-1+\mu}},$$

as given by Stephens (2000).

Similar arguments show that for $k \leq j_1 < j_2 < \dots < j_{b-1} < n$,

$$\mathbb{P}(J_1 = j_1, \dots, J_{b-1} = j_{b-1} \mid J_\Delta = k, \mathcal{U} \cap \mathcal{E}) = \binom{n-k}{b-1}^{-1}, \quad (8.1.8)$$

so that given $J_\Delta = k$, the places where the subtree has joins form a random (ordered) $(b-1)$ -subset of the integers $k, k+1, \dots, n-1$. Hence for $1 \leq i \leq b-1$ and $k \leq j_1 < \dots < j_i < n-i+b$,

$$\mathbb{P}(J_1 = j_1, \dots, J_i = j_i \mid J_\Delta = k, \mathcal{U} \cap \mathcal{E}) = \binom{n-j_i-1}{b-i-1} \binom{n-k}{b-1}^{-1}. \quad (8.1.9)$$

8.2 The distribution of T_Δ

We let J_Δ be the number of ancestors of the sample at the time the unique Δ mutation occurs. Clearly $J_0 < J_\Delta \leq J_1$. We can find the conditional distribution of the age T_Δ as follows. We have

$$\begin{aligned}
 & \mathbb{E}(e^{-\phi T_\Delta} \mid \mathcal{E}) \\
 &= \sum_{k=2}^{n-b+1} \mathbb{E}(e^{-\phi T_\Delta} \mathbb{1}(J_\Delta = k) \mid \mathcal{E}) \\
 &= \sum_{k=2}^{n-b+1} \mathbb{E}(\mathbb{E}(e^{-\phi T^{[k]}} \mathbb{1}(J_\Delta = k) \mid \mathbf{T}, \mathcal{E})) \\
 &= \sum_{k=2}^{n-b+1} \mathbb{E}(e^{-\phi T^{[k]}} \mathbb{P}(J_\Delta = k \mid \mathbf{T}, \mathcal{E})) \\
 &= \sum_{k=2}^{n-b+1} \mathbb{E}(e^{-\phi T^{[k]}} \frac{T_k \mu}{2} e^{-L_n \mu/2}) \frac{k(k-1)}{2} \binom{n-k}{b-1} \binom{n}{b+1}^{-1} \quad (8.2.1)
 \end{aligned}$$

where

$$T^{[k]} = T_n + \dots + T_{k+1} + UT_k, \quad (8.2.2)$$

and U is uniformly distributed on $(0, 1)$, independent of \mathbf{T} . The penultimate inequality comes from (8.1.6). This gives us:

Theorem 8.1 *The Laplace transform of the conditional distribution of the age T_Δ of a UEP observed b times in a sample of size n (where $0 < b < n$) is given by*

$$\begin{aligned}
 & \mathbb{E}(e^{-\phi T_\Delta} \mid \mathcal{U} \cap \mathcal{E}) \\
 &= \frac{\sum_{k=2}^{n-b+1} k(k-1) \binom{n-k}{b-1} \mathbb{E} \left[e^{-\phi T^{[k]}} T_k e^{-L_n \mu/2} \right]}{\sum_{k=2}^{n-b+1} k(k-1) \binom{n-k}{b-1} \mathbb{E} \left[T_k e^{-L_n \mu/2} \right]} \\
 &= \sum_{k=2}^{n-b+1} \mathbb{P}(J_\Delta = k \mid \mathcal{U} \cap \mathcal{E}) \frac{\mathbb{E}(e^{-\phi T^{[k]}} T_k e^{-L_n \mu/2})}{\mathbb{E}(T_k e^{-L_n \mu/2})}, \quad (8.2.3)
 \end{aligned}$$

where $T^{[k]}$ is defined in (8.2.2).

Proof. This follows from the previous steps and (8.1.7).

Remark. The representation in (8.2.3) provides a useful way to simulate observations from T_Δ ; this is exploited later. Note that the original random variables \mathbf{T} can be tilted by the size-biasing function $e^{-L_n \mu/2}$, so that

$$\mathbb{E}_\mu f(T_n, \dots, T_2) = \frac{\mathbb{E}(f(T_n, \dots, T_2) e^{-L_n \mu/2})}{\mathbb{E}(e^{-L_n \mu/2})}.$$

In what follows we refer to this as μ -biasing. The previous results can then be written in terms of these μ -biased times:

$$\mathbb{P}_\mu(J_\Delta = k \mid \mathcal{U} \cap \mathcal{E}) = \frac{k(k-1) \binom{n-k}{b-1} \mathbb{E}_\mu T_k}{\sum_{l=2}^{n-b+1} l(l-1) \binom{n-l}{b-1} \mathbb{E}_\mu T_l}, \quad (8.2.4)$$

and

$$\mathbb{E}_\mu(e^{-\phi T_\Delta} \mid \mathcal{U} \cap \mathcal{E}) = \sum_{k=2}^{n-b+1} \mathbb{P}_\mu(J_\Delta = k \mid \mathcal{U} \cap \mathcal{E}) \frac{\mathbb{E}_\mu(e^{-\phi T^{[k]}} T_k)}{\mathbb{E}_\mu T_k}. \quad (8.2.5)$$

8.3 The case $\mu = 0$

It is of great interest in practice to consider the limiting case in which the mutation rate at the special locus is extremely small. In this case rather more can be said about the age of a neutral mutation. An immediate specialization of Theorem 8.1 provides a proof of Griffiths and Tavaré’s (1998) result:

Lemma 8.2 *The Laplace transform of the conditional distribution of the age T_Δ of a UEP observed b times in a sample of size n has limit as $\mu \rightarrow 0$ given by*

$$\mathbb{E}(e^{-\phi T_\Delta} \mid \mathcal{U} \cap \mathcal{E}) = \frac{\sum_{k=2}^{n-b+1} k(k-1) \binom{n-k}{b-1} \mathbb{E}(e^{-\phi T^{[k]}} T_k)}{\sum_{k=2}^{n-b+1} k(k-1) \binom{n-k}{b-1} \mathbb{E} T_k}. \quad (8.3.1)$$

This result provides the distribution of T_Δ in reasonably explicit form. If we define

$$S_k = T_n + \cdots + T_k,$$

then

$$\begin{aligned} \mathbb{E}(T_k e^{-\phi(UT_k + T_{k+1} + \cdots + T_n)}) &= \mathbb{E} \left[\int_0^1 T_k e^{-\phi u T_k} du e^{-\phi(T_{k+1} + \cdots + T_n)} \right] \\ &= \mathbb{E} \left[\phi^{-1} (1 - e^{-\phi T_k}) e^{-\phi(T_{k+1} + \cdots + T_n)} \right] \\ &= \mathbb{E} \left[\phi^{-1} e^{-\phi(T_{k+1} + \cdots + T_n)} - \phi^{-1} e^{-\phi(T_k + \cdots + T_n)} \right] \\ &= \int_0^\infty e^{-\phi t} \{ \mathbb{P}(S_{k+1} \leq t) - \mathbb{P}(S_k \leq t) \} dt \\ &= \int_0^\infty e^{-\phi t} \mathbb{P}(A_n(t) = k) dt, \end{aligned}$$

the last equality following from the fact that the ancestral process $A_n(t) = k$ if, and only if, $S_k > t$ and $S_{k+1} \leq t$. Hence we have

Theorem 8.3 *Assuming the times T_j have continuous distributions, the density of the age T_Δ is given by*

$$f_\Delta(t) = \frac{\sum_{k=2}^{n-b+1} k(k-1) \binom{n-k}{b-1} \mathbb{P}(A_n(t) = k)}{\sum_{k=2}^{n-b+1} k(k-1) \binom{n-k}{b-1} \mathbb{E} T_k}, \quad t > 0. \quad (8.3.2)$$

Moments of T_Δ can be found in a similar way, and one obtains

$$\mathbb{E}(T_\Delta^j) = \frac{\sum_{k=2}^n k(k-1) \binom{n-k}{b-1} \frac{1}{j+1} \mathbb{E}(S_k^{j+1} - S_{k+1}^{j+1})}{\sum_{k=2}^n k(k-1) \binom{n-k}{b-1} \mathbb{E}(T_k)}, j = 1, 2, \dots, \quad (8.3.3)$$

from which the mean and variance of T_Δ can be obtained. For example, in the constant population size case, we obtain

$$\mathbb{E}(T_\Delta) = 2 \binom{n-1}{b}^{-1} \sum_{j=2}^n \binom{n-j}{b-1} \frac{n-j+1}{n(j-1)}. \quad (8.3.4)$$

The age of an allele in the population

To derive the population version of (8.3.3), we assume that $\{A_n(t), t \geq 0\}$ converges in distribution to a process $\{A(t), t \geq 0\}$ as $n \rightarrow \infty$, and that the time taken for $A(\cdot)$ to reach 1 is finite with probability 1. Then as $n \rightarrow \infty$, and $b/n \rightarrow x$, $0 < x < 1$, we see that

$$\mathbb{E}(T_\Delta^j) = \frac{\sum_{k=2}^\infty k(k-1)(1-x)^{k-2} \frac{1}{j+1} \mathbb{E}(S_k^{j+1} - S_{k+1}^{j+1})}{\sum_{k=2}^\infty k(k-1)(1-x)^{k-2} \mathbb{E}(T_k)}, j = 1, 2, \dots \quad (8.3.5)$$

In this population limit the density of the age of a mutant gene that has a relative frequency x is, from Theorem (8.2),

$$\begin{aligned} g_x(t) &= \frac{\sum_{k=2}^\infty k(k-1)(1-x)^{k-2} \mathbb{P}(A(t) = k)}{\sum_{k=2}^\infty k(k-1)(1-x)^{k-2} \mathbb{E}(T_k)} \\ &= \frac{\mathbb{E}(A(t)(A(t)-1)(1-x)^{A(t)-2})}{\sum_{k=2}^\infty k(k-1)(1-x)^{k-2} \mathbb{E}(T_k)}. \end{aligned} \quad (8.3.6)$$

The mean age of the mutation known to have frequency x in the population follows from (8.3.4) by letting $n \rightarrow \infty$, $b/n \rightarrow x$:

$$\mathbb{E}(T_\Delta) = \frac{-2x}{1-x} \log x. \quad (8.3.7)$$

Equation (8.3.4) is the well known formula derived by Kimura and Ohta (1973). The density (8.3.6) is also known in various forms (e.g. Watterson (1977) and Tavaré (1984)).

Remark. There have been numerous papers written about the ages of alleles over the years, mostly using diffusion theory and reversibility arguments. This section sets the problem in a coalescent framework (although the results are much more general than they seem!). Watterson (1996) discusses Kimura’s contribution to this problem. A modern perspective is given by Slatkin and Rannala (2000).

8.4 Simulating the age of an allele

An alternative to the analytical approach is to simulate observations from the joint conditional distribution of those features of the process that are of interest, for example the age T_Δ of the mutation Δ , and the time $T_{MRCA\Delta}$ to the MRCA of the individuals carrying Δ . In order to simulate such times, we can use the following algorithm based on Theorem 8.2.3 and (8.1.7).

Algorithm 8.1 To simulate from conditional distribution of T_Δ and $T_{MRCA\Delta}$.

1. Choose k according to the distribution of J_Δ in (8.1.7).
2. Choose j_1 from the conditional distribution of J_1 given $J_\Delta = k$ in (8.1.9) with $i = 1$.
3. Simulate an observation from the (unconditional) μ -biased joint distribution of the coalescence times T_n, \dots, T_{k+1} .
4. Conditional on the results of step 3, simulate from the random variable Z having the (standard) size-biased distribution of T_k and set $T^* = UZ$, where U is an independent $U(0,1)$ random variable.
5. Return $T_{MRCA\Delta} = T_n + \dots + T_{j_1+1}$, $T_\Delta = T_{MRCA\Delta} + T_{j_1} + \dots + T_{k+1} + T^*$.

Remark. Generating the appropriate size-biased distributions can be difficult when the population size varies. Another way to implement this is to replace steps 3 and 4 above with a rejection step:

- 3'. Generate $\mathbf{T} = (T_n, \dots, T_2)$ from the coalescent model, and compute $L_n = 2T_2 + \dots + nT_n$. Accept \mathbf{T} with probability

$$\frac{T_k \mu}{2} e^{-T_k \mu / 2} e^{-L_n \mu / 2}; \quad (8.4.1)$$

otherwise repeat.

- 4'. Set $T^* = UT_k$, where U is an independent $U(0,1)$ random variable.
- 5'. Return $T_{MRCA\Delta} = T_n + \dots + T_{j_1+1}$, $T_\Delta = T_{MRCA\Delta} + T_{j_1} + \dots + T_{k+1} + T^*$.

The extra factor of e comes from the fact that $\text{Po}(T_k \mu / 2)\{1\} \leq \text{Po}(1)\{1\}$. In the limiting case $\mu = 0$ an independence sampler can be used.

8.5 Using intra-allelic variability

Rannala and Slatkin (1997) discussed a method for estimating the age of an allele known to have frequency b in a sample of size n , given an estimate of the number of mutations, m , that have arisen in the region around the mutation locus. There are at least three versions of this problem, depending on where these new mutations are assumed to occur. For example, we might sequence in the region of the mutation Δ and find the number of additional segregating sites in the region. We suppose once more that these additional mutations occur at rate $\theta/2$ on the branches of the coalescent tree.

If one wants to simulate observations from the posterior distribution of trees and times conditional on the number m of segregating sites appearing in the b individuals carrying the mutation in a region completely linked to Δ , then a modification of Algorithm 8.1 can be used:

Algorithm 8.2 To simulate from conditional distribution of age of mutation and $T_{MRC\Delta}$ given m additional segregating sites in the Δ subtree.

1. Choose k according to the distribution of J_Δ in (8.1.7).
2. Choose j_1, j_2, \dots, j_{b-1} from the conditional distribution of J_1, J_2, \dots, J_{b-1} given $J_\Delta = k$ in (8.1.8).
3. Simulate an observation from the (unconditional) joint distribution of the coalescence times T_n, \dots, T_{k+1} , and use the indices in step 2 to compute the coalescence times T_b^*, \dots, T_2^* in the Δ -subtree, together with the length $L_{nb} = \sum_{j=2}^b jT_j^*$ of the Δ -subtree.
4. Accept these statistics with probability

$$\text{Po}(\theta L_{nb}/2)\{m\}/\text{Po}(m)\{m\},$$

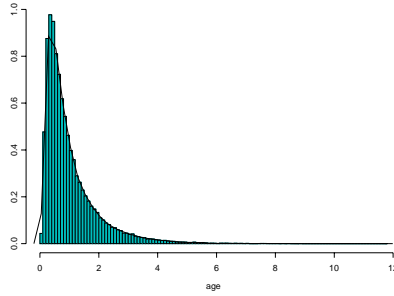
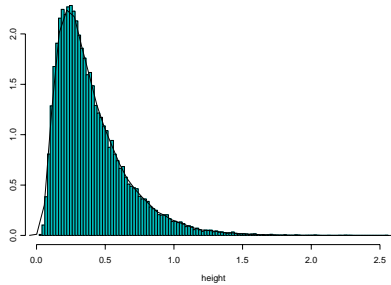
else return to step 1.

5. Conditional on the results of step 3, simulate from the random variable Z having the size-biased distribution of T_k and set $T^* = UZ$, where U is an independent $U(0,1)$ random variable.
6. Return $T_{MRC\Delta} = T_n + \dots + T_{j_1+1}$, $T_\Delta = T_{MRC\Delta} + T_{j_1} + \dots + T_{k+1} + T^*$.

Example

The conditional distribution of T_Δ and $T_{MRC\Delta}$ in the constant population size case were simulated using 50,000 runs of Algorithm 8.2 for the case $n = 200, b = 30, \theta = 4.0$ and $m = 5$ segregating sites observed in the subtree. The mean age was 1.01 with standard deviation 0.91, while the mean subtree height was 0.40 with a standard deviation of 0.25. Percentiles of the distributions are given below, together with the estimated densities. For further details and alternative simulation algorithms, see Griffiths and Tavaré (2003).

	2.5%	25%	50%	75%	97.5%
age	0.156	0.412	0.721	1.289	3.544
subtree height	0.099	0.218	0.334	0.514	1.056

Fig. 8.2. Density of age of mutation.**Fig. 8.3.** Density of height of subtree.

9 Markov Chain Monte Carlo Methods

In this section we introduce some models for DNA sequence data, and explore some computer intensive methods that can be used to estimate population parameters. The main inference technique discussed here is Markov chain Monte Carlo, introduced into this field by Kuhner *et al.* (1995, 1998).

We assume that mutations occur on the coalescent tree of the sample at rate $\theta/2$, independently in each branch of the tree. Here we study the case in which the type space E is finite, and we suppose that the mutation process is determined by

$$\gamma_{ij} = \mathbb{P}(\text{mutation results in type } j \mid \text{type was } i)$$

We write $\Gamma = (\gamma_{ij})$, and we note that γ_{ii} may be non-zero.

9.1 K -Allele models

One of the first models studied in any depth in this subject was the so-called K -allele model, in which $E = \{A_1, \dots, A_K\}$ corresponding to K possible alleles in the type space. Let $X_i(t)$ denote the fraction of the population that has allele A_i at time t . Many of the results concern the diffusion model for the process $\{(X_1(t), \dots, X_K(t)), t \geq 0\}$ with mutations determined according to the transition matrix Γ . The state space of the process is $\{\mathbf{x} = (x_1, \dots, x_K) \in [0, 1]^K : \sum_1^K x_i = 1\}$ and its generator has the form

$$L = \frac{1}{2} \sum_{i=1}^K \sum_{j=1}^K x_i(\delta_{ij} - x_j) \frac{\partial^2}{\partial x_i \partial x_j} + \sum_{j=1}^K \left(\sum_{i=1}^K x_i r_{ij} \right) \frac{\partial}{\partial x_j},$$

where

$$R = (r_{ij}) = \frac{\theta}{2}(\Gamma - I).$$

When the distribution of the type of a mutant is independent of its parental type, so that

$$\gamma_{ij} = \pi_j, \quad j \in E$$

where $\pi_j > 0, \sum_{j \in E} \pi_j = 1$, we recover the process studied in Section 3.1. The stationary distribution π of the diffusion is the Dirichlet distribution

$$\pi(x_1, \dots, x_K) = \frac{\Gamma(\theta)}{\Gamma(\theta\pi_1) \cdots \Gamma(\theta\pi_K)} x_1^{\theta\pi_1-1} \cdots x_K^{\theta\pi_K-1}. \tag{9.1.1}$$

Surprisingly perhaps, the distribution is known for essentially no other mutation matrices Γ . Suppose now that we take a sample of n genes from the stationary process with frequencies (X_1, \dots, X_K) . The sample comprises n_i genes of type $i, 1 \leq i \leq K$. Writing $\mathbf{n} = (n_1, \dots, n_K)$, the probability $q(\mathbf{n})$ that the sample has configuration \mathbf{n} is

$$q(\mathbf{n}) = \mathbb{E} \frac{n!}{n_1! \cdots n_K!} X_1^{n_1} \cdots X_K^{n_K}. \tag{9.1.2}$$

For the model (9.1.1), this gives

$$\begin{aligned} q(\mathbf{n}) &= \int \cdots \int \frac{n!}{n_1! \cdots n_K!} x_1^{n_1} \cdots x_K^{n_K} \pi(x_1, \dots, x_K) dx_1 \cdots dx_{K-1} \\ &= \frac{n! \Gamma(\theta) \Gamma(\theta\pi_1 + n_1) \cdots \Gamma(\theta\pi_K + n_K)}{n_1! \cdots n_K! \Gamma(\theta\pi_1) \cdots \Gamma(\theta\pi_K) \Gamma(\theta + n)} \\ &= \binom{\theta + n - 1}{n}^{-1} \prod_{j=1}^K \binom{\theta\pi_j + n_j - 1}{n_j}. \end{aligned} \tag{9.1.3}$$

In particular, the mean number of type i in the sample of size n is

$$\mathbb{E}(\text{number of allele } A_i) = n\mathbb{E}X_i = n\pi_i.$$

It is worth pointing out that a sample from any two-allele model can be described by (9.1.3), possibly after rescaling θ and Γ . To see this, suppose the matrix Γ has the form

$$\Gamma = \begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix}$$

Then the stationary distribution is $\boldsymbol{\pi} = \left(\frac{\beta}{\alpha+\beta}, \frac{\alpha}{\alpha+\beta}\right)$. Hence

$$\begin{aligned} R &\equiv \frac{\theta}{2}(\Gamma - I) \\ &= \frac{\theta}{2} \left(\begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix} - \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) \\ &= \frac{\theta}{2} \begin{pmatrix} -\alpha & \alpha \\ \beta & -\beta \end{pmatrix} \\ &= \frac{\theta}{2}(\alpha + \beta) \begin{pmatrix} -\frac{\alpha}{\alpha+\beta} & \frac{\alpha}{\alpha+\beta} \\ \frac{\beta}{\alpha+\beta} & -\frac{\beta}{\alpha+\beta} \end{pmatrix} \\ &= \frac{\theta}{2}(\alpha + \beta) \left(\begin{pmatrix} \frac{\beta}{\alpha+\beta} & \frac{\alpha}{\alpha+\beta} \\ \frac{\beta}{\alpha+\beta} & \frac{\alpha}{\alpha+\beta} \end{pmatrix} - \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right). \end{aligned}$$

We may therefore use the sampling formula (9.1.3) with $\theta\pi_1$ replaced by $\theta\beta$, and $\theta\pi_2$ replaced by $\theta\alpha$.

The number of real mutations

Suppose that mutations occur at rate $\nu/2$ on the coalescent tree (the switch from θ to ν will be explained shortly). At any mutation point, the current allele is changed according to the transition matrix Γ . We note that not all potential substitutions have to result in changes to the existing allele, as $\gamma_{jj} > 0$ is allowed. The *effective mutation rate* $\theta/2$ is defined to be the expected number of mutations per unit time that result in a change of allele:

$$\frac{\theta}{2} = \frac{\nu}{2} \sum_{j=1}^K \pi_j (1 - \gamma_{jj}), \quad (9.1.4)$$

where $\pi_j, j = 1, \dots, K$ denotes the stationary distribution of Γ .

Felsenstein's model

It is convenient to describe here one useful model for the case $K = 4$, corresponding to models for the base at a given site in a DNA sequence. Here, $E = \{A, G, C, T\}$. Because many of our applications focus on mitochondrial

DNA, in which transitions occur with much higher frequency than transversions, we use a model which allows for transition-transversion bias.

Suppose then that mutations arise at rate $\nu/2$. When a potential substitution occurs, it may be one of two types: *general*, in which case an existing base j is substituted by a base of type k with probability π_k , $1 \leq j, k \leq 4$; or *within-group*, in which case a pyrimidine is replaced by C or T with probability proportional to π_C and π_T respectively, and a purine is replaced by A or G with probability proportional to π_A and π_G respectively. The conditional probability of a general mutation is defined to be $1/(1 + \kappa)$, while the conditional probability of a within-group mutation is defined to be $\kappa/(1 + \kappa)$, where $\kappa \geq 0$ is the transition-transversion parameter. Thus the mutation matrix Γ is given by

$$\Gamma = \frac{1}{1 + \kappa} \Gamma_1 + \frac{\kappa}{1 + \kappa} \Gamma_2, \tag{9.1.5}$$

where $\Gamma_{1,ij} = \pi_j$, $j \in E$ and

$$\Gamma_2 = \begin{pmatrix} \frac{\pi_A}{\pi_A + \pi_G} & \frac{\pi_G}{\pi_A + \pi_G} & 0 & 0 \\ \frac{\pi_A}{\pi_A + \pi_G} & \frac{\pi_G}{\pi_A + \pi_G} & 0 & 0 \\ 0 & 0 & \frac{\pi_C}{\pi_C + \pi_T} & \frac{\pi_C}{\pi_C + \pi_T} \\ 0 & 0 & \frac{\pi_C}{\pi_C + \pi_T} & \frac{\pi_C}{\pi_C + \pi_T} \end{pmatrix}$$

In Γ_1 and Γ_2 , the states are written in order A, G, C, T . It is readily checked that the stationary distribution of Γ is $\boldsymbol{\pi} = (\pi_A, \pi_C, \pi_G, \pi_T)$. If we define

$$g = \frac{\nu}{2(1 + \kappa)}, \quad w = \kappa g, \tag{9.1.6}$$

then κ is the ratio of the within-class to general substitution rates. From (9.1.4), the effective mutation rate is given by

$$\frac{\theta}{2} = g \left(1 - \sum_{j \in E} \pi_j^2 \right) + 2w \left(\frac{\pi_A \pi_G}{\pi_A + \pi_G} + \frac{\pi_C \pi_T}{\pi_C + \pi_T} \right) \tag{9.1.7}$$

The transition matrix e^{Rt} of the mutation process with transition intensity matrix $R = \nu(\Gamma - I)/2$ is known. We denote the jk -th element by $r_{jk}(t)$; this is the probability that a base of type j has changed to a base of type k a time t later. Thorne *et al.* (1992) show that

$$r_{jk}(t) = \begin{cases} e^{-(g+w)t} + e^{-gt} (1 - e^{-wt}) \frac{\pi_k}{\pi_{H(k)}} + (1 - e^{-gt}) \pi_k & j = k \\ e^{-gt} (1 - e^{-wt}) \frac{\pi_k}{\pi_{H(k)}} + (1 - e^{-gt}) \pi_k, & H(j) = H(k) \\ (1 - e^{-gt}) \pi_k & H(j) \neq H(k) \end{cases}$$

where $\pi_R = \pi_A + \pi_G$, $\pi_Y = \pi_C + \pi_T$, and $H(i)$ denotes whether base i is a purine or a pyrimidine, so that $H(A) = H(G) = R$ and $H(C) = H(T) = Y$.

9.2 A biomolecular sequence model

Of particular interest to us is the case in which the types represent DNA or protein sequences of length s , say. Then the type space E has the form $E = E_0^s$, where E_0 is the type space of a single position, or site, in the sequence. The sites of the sequence may be labeled in many ways. The DNA alphabet $E_0 = \{A, C, G, T\}$ is one possibility, as is the 20 letter amino-acid sequence alphabet, or the 64 letter codon alphabet. Also common are the purine-pyrimidine alphabet, where $E_0 = \{Y, R\}$ and $Y = \{A, G\}$ denotes purines, $R = \{C, T\}$ the pyrimidines. In many evolutionary studies, transversions are not observed, and it might then be natural to think of sites as being binary, with $E_0 = \{A, G\}$ or $E_0 = \{C, T\}$. There are many possible models for the mutation process Γ , depending on what is assumed about the effects of mutation. Here we suppose that when a mutation occurs, it results in a substitution, the replacement of one element of E_0 by another one. The simplest version of this model supposes that the substitution occurs at site j with probability h_j , where

$$h_j \geq 0, \quad \sum_{j=1}^s h_j = 1. \quad (9.2.1)$$

The h_j are identical (and so equal to $1/s$) if there are no mutational hotspots, and h_j may be 0 if site j is invariable. Thus the h_j add some flexibility in modeling variable mutation rates across the sequences. A mutation occurring at site j produces substitutions according to transition matrix $P_j = (p_{lm}^{(j)})$. Thus substitutions change a sequence of type (i_1, \dots, i_s) to one of type (j_1, \dots, j_s) as follows:

$$(i_1, \dots, i_s) \rightarrow (i_1, \dots, i_{l-1}, j_l, i_{l+1}, \dots, i_s)$$

with probability $h_l p_{i_l j_l}^{(l)}$, $1 \leq l \leq s$. We may write Γ in the form

$$\Gamma = \sum_{l=1}^s h_l I \otimes \cdots \otimes I \otimes P_l \otimes I \otimes \cdots \otimes I \quad (9.2.2)$$

where I denotes the identity matrix, and \otimes denotes direct (or Kronecker) product: $A \otimes B = (a_{ij} B)$. Recall that if A, B, C, D are conformable matrices, then $(A \otimes B)(C \otimes D) = AC \otimes BD$. If π_l denotes the stationary distribution of P_l , and π denotes the stationary distribution of Γ , then it is easy to show that $\pi = \pi_1 \otimes \cdots \otimes \pi_s$.

Many properties of this process may be studied using the coalescent simulation described in Section 6.6. The previous result shows that for simulating sequences from a stationary population, the ancestral sequence may be generated by simulating independently at each site, according to the stationary distribution of each site.

9.3 A recursion for sampling probabilities

Return now to the K -allele model with mutation matrix $\Gamma = (\gamma_{ij})$, and $R = \frac{\theta}{2}(\Gamma - I)$. Let $q(\mathbf{n})$ be the probability that a sample of n genes has a type configuration of $\mathbf{n} = (n_1, \dots, n_K)$, and define $[K] = \{1, 2, \dots, K\}$. A fundamental recursion is given in

Theorem 9.1

$$\begin{aligned}
 q(\mathbf{n}) = & \frac{\theta}{n + \theta - 1} \left(\sum_{i=1}^K \frac{n_i}{n} \gamma_{ii} q(\mathbf{n}) + \sum_{i,j \in [K], n_j > 0, i \neq j} \frac{n_i + 1}{n} \gamma_{ij} q(\mathbf{n} + \mathbf{e}_i - \mathbf{e}_j) \right) \\
 & + \frac{n - 1}{n + \theta - 1} \sum_{j \in [K], n_j > 0} \frac{n_j - 1}{n - 1} q(\mathbf{n} - \mathbf{e}_j), \tag{9.3.1}
 \end{aligned}$$

where $\{\mathbf{e}_i\}$ are the K unit vectors. Boundary conditions are required to determine the solution to (9.3.1). These have the form

$$q(\mathbf{e}_i) = \pi_i^*, \quad i = 1, \dots, K, \tag{9.3.2}$$

where π_i^* is the probability that the most recent common ancestor is of type i .

Proof. To derive (9.3.1) consider the first event back in time that happened in the ancestral tree. Relative rates of mutation and coalescence for n genes are $n\theta/2 : n(n - 1)/2$, so the probability that the first event is a mutation is $\theta/(n + \theta - 1)$. To obtain a configuration of \mathbf{n} after mutation the configuration before must be either \mathbf{n} , and a transition $i \rightarrow i$ takes place for some $i \in [K]$ (the mutation resulted in no observable change), or $\mathbf{n} + \mathbf{e}_i - \mathbf{e}_j$, $i, j \in [K]$, $n_j > 0$, $i \neq j$ and a transition $i \rightarrow j$ take place. If a coalescence was the first event back in time, then to obtain a configuration \mathbf{n} the configuration must be $\mathbf{n} - \mathbf{e}_j$ for some $j \in [K]$ with $n_j > 0$ and the ancestral lines involved in the coalescence must be of type j . \square

The recursion in (9.3.1) is on n , the sample size. Given $\{q(\mathbf{m}); m < n\}$, simultaneous equations for the $\binom{n+K-1}{K-1}$ unknown probabilities $\{q(\mathbf{m}); m = n\}$ are non-singular, and in theory can be solved; cf. Lundstrom (1990). It is common to assume that

$$\pi_i^* = \pi_i, \quad i = 1, \dots, K, \tag{9.3.3}$$

where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$ is the stationary distribution of Γ . With this assumption, $q(\mathbf{n})$ is the stationary sampling distribution.

It is worth emphasizing that the probability $q(\mathbf{n})$ satisfying (9.3.1) is determined solely by the rate matrix R . Indeed, (9.3.1) can be rewritten in the form

$$\begin{aligned}
q(\mathbf{n}) = & \\
& \frac{2}{n(n-1)} \left(\sum_{i=1}^K n_i r_{ii} q(\mathbf{n}) + \sum_{i,j \in [K], n_j > 0, i \neq j} (n_i + 1) r_{ij} q(\mathbf{n} + \mathbf{e}_i + \mathbf{e}_j) \right) \\
& + \frac{1}{n-1} \sum_{j \in [K], n_j > 0} (n_j - 1) (\mathbf{n} - \mathbf{e}_j).
\end{aligned}$$

The point here is that different combinations of θ and Γ can give rise to the same R matrix. Nonetheless, we prefer to think of the model in terms of an overall rate θ and a matrix of substitution probabilities Γ . In practice, we often assume that Γ is known, and the aim might then be to estimate the single parameter θ , which reflects both the effective population size N and the mutation probability u .

Remark. The recursion in (9.3.1) has appeared in a number of guises in the literature, such as Sawyer *et al.* (1987) and Lundstrom *et al.* (1992). In the latter references, a quasi-likelihood approach for estimation of θ in the finitely-many-sites model is developed. The recursion (9.3.1) is used to find the probability distribution at each site, and the quasi-likelihood is computed by assuming independence across the sites.

Griffiths and Tavaré (1994) used the recursion for the finitely-many-sites model to find the likelihood. Conventional numerical solutions in this case are difficult to obtain because of the large number of equations. This prompted them to develop their Markov chain approach. See Forsythe and Leibler (1950) for an early application of Monte Carlo approaches to matrix inversion. We note here that early experience with the Griffiths-Tavaré method suggests it is not feasible for analyzing large amounts of sequence data. In the remainder of this section, we discuss a Markov chain Monte Carlo approach and give a number of examples of its use.

9.4 Computing probabilities on trees

For definiteness, assume we are dealing with DNA sequence data \mathcal{D} having s aligned sites in a sample of size n . We will use Λ to denote the (labeled) coalescent tree topology, and $\mathbf{T} = (T_2, \dots, T_n)$ to denote the coalescence times in the tree. For a given model of substitution at a particular site in the sequence, we will need to compute the probability of the bases in the sample, given a particular value of Λ and \mathbf{T} . This can be done using a recursive method, known as the *peeling algorithm*, described by Felsenstein (1973, 1981). The idea is to compute the probability of the bases b_1, \dots, b_n observed at a particular position in sequences $1, \dots, n$. Each node l in the tree is assigned a vector of length 4, the i -th entry of which gives the probability of the data below that node, assuming node l is base i . The algorithm is initialized by assigning the vector associated with a leaf i the vector with elements $\delta_{b_i, j}$, $j = 1, \dots, 4$.

The calculation now proceeds recursively. Imagine that the probability vectors (w_{u1}, \dots, w_{u4}) and (w_{v1}, \dots, w_{v4}) have been computed for the descendant nodes u and v respectively of node l . To compute the vector (w_{l1}, \dots, w_{l4}) at node l , we need to calculate the time t_{lu} along the branch from $l \rightarrow u$, and the time t_{lv} from $l \rightarrow v$. Then we calculate

$$w_{lz} = \left(\sum_x r_{zx}(t_{lu})w_{ux} \right) \cdot \left(\sum_y r_{zy}(t_{lv})w_{vy} \right),$$

where $r_{ij}(t)$ is the probability that base i has mutated to base j a time t later.

This scheme allows us to recurse up to the root of the tree. That node has label $l = 2n - 1$ and descendant nodes u and v . We finish the computation of the probability L of the configuration at that site by computing

$$L = \sum_z \pi_z^0 w_{uz} w_{vz}$$

where $\pi_z^0, z = 1, \dots, 4$ is the distribution of the ancestral base.

Once the likelihood at a single base position is calculated, the likelihood of the set of n sequences can be calculated using the fact that for the mutation model in Section 9.2 the sites evolve independently, conditional on Λ and \mathbf{T} . Hence if L_i denotes the likelihood of the i -th site, the overall likelihood is

$$\mathbb{P}(\mathcal{D} \mid \Lambda, \mathbf{T}) = \prod_{i=1}^s L_i. \quad (9.4.1)$$

9.5 The MCMC approach

Here we discuss a version of the Metropolis-Hastings algorithm, due originally to Metropolis *et al.* (1953) and Hastings (1970) that will be exploited for inference on coalescent trees. Our presentation follows that of Markovtsova (2000). The algorithm produces correlated samples from a posterior distribution π of interest, in our case $\pi(G) \equiv f(G \mid \mathcal{D})$, where $G \equiv (\Lambda, \mathbf{T}, M)$, M representing the mutation parameters and \mathcal{D} representing the sequence data. We use these samples to make inferences about parameters and statistics of interest. Examples include the effective mutation rate θ , the time to the most recent common ancestor, ages of a particular event in the sample, or population growth rates. We can write

$$f(G \mid \mathcal{D}) = \mathbb{P}(\mathcal{D} \mid G)g_1(\Lambda)g_2(\mathbf{T})g_3(M)/f(\mathcal{D}). \quad (9.5.1)$$

The first term on the right can be computed using the peeling algorithm described in the last section and an appropriate model for mutation among the sequences. The term $g_1(\Lambda)$ on the right of (9.5.1) is the coalescent tree topology distribution, $g_2(\mathbf{T})$ is the density of the coalescence times \mathbf{T} , and $g_3(M)$ is the

prior distribution for the mutation parameters M . The normalizing constant $f(\mathcal{D})$ is unknown and hard to compute. The algorithm starts with an arbitrary choice of Λ , \mathbf{T} and M . New realizations of G are then proposed, and accepted or rejected, according to the following scheme.

Algorithm 9.1 Basic Metropolis-Hastings method:

1. Denote the current state by $G = (\Lambda, \mathbf{T}, M)$.
2. Output the current value of G .
3. Propose $G' = (\Lambda', \mathbf{T}', M')$ according to a kernel $Q(G \rightarrow G')$.
4. Compute the Hastings ratio

$$h = \min \left\{ 1, \frac{\pi(G')Q(G' \rightarrow G)}{\pi(G)Q(G \rightarrow G')} \right\}. \quad (9.5.2)$$

5. Accept the new state G' with probability h , otherwise stay at G .
6. Return to step 1.

Let $X(t)$ denote the state of this chain after t iterations. Once $X(t)$ has ‘reached stationarity’ its values represent samples from the distribution $\pi(G) = \pi(\Lambda, \mathbf{T}, M)$. The nature of the algorithm is such that consecutive outputs will be correlated. For many problems this might be not a bad thing, however one should be careful with using the output for calculating standard errors. But in some cases it is desirable to simulate approximately independent samples from the posterior distribution of interest, in which case we use output from every m^{th} iteration, for a suitable choice of m .

Current methods

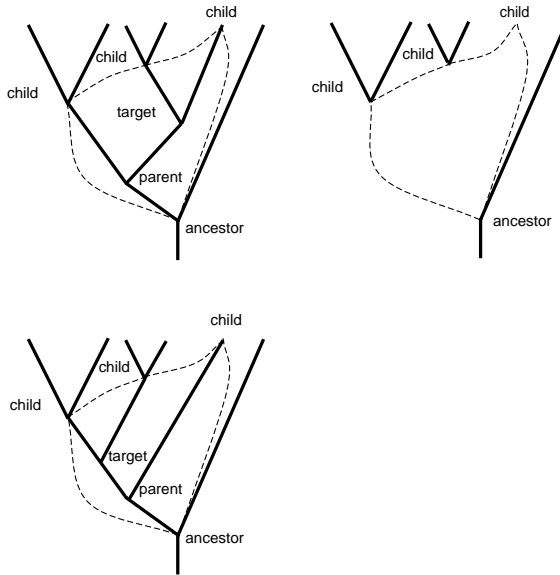
In this section we describe some methods of sampling genealogies. Most of these algorithms are very similar and often differ only in tree representation and useful tricks to speed up the computations. All of them start with an initial genealogy (random or UPGMA) and make small modifications to it. Choices among possible modifications may be random or deterministic.

The first is due to Kuhner *et al.* (1995). As before, the genealogy consists of two parts: the tree topology and a set of times between coalescent events, but time is rescaled in terms of the overall mutation rate in such a way that in one unit of time the expected number of mutations per site is 1. Figure 9.1 shows the updating process: choosing a neighborhood (the region of genealogy to be changed), rearranging the topology in that neighborhood, and choosing new branch lengths within the neighborhood. This fundamental operation is applied repeatedly. To make rearrangements, a node is chosen at random from among all nodes that have both parents and children (i.e., are neither leaves nor the bottom-most node of the genealogy). This node is referred to as the target. The neighborhood of rearrangement consists of the target node, its

child, parent, and parent's other child. A rearrangement makes changes of two types: reassorts the tree children among target and parent, and modifies the branch length within the neighborhood. The lineages to be redrawn are referred to as active lineages, and the lineages outside of the neighborhood as inactive lineages.

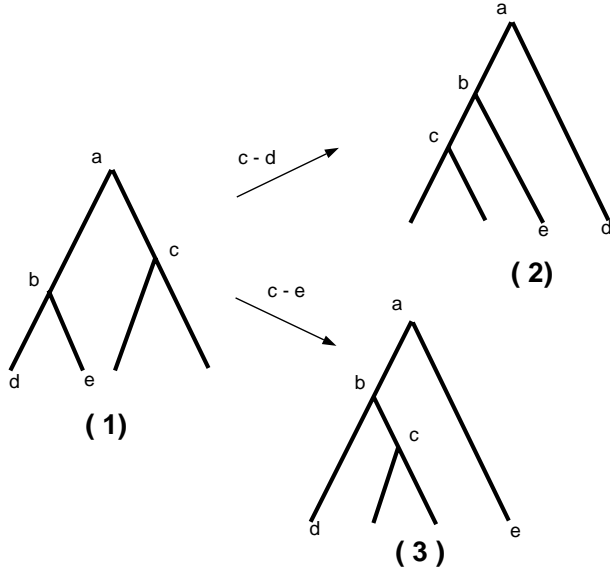
The times of the target and parent nodes are drawn from a conditional coalescent distribution with the given mutation rate, conditioned on the number of inactive lineages. For each time interval, the probability of coalescence among the active lineages depends on the number of active and inactive lineages present in the genealogy during that time interval. A random walk, weighted by these probabilities, is used to select a specific set of times.

Fig. 9.1. Steps in rearranging a genealogy. Top left: selecting a neighborhood. Top right: erasing the active lineages. Bottom: redrawing the active lineages.



Yang and Rannala (1997) use a stochastic representation of the nearest neighbor interchange (NNI) algorithm as a core of the transition kernel. This algorithm generates two neighboring topologies for each interior branch (see Figure 9.2). Consider an interior branch $a - b$, where a is the ancestral node and b is the descendant node. Node c is the other descendant of a , and nodes d and e are descendants of b . The two neighbors of tree 1 are generated by interchanging node c with node d (tree 2), and node c with node e (tree 3).

Equal probabilities are assigned to each of the neighboring topologies. The NNI algorithm modifies the topology but ignores the ordering of the nodes

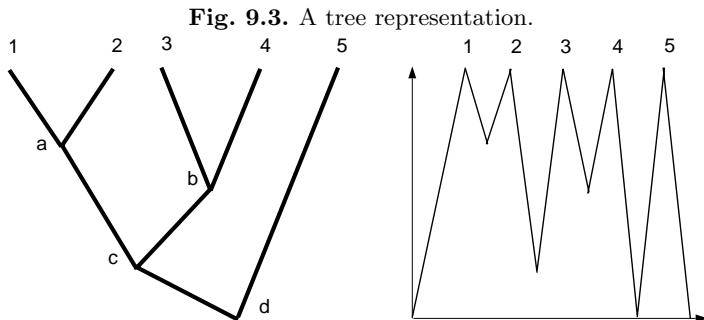
Fig. 9.2. NNI algorithm for a rooted binary tree topology.

(i.e., labeled history). To modify the NNI algorithm so that the chain moves between labeled histories, they assign an equal probability to each of the possible labeled histories for a nominated topology. This involves enumerating and recording all the labeled histories for that topology. The move to another labeled history that belongs to the current tree topology is allowed with the specified probability if the topology has more than one labeled history. Yang and Rannala use this transition kernel in the study of species data; the time to the MRCA is scaled to be 1 and times between speciation events have different distributions than those specified by the coalescent.

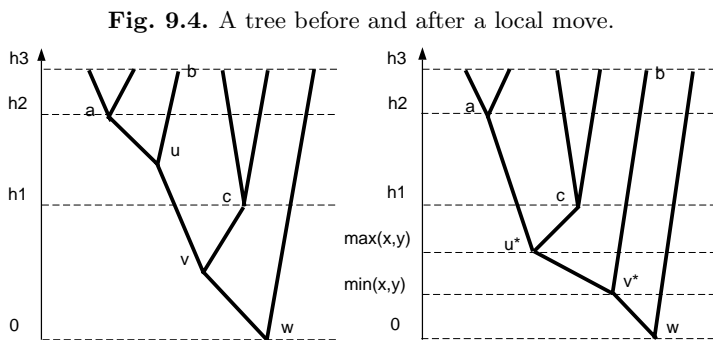
Wilson and Balding (1998) designed an algorithm to deal with microsatellite (or short tandem repeat) data. A step-wise model is chosen for the changes in repeat number at each mutation event. Although calculation of the likelihood via peeling is feasible for problems of moderate size, increasing the dimension of the parameter space by introducing the allelic state of the internal nodes permits much faster likelihood calculations. The algorithm uses a very simple method for generating candidate trees. It involves removing a branch from the tree at random and adding it anywhere in the tree, but locations close to similar allelic types are preferentially chosen.

Larget and Simon (1999) use an algorithm for moving in a tree space that is very close to the one developed by Mau *et al.* (1999). It uses the fact that for a given choice of ordering all sub-trees from left to right there is a unique in-order traversal of the tree. Each internal node is adjacent to two leaves in this traversal, the right-most leaf of its left sub-tree and the left most leaf

of its right sub-tree. Given the ordering of the nodes and distances between adjacent nodes, the tree topology and branch lengths are uniquely determined. Each taxon appears at a peak of the graph, and each internal node is a valley (see Figure 9.3).



The transition kernel consists of two different moves: global and local. For a global move one representation of the current tree is selected uniformly at random by choosing left/right orientation of the two sub-trees with equal probability for each internal node. Then the valley depths are simultaneously and independently modified by adding to each a perturbation in either direction, keeping the depth between 0 and a specified maximum. The local move modifies a tree only in a small neighborhood of a randomly chosen internal branch, leaving the remainder of the tree unchanged. Let u and v be the nodes joined by the randomly chosen edge (see Figure 9.4).



Leaving positions of a , b , c , and w fixed, new positions for nodes u and v are picked. Let $h_1 < h_2 < h_3$ be the distances between c and w , a and w , and

b and w correspondingly. In the local move, x is chosen uniformly at random from $[0, h_2]$, and y is chosen uniformly at random from $[0, h_1]$. Proposed nodes u^* and v^* will be distances $\max(x, y)$ and $\min(x, y)$ from w , respectively. If $\max(x, y) < h_1$, there are three possible tree topologies. One of the children, a , b , and c , is randomly chosen to be joined to v^* , with the others becoming children of u^* . If v is the root of the tree, the distances between v and the children a , b , and c are changed and the new location of u is chosen. The local move is very similar in character to the method of Kuhner *et al.* (1995).

9.6 Some alternative updating methods

We have some freedom in choosing the proposal kernel $Q(\cdot, \cdot)$. Ideally $Q(\cdot, \cdot)$ is relatively easy to calculate since the scheme above may need to iterated many times in order to converge to stationarity. Furthermore we have to demonstrate that the chain $X(t)$ satisfies the conditions of irreducibility and positive recurrence in order to show that the ergodic theorem applies and so the limiting distribution is indeed $f(A, \mathbf{T}, M \mid \mathcal{D})$.

We define level l of the genealogy to be the first point at which there are l distinct ancestors of the sample. The bottom of a genealogy of n individuals is therefore referred to as level n , and the MRCA of the sample is level 1. Recall that T_l denotes the time between levels l and $l - 1$. To propose a new graph (A', \mathbf{T}') we considered three different proposal kernels.

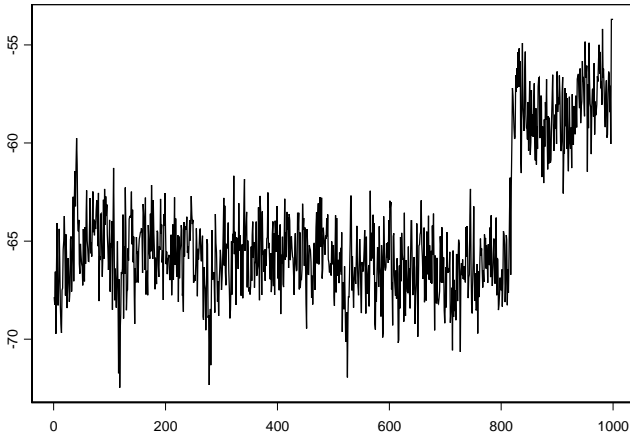
A bad sampler

Here is a simple algorithm:

1. Pick a level, l say ($l = n, n - 1, \dots, 2$), according to an arbitrary distribution F .
2. Delete upper part of the tree starting from level l .
3. Attach a new top of the tree generated according to the coalescent prior for a sample of l individuals.
4. Generate a new time T'_l , to replace the old T_l according to an exponential distribution with parameter $l(l - 1)/2$.

This algorithm works poorly, mainly because the suggested changes were too global. If we chose level l close to the bottom of the tree and attach a random top to it, then the new tree will be very different from the old one and has small chance of being accepted. As a result our sample will consists of trees with similar topologies and almost the same likelihood. But sometimes quite a different tree might be accepted and our Markov chain would move to other part of state space and stay there for long time. Figure 9.5 is an example of such a chain. This algorithm seems not to be very efficient in exploring the state space of trees.

The following algorithm looks simple and is easy to implement. It makes changes which are more local than the algorithm described above.

Fig. 9.5. Time series plot of log-likelihood

1. Pick a level, l say ($l = n, n - 1, \dots, 2$), according to an arbitrary distribution F .
2. Label the l lines $1, 2, \dots, l$.
3. Let L_i and L_j be the two lines which coalesce.
4. With probability $1/2$ replace this coalescence by one between L_i and a randomly chosen line (possibly resulting in the same topology as before).
5. Otherwise replace this coalescence by one between L_j and a randomly chosen line (also possibly resulting in the same topology as before).
6. Generate a new time T'_l , to replace the old T_l according to an exponential distribution with parameter $l(l - 1)/2$.

An example of a possible move, for a genealogy of five individuals, is shown in Figure 9.6.

This algorithm also does not work well, primarily because it is relatively hard to switch the order of two coalescence events. For example, we need several iterations of the algorithm to move from G to G' as illustrated in Figure 9.7.

Theoretically, this kernel has all the required properties, but it is simply not efficient. We might try other distributions for the choice of level l , or for the new time T'_l , but it is doubtful these would help. Our experience was that the algorithm became stuck in local maxima which required a re-ordering of coalescences in order to escape.

Fig. 9.6. A move in the sampler

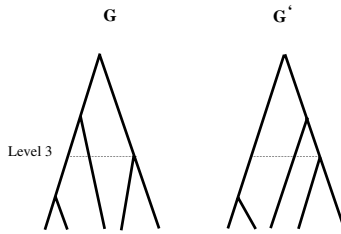
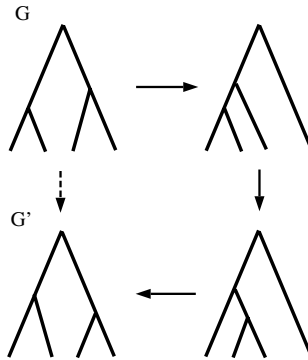


Fig. 9.7. Change of order of two coalescences



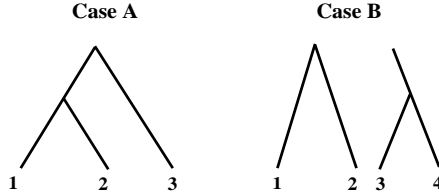
A good sampler

Lack of success with first two algorithms leads to the following approach, described in Markovtsova *et al.* (2000).

Algorithm 9.2 Local updating method.

1. Pick a level, l say ($l = n, n - 1, \dots, 3$), according to an arbitrary distribution F .
2. For the chosen l observe the pattern of coalescence at levels $l - 1$ and $l - 2$. This pattern falls into two cases, according to whether the coalescence at level $l - 2$ involves the line which results from the coalescence at level $l - 1$. These are illustrated in Figure 9.8. In Case A our kernel randomly generates a new topology involving the same three lines of ancestry; this new topology will also be Case A and may be the same topology with which we began. These are illustrated in Figure 9.9. In Case B we change

Fig. 9.8. Two possible coalescence patterns

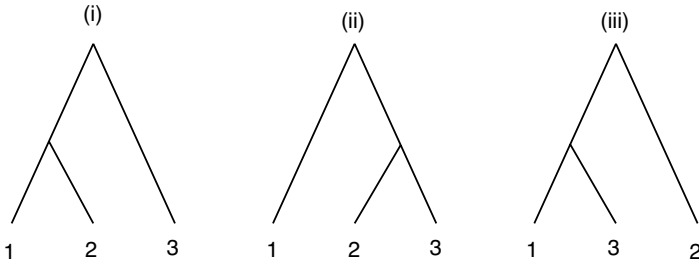


the order of the two coalescence events. So, for the example drawn above, we move to the state shown in Figure 9.10.

3. Generate new times T'_l and T'_{l-1} according to an arbitrary distribution, and leave other times unchanged. Thus we only alter the times corresponding to the levels at which the topology has been changed. This ensures that (A', T') is similar to (A, T) and therefore has a reasonable probability of being accepted.

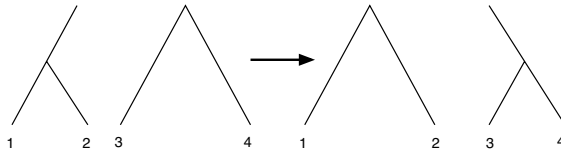
There are several variants of Step 2 of the above scheme. For example, one can allow the topology to remain the same in Case B, but not in Case A. We also tried a variant of Case B in which we proposed a new Case B topology uniformly from the six possible choices in which the four lines are paired randomly. None of these variations impacts significantly on the results.

Fig. 9.9. Possible moves in Case A



There are many possible choices for the updating times T'_l and T'_{l-1} . One might propose new values of T'_j from the pre-data coalescent distribution as it was done in first two algorithms. Second, one might generate times from a Normal distribution with mean equal to the currently accepted value T_j . We chose to truncate the Normal distribution in order to ensure that negative times were not proposed. The variances of the Normal distributions are parameters that can be tuned to get good mixing properties. Unfortunately,

Fig. 9.10. Possible moves in Case B



the optimal choice of variance appears to be highly data-dependent. In principle all choices are valid, but the rate of approach to stationarity, and the correlation between consecutive iterations, can vary significantly. The second approach might work better when trees are much shorter, or longer, than would be expected *a priori*.

Finally, we update the mutation parameter $M = (g)$ every k iterations. There are several ways to do it. First one is to propose new value g' from prior distribution. This updating mechanism works well in the case when the prior for g is very concentrated, i.e. a uniform with narrow support or a Normal distribution with small variance. This approach might be used when some external information is available. Second one is to generate new value g' according to truncated Normal distribution with mean g . The variance of this distribution requires some tuning to ensure well-behaved, i.e. uncorrelated, output. This approach works fine in case of uninformative prior or prior with wide support.

The Hastings ratio

Writing $G = (\Lambda, \mathbf{T}, M)$, the kernel Q can be expressed as the product of three terms:

$$Q(G \rightarrow G') = Q_1(\Lambda \rightarrow \Lambda') Q_2(\mathbf{T} \rightarrow \mathbf{T}' \mid \Lambda \rightarrow \Lambda') Q_3(M \rightarrow M').$$

Consequently the Hastings ratio can be written in the form

$$h = \min \left\{ 1, \frac{\mathbb{P}(\mathcal{D} \mid G') \frac{g_1(\Lambda') g_2(\mathbf{T}') g_3(M')}{\mathbb{P}(\mathcal{D} \mid G) \frac{g_1(\Lambda) g_2(\mathbf{T}) g_3(M)}} \times \frac{Q_1(\Lambda' \rightarrow \Lambda) Q_2(\mathbf{T}' \rightarrow \mathbf{T} \mid \Lambda' \rightarrow \Lambda) Q_3(M' \rightarrow M)}{Q_1(\Lambda \rightarrow \Lambda') Q_2(\mathbf{T} \rightarrow \mathbf{T}' \mid \Lambda \rightarrow \Lambda') Q_3(M \rightarrow M')} \right\}, \quad (9.6.1)$$

the unknown term $f(\mathbf{D})$ cancelling. We can further simplify (9.6.1) by noting that, since pairs of lines are chosen uniformly to coalesce, all topologies are, *a priori*, equally likely. Hence $g_1(\Lambda') = g_1(\Lambda)$. Furthermore, our transition

kernel changes only two of the times on the tree, T_l and T_{l-1} say. Finally, it is easy to show that $Q_1(A \rightarrow A') = Q_1(A' \rightarrow A)$, reducing (9.6.1) to

$$h = \min \left\{ 1, \frac{\mathbb{P}(\mathcal{D} \mid G') g_2(\mathbf{T}') g_3(M') f_l(t_l) f_{l-1}(t_{l-1}) Q_3(M' \rightarrow M)}{\mathbb{P}(\mathcal{D} \mid G) g_2(\mathbf{T}) g_3(M) f_l(t'_l) f_{l-1}(t'_{l-1}) Q_3(M \rightarrow M')} \right\}, \tag{9.6.2}$$

where $f_l(\cdot)$ and $f_{l-1}(\cdot)$ are the densities of the time updating mechanism at levels l and $l - 1$.

If one uses a transition kernel which proposes new times that are exponential with parameter $l(l - 1)/2$ at level l , (*i.e.* the unconditional coalescent distribution for times), then further cross-cancellation reduces (9.6.2) to

$$h = \min \left\{ 1, \frac{\mathbb{P}(\mathcal{D} \mid G') g_3(M') Q_3(M' \rightarrow M)}{\mathbb{P}(\mathcal{D} \mid G) g_3(M) Q_3(M \rightarrow M')} \right\}. \tag{9.6.3}$$

A similar simplification also follows if one proposes new mutation rates independently of the currently accepted rate and

$$h = \min \left\{ 1, \frac{\mathbb{P}(\mathcal{D} \mid G')}{\mathbb{P}(\mathcal{D} \mid G)} \right\}. \tag{9.6.4}$$

In order to test the algorithm for moving around tree space, we can use a simple mutation model for which there are alternative algorithms. One obvious choice is the infinitely-many-sites model, for which we have already developed some theory in Section 7. The data take the form of the number of segregating sites in the sample, and Algorithm 7.3 can be used to generate observations from the posterior distribution of features of the tree, conditional on the number of segregating sites observed.

9.7 Variable population size

The methods discussed in above can easily be adapted to model populations which are not of a fixed constant size. As in Section 2.4, let $N(t)$ denote the population size at time t , where time is measured in units of $N = N(0)$ generations, and write

$$N(t) = f(t)N(0), \quad \Lambda(t) = \int_0^t \frac{1}{f(u)} du.$$

If $A_n(t)$ is the ancestral process for a sample of size n evolving in a population of constant size, $A_n^v(t) = A_n(\Lambda(t))$ is the coalescent process appropriate for the population of varying size.

We let T_k^v record the coalescent time spent with k lines of descent in a growing population. The algorithm works by manipulating the underlying *coalescent* times, $\{T_i\}$, defined on the original coalescent time-scale, and subsequently transforming them to times in the varying population while calculating probability of data given tree.

Define $S_i = \sum_{j=i+1}^n T_j$. S_i represents the amount of standard coalescent time taken to get to a level with i lines of descent present. Similarly, $S_i^v = \sum_{j=i+1}^n T_j^v$ in the varying population. We transform the S_i to the S_i^v via $S_i^v = \min \{s : A(s) = S_i\}$. The proposal kernel works by manipulating the underlying coalescent times, $\{T_i\}$. Assuming we have picked level l in our updating step, new times T_l^v, T_{l-1}^v are proposed as follows. We begin by generating new times $T_l' = t_l'$ and $T_{l-1}' = t_{l-1}'$. Having done so, we recalculate S_k for all $k \leq l$. From these values we derive the new $\{S_i^v\}$, noting that $S_i^v' = S_i^v$ for $i > l$.

9.8 A Nuu Chah Nulth data set

We illustrate our approach with a sample of mitochondrial sequences from the Nuu Chah Nulth obtained by Ward *et al.* (1991). The data D are 360 bp sequences from region I of the control region obtained from a sample of $n = 63$ individuals. The observed base frequencies are $(\pi_A, \pi_G, \pi_C, \pi_T) = (0.3297, 0.1120, 0.3371, 0.2212)$. The data have 26 segregating sites and a mean heterozygosity of 0.0145 per site. There are 28 distinct haplotypes with a haplotype homozygosity of 0.0562. We fit two models to these data, both of which are variants of Felsenstein's model described in Section 9.2:

Model 1. All sites mutate at the same rate, so that $g_i \equiv g$ for all sites i . Here $M = (g, \kappa)$.

Model 2. The special case of Model 1 in which κ is assumed known, so that $M = (g)$.

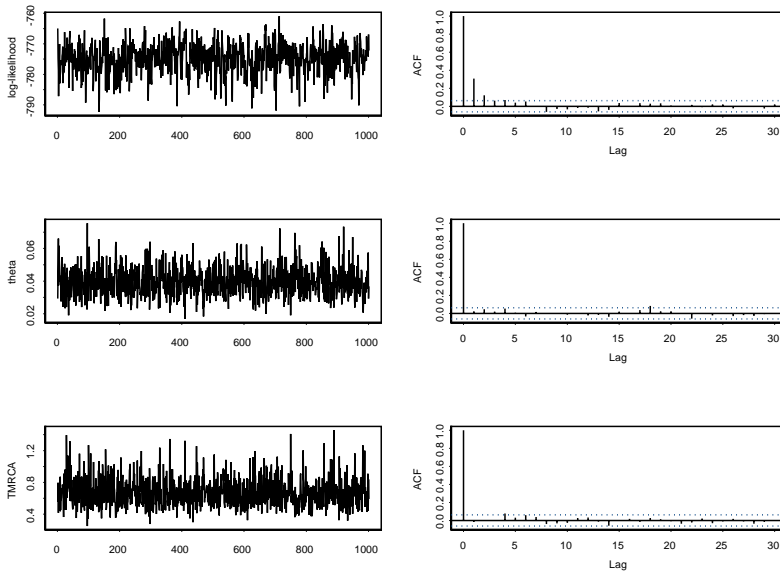
Model 2 above serves as the simplest description of mutation in hypervariable region I of mtDNA. It was used by Kuhner *et al.* (1995) in their analysis of the same data set.

We implemented the MCMC approach described in Algorithm 9.2. One should begin to sample from the process $X(\cdot)$ once it has "reached stationarity". There are many heuristic tests for this, none of which is infallible. For a critique see Gilks *et al.* (1996). Some simple diagnostics are functions of the statistics of interest such as autocorrelations and moving averages. It is also valuable to run the chain from several different, widely spaced, starting points, and compare the long-term behavior.

The output typically appeared to be non-stationary for up to 200,000 iterations of the algorithm. We sampled every 10,000th iteration in order to approximate a random sample from the stationary distribution. In a bid to be very conservative, and since the algorithms run rapidly, we generally discarded the first 2500 samples. After this, our output is typically based on 5000 samples. The acceptance rate was typically around 80%. For runs in which, for example, we needed to tune the variance parameter, the burn-in length varied but the estimated parameter values were unchanged for the different variances we tried.

Figure 9.11 shows the resultant time series for the log-likelihood, the mutation parameter θ , the time to the MRCA and their associated autocorrelation functions. These appear fine, with the proviso that the time series of log-likelihoods is correlated for several lags. While this is not in itself a problem it means one must interpret standard errors with care. As a further check for convergence to stationarity we used the package of diagnostics provided in CODA (Best *et al.* (1995)). All tests were passed.

Fig. 9.11. Example diagnostics



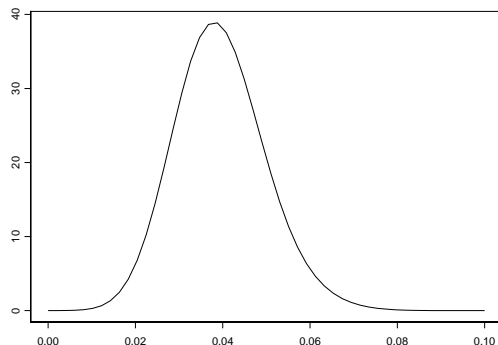
Some time can be saved by starting the process from a genealogy (A, T) for which $\mathbb{P}(A, T \mid \mathcal{D})$ is relatively high. The rationale for this is that it is sensible to start from a region of the state-space which is well supported by the data. As an example of this one might use the UPGMA tree for the data-set, as described in Kuhner *et al.* (1995). However, we prefer to start from random tree topologies since convergence from different starting points is potentially a useful diagnostic for stationarity.

The analysis of Model 1 gave a median for κ of 65.1, with 25th and 75th percentiles of 32.7 and 162.7 respectively. Note that the data are consistent with no transversions having occurred during the evolution of the sample. Consequently, the posterior distribution for κ has a very long right tail and statistics for the mean, which are strongly influenced by outliers, are poten-

tially misleading and are therefore not presented. The median value of g was 6.87×10^{-4} and the median value for w was 4.47×10^{-2} . These results show that the data are consistent with a value of $\kappa = 100$, as assumed by Kuhner *et al.* (1995).

In what follows we also took $\kappa = 100$, and a uniform prior on $(0, 100)$ for θ . The posterior distribution of the effective mutation rate has a median of 0.038, mean 0.039 and 25th and 75th percentiles of 0.033 and 0.045 respectively. Figure 9.12 shows the posterior distribution of θ .

Fig. 9.12. Posterior density of per site effective mutation rate θ



Since the posterior density of θ is proportional to the likelihood in this case, we may use an estimate of the posterior density to find the maximum likelihood estimate of θ . From the density shown in Figure 9.12, we obtained an MLE of $\hat{\theta} = 0.038$. Kuhner *et al.* (1995) obtained the value $\hat{\theta} = 0.040$ for these data, using the same value of κ . Presumably the difference in the estimates arises from both the parameters chosen for the density estimation, and the different approaches to the optimization. From an estimate of the curvature of the log-density we get an estimate of the standard error of $\hat{\theta}$ of 0.010, resulting in an approximate 95% confidence interval of (0.018, 0.058).

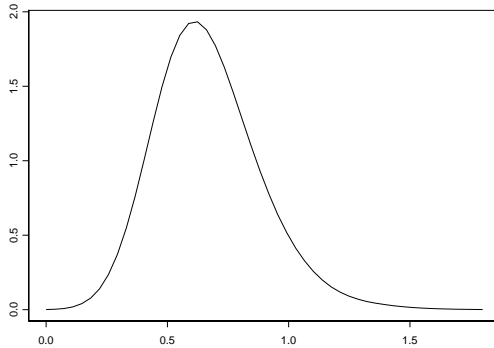
Remark. Estimates of standard errors based on curvature of the log-density should be treated as heuristic. In problems such as these, θ cannot be estimated consistently so the standard theory does not apply.

For comparison, the Watterson estimator (5.3.7) of θ , based on 26 segregating sites in the data, is 0.015 with an estimated standard error of 0.005; the 95% confidence interval for θ is then (0.005, 0.025). The lower MLE obtained

using the Watterson estimator is expected, because multiple mutations at the same site are ignored.

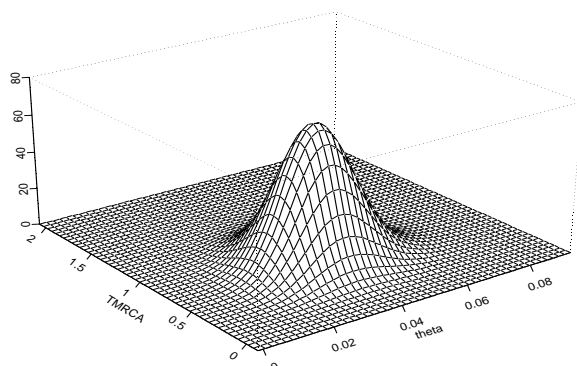
The prior distribution of the time to MRCA of a sample of $n = 63$ has a mean of $2(1 - 1/63) = 1.97$. With an effective size of $N = 600$, a 20 year generation time and a value of $\sigma^2 = 1$ for the variance of the offspring distribution, this is about 23,600 years. The posterior distribution of the time T_{MRCA} to the MRCA (in years) has median 7700, mean 8100 and 25th and 75th percentiles of 6500 and 9300 respectively. The corresponding posterior density appears in Figure 9.13. The joint posterior density of T_{MRCA} and θ is given in Figure 9.14. For a frequentist approach to inference about T_{MRCA} , see Tang *et al.* (2002).

Fig. 9.13. Posterior density of time to MRCA



Testing goodness-of-fit

The adequacy of the fit of models like these can be assessed using the Bayesian posterior predictive distribution. To implement this, we use a variant of the parametric bootstrap. The idea is to simulate observations from the *posterior* distribution of (Λ, \mathbf{T}, M) , and then for each of the trees (Λ, \mathbf{T}) to simulate the mutation process with parameters specified by M . The distribution of certain summary statistics observed in the simulated data is found, and the values of the statistics actually observed in the data are compared to these distributions. We chose to use the number of haplotypes, the maximal haplotype frequency, the haplotype homozygosity, the number of segregating sites and a measure of nucleotide diversity. In practice, we use the output from the MCMC runs to generate the observations on (Λ, \mathbf{T}, M) . In Table 11 we give the results of this comparison for Model 2 using 4000 values from each

Fig. 9.14. Joint posterior density of TMRCA and θ 

posterior distribution. There is some evidence that the constant rate model does not fit well, particularly regarding the haplotype distribution. The total number of segregating sites observed in the bootstrap samples gives some evidence of lack-of-fit; the model predicts more segregating sites than are seen in the data. One explanation for this apparent discrepancy might be that the model is not allowing for rate heterogeneity, and therefore does not typically produce enough recurrent mutations. This will lead to a tendency for the mutations which do occur to be spread over a greater number of sites. A model that allows for multiple classes of rates appears in Markovtsova *et al.* (2000b).

Remark. For an implementation of Bayesian methods for the coalescent (and many other species tree problems), using Metropolis-coupled MCMC, see Huelsenbeck and Ronquist's *MrBayes* program, at

<http://morphbank.ebc.uu.se/mrbayes/info.php>

9.9 The age of a UEP

In this section we provide an MCMC approach that can be used to find the posterior distribution of the age of a unique event polymorphism (UEP). As in the introduction of Section 8, there are several versions of this problem. For the most part, we assume that we have sequenced a region of DNA, and have determined for each of them whether or not the UEP is present. The key figure is given in Figure 8.1. Let \mathcal{U} denote the single event that causes the UEP mutation Δ . The scaled mutation rate at the UEP locus is $\mu/2$. The event that the coalescent tree has the UEP property is, once again, denoted by \mathcal{E} . For definiteness we assume that the sequences are evolving according to

Table 11. Assessing goodness-of-fit of Model 2

Statistic	Observed value	Model 2 Fraction of simulations \leq observed value
# haplotypes	28	0.83
max. haplotype frequency	9	0.36
homozygosity	0.0562	0.12
heterozygosity per site	0.0145	0.36
# segregating sites	26	0.05

Felsenstein’s model. The material in this section comes from Markovtsova *et al.* (2000a).

Modification of Markov chain Monte Carlo method

The event \mathcal{U} corresponds to a single mutation arising on the branch indicated in Figure 8.1 and no other mutations on the rest of the coalescent tree. Let A denote the age of the UEP, and denote the mutation parameters by $M = (g, \kappa, \mu)$. In what follows we assume a prior distribution for M , and apply an MCMC method for generating observations from the conditional density $f(A, G \mid \mathcal{D}, \mathcal{E} \cap \mathcal{U})$ of A and $G = (A, \mathbf{T}, M)$ given \mathcal{D}, \mathcal{E} and \mathcal{U} . To do this we express the required conditional density as a product of simpler terms and describe how each can be calculated.

First we note that

$$f(A, G \mid \mathcal{D}, \mathcal{U} \cap \mathcal{E}) = f(A \mid G, \mathcal{D}, \mathcal{U} \cap \mathcal{E})f(G \mid \mathcal{D}, \mathcal{U} \cap \mathcal{E}). \tag{9.9.1}$$

The first term on the right of (9.9.1) can be evaluated by considering Figure 8.1 once more. Given that a single mutation occurs on the indicated branch, the Poisson nature of the mutation process for the UEP means that the location of the mutation is uniformly distributed over that branch. Thus we can simulate observations from the conditional distribution of A by simulating from the second term on the right of (9.9.1), reading off the length of the branch on which the UEP mutation occurs, and adding a uniformly distributed fraction of that length to the height of the subtree containing all the chromosomes carrying the UEP. Our task is therefore reduced to simulating from the second term on the right of (9.9.1).

Let $g_1(A \mid \mathcal{E})$ denote the conditional distribution of the coalescent tree A given \mathcal{E} , $g_2(\mathbf{T})$ the density of the coalescence times \mathbf{T} , and $g_3(M)$ the prior for the mutation rates $M = (g, \kappa, \mu)$. We can then write

$$f(G | \mathcal{D}, \mathcal{U} \cap \mathcal{E}) = \mathbb{P}(\mathcal{D}, \mathcal{U} | G, \mathcal{E}) g_1(A | \mathcal{E}) g_2(\mathbf{T}) g_3(M) / \mathbb{P}(\mathcal{D}, \mathcal{U} | \mathcal{E}). \quad (9.9.2)$$

The term $\mathbb{P}(\mathcal{D}, \mathcal{U} | G, \mathcal{E})$ is the product of two terms,

$$\mathbb{P}(\mathcal{D}, \mathcal{U} | G, \mathcal{E}) = \mathbb{P}(\mathcal{D} | G, \mathcal{E}) \mathbb{P}(\mathcal{U} | G, \mathcal{E}).$$

The first of these, the likelihood of \mathcal{D} , can be computed using the peeling algorithm and the mutation model described above, while the second is

$$\frac{\mu S}{2} e^{-\mu S/2} \times e^{-\mu(L_n - S)/2} = \frac{\mu S}{2} e^{-\mu L_n/2}, \quad (9.9.3)$$

where S is the length of the branch on which the single UEP mutation must occur, and $L_n = \sum_{i=2}^n iT_i$ is the total length of the tree. The normalizing constant $\mathbb{P}(\mathcal{D}, \mathcal{U} \cap \mathcal{E})$ is unknown, and hard to compute. As a consequence, we use a version of the Metropolis-Hastings algorithm to simulate from the required conditional distribution.

Proposal kernel

We make a minor modification to Algorithm 9.2 in order to ensure that new trees are also consistent with the event \mathcal{E} . If, when we pick a level, we find we are in case A, and exactly two of the lines carry the UEP, then we cannot change the order in which the two coalescences occur, since such a change would produce a new tree topology which is inconsistent with \mathcal{E} . In such a situation we leave the topology unchanged.

Having constructed a new topology, which may be the same as the existing topology, we generate a new set of times in the same way as it was described in Section 9.5. We found that a kernel which proposes new values of T'_l and T'_{l-1} having the pre-data coalescent distribution worked well.

Finally, we update $M = (g, \kappa, \mu)$, where g and κ are the rate parameters for the sequence model and μ is the rate parameter for the UEP. The parameters g and κ were updated every tenth iteration, and μ was updated on each iteration for which g was not updated. These were updated using truncated Normals, whose variances require some tuning.

The Hastings ratio

Writing $G = (A, \mathbf{T}, M)$, the kernel Q can be expressed as the product of three terms:

$$Q(G' \rightarrow G) = Q_1(A' \rightarrow A) Q_2(\mathbf{T}' \rightarrow \mathbf{T} | A' \rightarrow A) Q_3(M' \rightarrow M).$$

Using (9.9.1), (9.9.2) and (9.9.3), the Hastings ratio (the probability with which we accept the new state) can be written in the form

$$h = \min \left\{ 1, \frac{\mathbb{P}(\mathcal{D} \mid G', \mathcal{E})}{\mathbb{P}(\mathcal{D} \mid G, \mathcal{E})} \frac{\mathbb{P}(U \mid G', \mathcal{E})}{\mathbb{P}(U \mid G, \mathcal{E})} \frac{g_1(A' \mid \mathcal{E})}{g_1(A \mid \mathcal{E})} \frac{g_2(\mathbf{T}')}{g_2(\mathbf{T})} \frac{g_3(M')}{g_3(M)} \right. \\ \left. \times \frac{Q_1(A' \rightarrow A)}{Q_1(A \rightarrow A')} \frac{Q_2(\mathbf{T}' \rightarrow \mathbf{T} \mid A' \rightarrow A)}{Q_2(\mathbf{T} \rightarrow \mathbf{T}' \mid A \rightarrow A')} \frac{Q_3(M' \rightarrow M)}{Q_3(M \rightarrow M')} \right\},$$

the unknown term $\mathbb{P}(\mathcal{D}, \mathcal{U} \cap \mathcal{E})$ cancelling. For our choice of transition kernel Q , it can be shown that $g_1(A' \mid \mathcal{E}) = g_1(A \mid \mathcal{E})$. We also have $Q_1(A \rightarrow A') = Q_1(A' \rightarrow A)$, and we note that Q changes only two of the times associated with T or T' . Hence h reduces to

$$h = \min \left\{ 1, \frac{\mathbb{P}(\mathcal{D} \mid G', \mathcal{E})}{\mathbb{P}(\mathcal{D} \mid G, \mathcal{E})} \frac{\mathbb{P}(U \mid G', \mathcal{E})}{\mathbb{P}(U \mid G, \mathcal{E})} \frac{g_2(\mathbf{T}')g_3(M')}{g_2(\mathbf{T})g_3(M)} \right. \\ \left. \times \frac{f_l(t_l)f_{l-1}(t_{l-1})}{f_l(t'_l)f_{l-1}(t'_{l-1})} \frac{Q_3(M' \rightarrow M)}{Q_3(M \rightarrow M')} \right\}, \tag{9.9.4}$$

where $f_l(\cdot)$ and $f_{l-1}(\cdot)$ are the densities of the time updating mechanism given that changes occur to the tree A at levels l and $l - 1$.

In Section 8 we derived a number of theoretical results concerning the age of a UEP given its frequency in the sample in the limiting case $\mu \rightarrow 0$. In order to compare these results with those obtained by including the sequence information, we modified our algorithm to allow $\mu = 0$. Assuming κ is known, the mutation parameter M is now one-dimensional: $M = (g)$. The other change occurs to the conditional probability in (9.9.3), since now $\mathbb{P}(U \mid G, \mathcal{E}) \propto S$, the length of the branch on which the UEP mutation must occur. This change appears in the Hastings ratio (9.9.4), where

$$\frac{\mathbb{P}(U \mid G', \mathcal{E})}{\mathbb{P}(U \mid G, \mathcal{E})} = \frac{S'}{S}.$$

In order to check tree moves, we can again use the infinitely-many-sites model of mutation. We compare distributions of time to the most recent common ancestor of the group of individuals carrying a specific mutation, the length of the corresponding sub-tree and the time to the mutation generated by the rejection method described in Algorithm 8.2 for the $\mu = 0$ case, and the modified version of our general MCMC scheme.

9.10 A Yakima data set

To illustrate the method we find the conditional distribution of the age of the 9 basepair mitochondrial region V deletion in a sample of Yakima described by Shields *et al.* (1993) The sample comprise $n = 42$ individuals, of whom $b = 26$ have the deletion. The data \mathcal{D} comprise 360 basepairs from hyper-variable region I of the control region, sequenced for all 42 individuals. The observed base frequencies are $(\pi_A, \pi_G, \pi_C, \pi_T) = (0.328, 0.113, 0.342, 0.217)$. We note that all individuals having a given control region sequence had the

same deletion status, as might be expected if the deletion arose once quite recently.

For the analysis discussed here, the output typically appeared to be non-stationary for at least 200,000 iterations of the algorithm. We generally discarded the first 25 million iterations. After this, we sampled every 5,000th iteration. Our output is typically based on 5000 samples from our stationary process. The acceptance rate was generally around 70%.

Preliminary analysis of the sequence data (without regard to presence or absence of the deletion) was performed using the approach outlined in Section 9.5. For the present mutation model, we took uninformative priors (in the form of uniform densities having wide but finite support) for the mutation rates g and w and examined the posterior distribution of $\kappa = w/g$. The posterior median was 65.9, the distribution having 25th percentile of 34.0 and 75th percentile of 160.2. The data are certainly consistent with the value of $\kappa = 100$ we used in the Nuh Chah Nulth example in Section 9.8. We therefore treat $\kappa = 100$ as fixed in the subsequent analyses; from (9.1.7) we find that $\theta = 88.17g$.

We repeated the analysis with an uninformative prior, uniform on $(0, 0.1)$, for the single parameter g . This resulted in the posterior density for θ given in Figure 9.15. Summary statistics are shown in Table 12. Our approach also provides a way to find the maximum likelihood estimator of θ , since with a flat prior the posterior is proportional to the likelihood. From a kernel density estimate we obtained an MLE of $\hat{\theta} = 0.039$ with an estimated standard error of 0.010. This is consistent with the estimate of θ we found for the Nuu Chah Nulth data. Since the base frequencies in both data sets are similar and the mutation rates are likely to be the same, we conclude that the effective sizes of the two populations are also approximately equal. The effective population size of the Nuu Chah Nulth was estimated from anthropological data by Ward *et al.* (1991) to be about $N = 600$, a number we take for the Yakima as well.

Under the pre-data coalescent distribution, the mean time to the MRCA of a sample of $n = 42$ is $2(1 - 1/42) = 1.95$. With an effective size of $N = 600$ and a 20 year generation time, this is about 23,500 years. The posterior density of the time to the MRCA given the control region data D is shown in Figure 9.16. The posterior mean is 0.72, or about 8,600 years. Summary statistics are given in Table 13. The posterior distribution of the total tree length $L_{42} = \sum_{j=2}^{42} jT_j$ has mean 5.68.

We turn now to the deletion data. We ran our MCMC algorithm using a uniform $(0, 10)$ prior for μ , and a uniform $(0, 0.1)$ prior for g . The posterior density of θ is shown in Figure 9.15. Summary statistics are presented in Table 12. The distribution is qualitatively the same as that obtained by ignoring the deletion data. The posterior distribution of the deletion parameter μ has mean 0.75 and median 0.61; the 25th percentile is 0.34 and the 75th percentile is 0.99. The posterior density of the time to the MRCA of the group carrying the deletion is shown in Figure 9.17. The summary statistics are found in Table 14.

Fig. 9.15. Posterior density of mutation rate θ

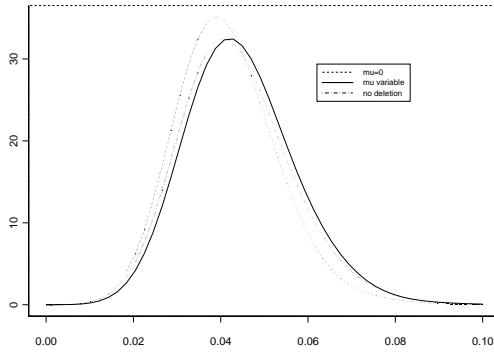


Fig. 9.16. Posterior density of TMRCA

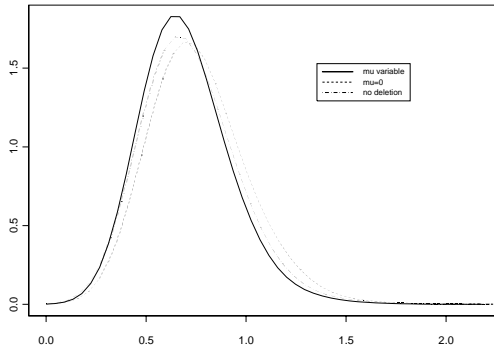
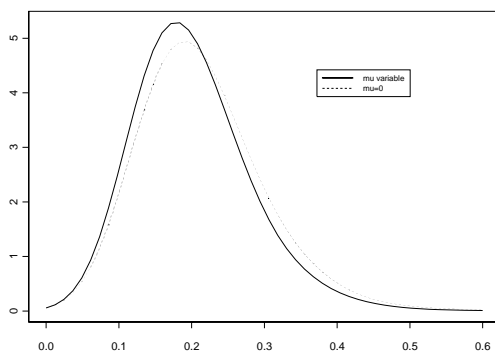


Table 12. Summary statistics for θ

θ	no deletion	μ variable	$\mu = 0$
mean	0.044	0.045	0.041
median	0.042	0.043	0.040
25th percentile	0.036	0.037	0.034
75th percentile	0.050	0.051	0.047

Table 13. Summary statistics for time to MRCA of the sample.

Time to MRCA	no deletion	μ variable	$\mu = 0$
mean	0.72 (8,600 yrs)	0.70 (8,400 yrs)	0.76 (9,200 yrs)
median	0.69 (8,300 yrs)	0.67 (8,000 yrs)	0.73 (8,800 yrs)
25th percentile	0.57 (6,800 yrs)	0.56 (6,700 yrs)	0.61 (7,300 yrs)
75th percentile	0.84 (10,100 yrs)	0.81 (9,700 yrs)	0.88 (10,600 yrs)

Fig. 9.17. Posterior density of TMRCA of deletion

The deletion arises uniformly on the branch indicated in Figure 8.1, so that the age of the mutation is the time to the MRCA of the deletion group plus a uniform fraction of the mutation branch length. The posterior distribution of the age is given in Figure 9.18, and summary statistics in Table 15.

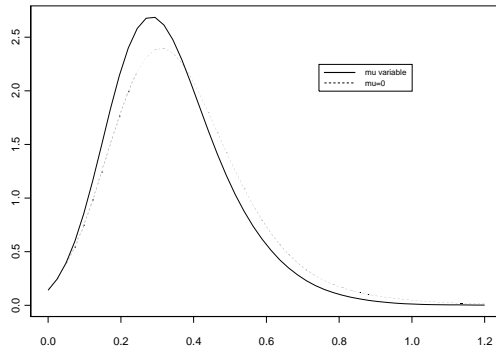
We also looked at the time to the MRCA of the entire sample when the deletion status of each sequence is included. The posterior density of this time is shown in Figure 9.16, with summary statistics given in Table 13. For these data the inclusion of deletion status has little effect on the posterior distribution.

The output from the MCMC runs can be used to assess whether the UEP assumption is reasonable. We first generated 5000 observations of the tree length L_{42} conditional on the data \mathcal{D} ; as noted above, the sample mean is

Table 14. Summary statistics for the time to MRCA of the group carrying the deletion.

Time to MRCA	μ variable	$\mu = 0$
mean	0.20 (2400 yrs)	0.21 (2600 yrs)
median	0.19 (2300 yrs)	0.20 (2400 yrs)
25th percentile	0.15 (1800 yrs)	0.16 (1900 yrs)
75th percentile	0.24 (2900 yrs)	0.25 (3100 yrs)

Fig. 9.18. Posterior density of age of deletion



5.68. The modal posterior value of μ is 0.30, a value that we treat as a point estimate of μ . The expected number of deletions arising on the coalescent tree is then $0.30 \mathbb{E}(L_{42}|\mathcal{D})/2$, which we estimate from the posterior mean tree length as $0.30 \times 5.68/2 = 0.85$. We can also use this value of μ and the simulated values of L_{42} to estimate the probability that exactly one mutation would occur on such a tree; we obtained an estimate of 0.36. Similarly, we estimated the probability of at least one mutation occurring as 0.57, so that the conditional probability that the mutation occurred once, given it occurred at least once, is estimated to be 0.63. Thus it is not unreasonable to assume that the deletion arose just once.

When $\mu = 0$, the posterior density of θ is shown in Figure 9.15, with summary statistics given in Table 12; there is little difference from the case where

Table 15. Summary statistics for age of the deletion.

Age of deletion	μ variable	$\mu = 0$
mean	0.34 (4100 yrs)	0.36 (4400 yrs)
median	0.31 (3700 yrs)	0.33 (4000 yrs)
25th percentile	0.23 (2800 yrs)	0.25 (3000 yrs)
75th percentile	0.41 (5000 yrs)	0.44 (5300 yrs)

μ is allowed to vary. The posterior density of the time to the MRCA is given in Figure 9.16, with summary statistics in Table 13. The mean time of 0.76 (or about 9,100 years) stands in marked contrast to the value of 2.68 (about 32,200 years) obtained from Griffiths and Marjoram (1996). The summary statistics for the posterior distribution of the time to the MRCA of the group carrying the deletion are given in Table 14. The results are qualitatively the same as the case of variable μ . The posterior density of the age of the deletion appears in Figure 9.18, with summary statistics shown in Table 15. The posterior mean is 0.36 (or about 4,400 years), compared to the value of 1.54 (or about 18,500 years) obtained from equation (8.3.4) when the sequence data are ignored. As expected, the mean age is higher than it is when μ is non-zero.

10 Recombination

In this section we study the generalization of the coalescent to the case of recombination. The basic groundwork of the subject comes from the seminal paper of Hudson (1983) and the ancestral recombination graph described by Griffiths (1991). We study the two locus model first, and then generalize to a model with arbitrary recombination rates. Later in the section we discuss methods for estimating the recombination rate, the behavior of measures of linkage disequilibrium, and uses of the coalescent for fine-scale mapping of disease genes.

10.1 The two locus model

Consider two linked loci, A and B , in a population of fixed size N chromosomes; neutrality, random mating and constant population size are assumed as before. For convenience, suppose the population reproduces according to a Wright-Fisher model with recombination: independently across offspring, in the next generation

- (i) with probability $1 - r$ the individual chooses a chromosome from the previous generation and inherits the genes at the A and B loci.
- (ii) with probability r the individual chooses 2 chromosomes from the previous generation and inherits the gene at the A locus from one and the gene at the B locus from the other.

In this model recombination is possible only between the two loci. If we focus on either of the two loci alone, we are watching a Wright-Fisher process evolve. It follows that the genealogical tree of a sample from one of the loci is described by the coalescent. There is thus a genealogical tree for each of the two loci. The effect of recombination is to make these two trees correlated. If $r = 0$, the loci are completely linked and the trees at each locus are identical. Early results for this model were obtained by Strobeck and Morgan (1978) and Griffiths (1981).

We consider the case in which N is large and r is of order N^{-1} ; this balances the effects of drift and recombination. We define the (scaled) recombination rate ρ by

$$\rho = \lim_{N \rightarrow \infty} 2Nr \tag{10.1.1}$$

The ancestral process

Just as in the earlier model, we can calculate the chance that if there are currently k ancestors of the sample then in the previous generation there are also k . To the order of approximation we need, this occurs only if there are no recombination events in the k ancestors as they choose their parents, and the k also chose distinct parents. This event has probability

$$(1-r)^k \left(1 - \frac{1}{N}\right) \cdots \left(1 - \frac{k-1}{N}\right),$$

which, in the light of (10.1.1) is just

$$1 - \frac{k\rho}{2N} - \frac{k(k-1)}{2N} + O(N^{-2}).$$

In a similar way, we can compute the probability that the number of distinct parents chosen in the previous generation increases from k to $k+1$. To the order we need, this occurs if precisely one recombination event occurs and the other $k-1$ ancestors choose distinct parents. A straightforward calculation shows that this probability is

$$\frac{k\rho}{2N} + O(N^{-2}).$$

Finally we can compute the chance that the number of ancestors goes down by 1, from k to $k-1$. The same sort of calculation shows this is

$$\frac{k(k-1)}{2N} + O(N^{-2}).$$

All other possibilities have smaller order. Thus we conclude that the number $A_n^N(Nt)$ behaves in the limit as $N \rightarrow \infty$ like continuous time birth and death process in which the transition rates are

$$\begin{aligned} k &\rightarrow k+1 && \text{at rate } k\rho/2 \\ k &\rightarrow k-1 && \text{at rate } k(k-1)/2 \end{aligned}$$

starting from state n . Because of the quadratic death rate compared to the linear growth rate, it is clear that the process will visit the value 1 infinitely often. The first occurrence of 1 corresponds to an MRCA.

A number of properties of the ancestral process $A_n^\rho(\cdot)$ can be found simply. Let M_n denote the maximum number of ancestors of the sample before it reaches its MRCA, and let τ_n denote the time to this MRCA. Griffiths (1991) proved:

Lemma 10.1 *The expected TMRCA is given by*

$$\mathbb{E}\tau_n = \frac{2}{\rho} \int_0^1 \left(\frac{1-v^{n-1}}{1-v} \right) (e^{\rho(1-v)} - 1) dv, \quad (10.1.2)$$

and the distribution of M_n is given by

$$\mathbb{P}(M_n \leq k) = \frac{\sum_{j=n-1}^{k-1} j! \rho^{-j}}{\sum_{j=0}^{k-1} j! \rho^{-j}}, \quad k \geq n. \quad (10.1.3)$$

Proof. The expected height follows from standard results for birth-and-death processes. Define

$$\rho_i = \frac{\mu_2 \cdots \mu_{i-1}}{\lambda_2 \cdots \lambda_i}, \quad i \geq 2.$$

For the ancestral process, it can be checked that $\rho_i = 2\rho^{i-2}/i!$. The waiting time to reach 1 has mean given by

$$\mathbb{E}\tau_n = \sum_{r=1}^{n-1} \left(\prod_{k=2}^r \frac{\mu_k}{\lambda_k} \right) \sum_{j=r+1}^{\infty} \rho_j,$$

where empty products have value 1 by convention. In our setting, this reduces to

$$\begin{aligned} \mathbb{E}\tau_n &= 2 \sum_{m=2}^n \sum_{l \geq 0} \rho^l \frac{(m-2)! \Gamma(m+2)}{(l+m)!(m+1)!} \\ &= \frac{2}{\rho} \int_0^1 \left(\frac{1-v^{n-1}}{1-v} \right) (e^{\rho(1-v)} - 1) dv. \end{aligned}$$

To find the distribution of M_n , define $p_n(k) = \mathbb{P}(M_n \leq k)$, with $p_1(k) = 1, k \geq 1$ and $p_n(k) = 0$ if $n > k$. By considering whether a coalescence or a recombination occurs first in the ancestry of the sample, we see that

$$p_n(k) = \frac{n-1}{\rho+n-1} p_{n-1}(k) + \frac{\rho}{\rho+n-1} p_{n+1}(k),$$

and it may readily be checked by induction that the solution is given by (10.1.3). \square

As $\rho \downarrow 0$, we see from (10.1.2) that $\mathbb{E}\tau_n \rightarrow 2 \int_0^1 (1-v^{n-1}) dv = 2(1-1/n)$, as expected from our study of the coalescent. As $\rho \rightarrow \infty$, $\mathbb{E}\tau_n \rightarrow \infty$ also. When $n = 2$, we have

$$\mathbb{E}\tau_2 = 2\rho^{-2}(e^\rho - 1 - \rho),$$

and as $n \rightarrow \infty$,

$$\mathbb{E}\tau_\infty = \frac{2}{\rho} \int_0^1 v^{-1}(e^{\rho v} - 1) dv.$$

This last can be interpreted as the time taken for the whole population to be traced back to its common ancestor.

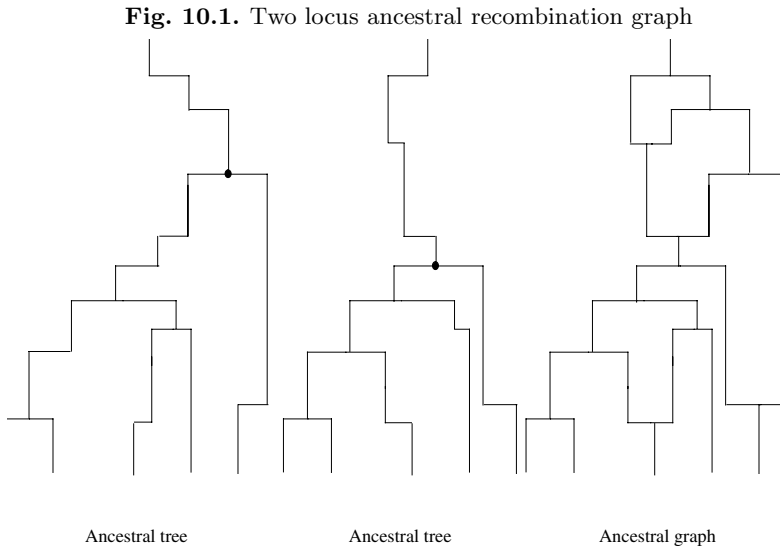
It follows from (10.1.3) that $M_n/n \rightarrow 1$ in probability as $n \rightarrow \infty$, showing that the width of the graph does not exceed n by very much.

The ancestral recombination graph

We have seen that the ancestral process starts from $A_n^\rho(0) = n$, and has the property that if there are currently k ancestors then

- (i) Any particular pair of branches coalesce at rate 1.
- (ii) Any given branch splits into two at rate $\rho/2$.

The ancestral process $A_n^\rho(\cdot)$ is of limited use on its own; just as in the coalescent setting it is the way these individuals are related that matters. This leads to the idea of the *ancestral recombination graph* (or ARG). We construct such an ancestral recombination graph in such a way that when two edges are added at a recombination event, the genes represented by the left branch correspond to the A locus, and the right edges correspond to the B locus. In this way the ancestry of the A locus may be traced by following the left branch at each split, and the ancestry of the B locus by following the right branch. The ancestry of the A locus is a coalescent tree \mathcal{T}_A , and the ancestry of the B locus is a coalescent tree \mathcal{T}_B . These trees are dependent. Each tree has its own MRCA (which might be the same). An example of the ancestral graph, together with the two subtrees \mathcal{T}_A and \mathcal{T}_B is given in Figure 10.1. The MRCA at each locus marginally is denoted by a \bullet .



Note that τ_n may now be interpreted as the height of the ARG, and M_n may be interpreted as its width. Of course, τ_n is at least as great as the time taken to find the MRCA at the A locus and at the B locus.

The structure of the ARG

In this section, we study the structure of the genealogical graph \mathcal{G} in more detail. The graph includes the coalescent tree \mathcal{T}_A of the A locus and the

coalescent tree \mathcal{T}_B of the B locus. Denote the edge set of a graph by $\mathcal{E}(\cdot)$. It is useful to partition the edges $\mathcal{E}(\mathcal{G})$ into four disjoint sets:

$$\begin{aligned} \mathcal{A} &= \mathcal{E}(\mathcal{T}_A) \cap \mathcal{E}(\mathcal{T}_B)^c; \\ \mathcal{B} &= \mathcal{E}(\mathcal{T}_A)^c \cap \mathcal{E}(\mathcal{T}_B); \\ \mathcal{C} &= \mathcal{E}(\mathcal{T}_A) \cap \mathcal{E}(\mathcal{T}_B); \\ \mathcal{D} &= \mathcal{E}(\mathcal{G}) \cap \mathcal{E}(\mathcal{T}_A)^c \cap \mathcal{E}(\mathcal{T}_B)^c. \end{aligned}$$

Those edges in \mathcal{A} represent ancestors who contribute to the genetic material of the sample at the A locus *only*, and similarly for \mathcal{B} and the B locus. Edges in \mathcal{C} correspond to ancestors that contribute genetic material at both loci, and those in \mathcal{D} contribute no genetic material to the sample.

At any given time t , the ancestors of the sample (i.e. the edges $\mathcal{E}(\mathcal{G}_t)$ of the ancestral graph \mathcal{G}_t of a cross section of \mathcal{G} taken at time t) can be divided into these four types. Define

$$\begin{aligned} n_{\mathcal{A}}(t) &= |\mathcal{E}(\mathcal{G}_t) \cap \mathcal{A}| \\ n_{\mathcal{B}}(t) &= |\mathcal{E}(\mathcal{G}_t) \cap \mathcal{B}| \\ n_{\mathcal{C}}(t) &= |\mathcal{E}(\mathcal{G}_t) \cap \mathcal{C}| \\ n_{\mathcal{D}}(t) &= |\mathcal{E}(\mathcal{G}_t) \cap \mathcal{D}|, \end{aligned}$$

where $|\cdot|$ denotes the number of elements in a set. Clearly

$$n_{\mathcal{A}}(t) + n_{\mathcal{B}}(t) + n_{\mathcal{C}}(t) + n_{\mathcal{D}}(t) = |\mathcal{E}(\mathcal{G}_t)| \equiv A_n^{\rho}(t),$$

where $A_n^{\rho}(t)$ is the ancestral process of the ARG. Furthermore,

$$n_{\mathcal{A}}(t) + n_{\mathcal{C}}(t) = |\mathcal{E}(\mathcal{T}_A(t))| \equiv A_n(t), \tag{10.1.4}$$

and

$$n_{\mathcal{B}}(t) + n_{\mathcal{C}}(t) = |\mathcal{E}(\mathcal{T}_B(t))| \equiv B_n(t), \tag{10.1.5}$$

where $A_n(\cdot)$ and $B_n(\cdot)$ are the marginal ancestral processes for the A and B loci respectively.

Of interest is the evolution of the process

$$\mathbf{m}(t) = (n_{\mathcal{A}}(t), n_{\mathcal{B}}(t), n_{\mathcal{C}}(t), n_{\mathcal{D}}(t)), \quad t \geq 0.$$

One way to think of the process \mathbf{m} is to label edges as $(1,0)$, $(0,1)$, $(1,1)$, or $(0,0)$ according as the edge is in \mathcal{A} , \mathcal{B} , \mathcal{C} , or \mathcal{D} respectively. When a coalescence occurs to two edges of type (α, β) and (γ, δ) the resultant ancestor is of type $(\max(\alpha, \gamma), \max(\beta, \delta))$, and if a recombination occurs to an edge of type (α, β) , the two new edges are of type $(\alpha, 0)$ and $(0, \beta)$.

Ethier and Griffiths (1990a) show that the process is Markovian. If the current state is (a, b, c, d) , the next state and its transition rate are given by

$(a + 1, b + 1, c - 1, d)$		$c\rho/2$
$(a - 1, b - 1, c + 1, d)$		ab
$(a - 1, b, c, d)$		$ac + a(a - 1)/2$
$(a, b - 1, c, d)$	at rate	$bc + b(b - 1)/2$
$(a, b, c - 1, d)$		$c(c - 1)/2$
$(a, b, c, d + 1)$		$(a + b + d)\rho/2$
$(a, b, c, d - 1)$		$d(a + b + c) + d(d - 1)/2$.

To see this, consider first the transition $(a, b, c, d) \rightarrow (a + 1, b + 1, c - 1, d)$: this occurs if a recombination event occurs on an edge of type (1,1). This results in loss of a (1,1) edge, and the addition of one (1,0) edge and one (0,1) edge. The rate of such changes is $c\rho/2$. Considering the change $(a, b, c, d) \rightarrow (a - 1, b, c, d)$ for example, we see that this results from a coalescence of a (1,0) edge and a (1,1) edge, or the coalescence of two (1,0) edges. Both possibilities result in the net loss of a (1,0) edge. The first type of change occurs at rate ac and the second sort at rate $a(a - 1)/2$. In a similar way the other transitions and their rates may be verified. The overall transition rate is the sum of these rates; if $a + b + c + d = n$, this rate is given by $d_n \equiv c\rho/2 + n(n - 1)/2$.

There is a reduced version of the Markov chain $\mathbf{m}(\cdot)$ that records only the first three coordinates:

$$\mathbf{n}(t) = (n_A(t), n_B(t), n_C(t)), \quad t \geq 0.$$

Examining the transition rates of $\mathbf{m}(\cdot)$ given above shows that $\mathbf{n}(\cdot)$ is also Markovian, and from a state of the form (a, b, c) its transitions are to

$$(a_1, b_1, c_1) = \begin{cases} (a + 1, b + 1, c - 1) & r_1 = c\rho/2 \\ (a - 1, b - 1, c + 1) & r_2 = ab \\ (a - 1, b, c) & \text{at rate } r_3 = ac + a(a - 1)/2 \\ (a, b - 1, c) & r_4 = bc + b(b - 1)/2 \\ (a, b, c - 1) & r_5 = c(c - 1)/2 \end{cases} \quad (10.1.6)$$

Note that recombination takes place only on the edges in \mathcal{C} . The rate of change from a state (a, b, c) with $n = a + b + c$ is given by

$$d_n \equiv \frac{c\rho}{2} + \frac{n(n - 1)}{2}. \quad (10.1.7)$$

Since the values of both $n_A(t) + n_C(t)$ and $n_B(t) + n_C(t)$ cannot increase as t increases, and eventually both must have the value 1, we see that the reduced process has absorbing states at $(1, 0, 0)$, $(0, 1, 0)$ and $\{(1, 1, 0), (0, 0, 1)\}$. It starts from $\mathbf{n}(0) = (0, 0, n)$. We might also consider the case in which only some of the genes, say $c < n$, are typed at both loci, while a are typed only at the A locus and the remaining $b = n - a - c$ are typed only at the B locus. In this case, $\mathbf{n}(0) = (a, b, c)$.

10.2 The correlation between tree lengths

In this section, we derive a recursion satisfied by the covariance of the tree lengths L^A and L^B of the marginal trees \mathcal{T}_A and \mathcal{T}_B respectively. The development here follows that of Pluzhnikov (1997).

For an initial configuration $\mathbf{n}(0) = (a, b, c)$ define $F(a, b, c; \rho)$ to be the covariance between L^A and L^B . Thus $F(0, 0, n; \rho)$ is the covariance of the marginal tree lengths for a sample of size n typed at both loci. We watch the Markov chain $\mathbf{n}(\cdot)$ only at the points it changes state. The resulting jump chain is denoted by $\mathbf{N}(\cdot)$. Let Z be a random variable that gives the outcome of a one-step jump of the chain $\mathbf{N}(\cdot)$ starting from (a, b, c) , and let $Z = z_1$ correspond to the move to $(a + 1, b + 1, c - 1)$, $Z = z_2$ correspond to the move to $(a - 1, b - 1, c + 1)$ and so on, in the order given in (10.1.6). The jump probabilities are

$$p_i = \mathbb{P}(Z = z_i) = r_i/d_n, \quad i = 1, \dots, 5. \tag{10.2.1}$$

Pluzhnikov (1997) established the following representation, which follows immediately from the properties of the coalescent trees \mathcal{T}_A and \mathcal{T}_B and the ARG.

Lemma 10.2 *Conditional on $Z = (a_1, b_1, c_1)$, we have*

$$L^A = X_A + T_A \tag{10.2.2}$$

where

- (i) $X_A \sim L_{a_1+c_1}$, where L_m denotes the length of an m -coalescent tree;
- (ii) $T_A \sim n_1 T$, where T is exponential(d_n) and $n_1 = a + c$;
- (iii) X_A and T_A are independent.

Furthermore, a similar representation holds for L^B given Z :

$$L^B = X_B + T_B \sim L_{b_1+c_1} + n_2 T, \tag{10.2.3}$$

where $n_2 = b + c$. In addition, X_B and T_A are independent, as are X_A and T_B .

This leads to the main result of this section, derived originally in somewhat different form by Kaplan and Hudson (1985).

Theorem 10.3 *For any $\rho \in [0, \infty)$, the covariance $\text{Cov}(L^A, L^B) := F(a, b, s; \rho)$ satisfies the linear system*

$$\begin{aligned} d_n F(a, b, c; \rho) &= r_1 F(a + 1, b + 1, c - 1; \rho) + r_2 F(a - 1, b - 1, c + 1; \rho) \\ &\quad + r_3 F(a - 1, b, c; \rho) + r_4 F(a, b - 1, c; \rho) \\ &\quad + r_5 F(a, b, c - 1; \rho) + R_n \end{aligned} \tag{10.2.4}$$

where $n = a + b + c$, $n_1 = a + c$, $n_2 = b + c$, $d_n = (n(n - 1) + c\rho)/2$, the r_i are given in (10.1.6), and $R_n = 2c(c - 1)/((n_1 - 1)(n_2 - 1))$. The system (10.2.4) has a unique solution satisfying the boundary conditions

$$F(a, b, c; \rho) = 0 \text{ whenever } n_1 < 2, \text{ or } n_2 < 2, \text{ or } a < 0, \text{ or } b < 0, \text{ or } c < 0. \tag{10.2.5}$$

Proof. The proof uses the formula for conditional covariances, namely

$$\text{Cov}(X, Y) = \mathbb{E}(\text{Cov}(X, Y \mid Z)) + \text{Cov}(\mathbb{E}(X \mid Z), \mathbb{E}(Y \mid Z)),$$

with $X = L^A$, $Y = L^B$ and Z as defined above. Clearly,

$$\mathbb{E}(\text{Cov}(X, Y \mid Z)) = \sum_{i=1}^5 p_i \text{Cov}(X, Y \mid Z = z_i),$$

where the p_i are defined in (10.2.1). Now

$$\begin{aligned} \text{Cov}(X, Y \mid Z = z_1) &= \text{Cov}(X_A + T_A, X_B + T_B) \\ &= \text{Cov}(X_A, X_B) + \text{Cov}(T_A, T_B) \\ &= F(a + 1, b + 1, c - 1, ; \rho) + n_1 n_2 \text{Var}(T) \\ &= F(a + 1, b + 1, c - 1, ; \rho) + n_1 n_2 d_n^{-2} \end{aligned} \tag{10.2.6}$$

Using similar arguments gives

$$\begin{aligned} \mathbb{E}(\text{Cov}(X, Y \mid Z)) &= r_1 F(a + 1, b + 1, c - 1; \rho) + r_2 F(a - 1, b - 1, c + 1; \rho) \\ &\quad + r_3 F(a - 1, b, c; \rho) + r_4 F(a, b - 1, c; \rho) \\ &\quad + r_5 F(a, b, c - 1; \rho) + n_1 n_2 d_n^{-2}. \end{aligned} \tag{10.2.7}$$

Next, recall that

$$\text{Cov}(\mathbb{E}(Y \mid Z), \mathbb{E}(Y \mid Z)) = \mathbb{E}[(\mathbb{E}(X \mid Z) - \mathbb{E}(X))(\mathbb{E}(Y \mid Z) - \mathbb{E}(Y))].$$

Using basic properties of the regular coalescent, we can derive the distributions of $f(Z) = \mathbb{E}(X \mid Z) - \mathbb{E}(X)$ and $g(Z) = \mathbb{E}(Y \mid Z) - \mathbb{E}(Y)$; these are given in Table 16. Hence we find that

$$\begin{aligned} \text{Cov}(\mathbb{E}(X \mid Z), \mathbb{E}(Y \mid Z)) &= \mathbb{E}(f(Z)g(Z)) \\ &= \sum_{i=1}^5 p_i f(z_i)g(z_i) \\ &= -\frac{n_1 n_2}{d_n^2} + \frac{2c(c - 1)}{d_n(n_1 - 1)(n_2 - 1)} \end{aligned} \tag{10.2.8}$$

Adding (10.2.7) and (10.2.8) yields (10.2.4).

Table 16. The probability distribution of $f(Z)$ and $g(Z)$

Z	$f(Z)$	$g(Z)$	$\mathbb{P}(Z = z_i)$
$(a + 1, b + 1, c - 1)$	n_1/d_n	n_2/d_n	p_1
$(a - 1, b - 1, c + 1)$	n_1/d_n	n_2/d_n	p_2
$(a - 1, b, c)$	$n_1/d_n - 2/(n_1 - 1)$	n_2/d_n	p_3
$(a, b - 1, c)$	n_1/d_n	$n_2/d_n - 2/(n_2 - 1)$	p_4
$(a, b, c - 1)$	$n_1/d_n - 2/(n_1 - 1)$	$n_2/d_n - 2/(n_2 - 1)$	p_5

The boundary conditions follow from the restriction that the ancestral process for each locus be considered no further back than its MRCA. \square

Equations like (10.2.4) can be solved by observing that if the degree of $F(a, b, c)$ is defined as $a + b + 2c$, then the degree on the right is at most the degree on the left; knowing lower degree terms allows the higher degree terms to be found by solving a lower triangular system of equations. Ethier and Griffiths (1990) developed an efficient computational method for solving such systems. The solution is known explicitly in very few cases, among them Griffiths’ (1981) result

$$F(0, 0, 2; \rho) = \frac{4(\rho + 18)}{\rho^2 + 13\rho + 18}. \tag{10.2.9}$$

Some other examples

The equation in (10.2.4) can be written in the form

$$F(a, b, c; \rho) = \mathcal{L}F + g(a, b, c; \rho) \tag{10.2.10}$$

where in (10.2.4) we had $g(a, b, c; \rho) = d_n^{-1}R_n$. The same type of equation arises in studying many properties of the ARG. We mention two of them, derived by Griffiths (1991).

Define the time $W_n = \max(T_A, T_B)$ by which the sample of size n has a common ancestor at both the A and B loci. This is the time taken to reach the states $\{(1, 1, 0), (0, 0, 1)\}$ starting from $(0, 0, n)$. Starting from a configuration of (a, b, c) with $a + b + c = n$, the expected waiting time $f(a, b, c; \rho)$ satisfies (10.2.10) with

$$g(a, b, c; \rho) = d_n^{-1},$$

and boundary conditions

$$f(1, 0, 0; \rho) = 0, \quad f(0, 1, 0; \rho) = 0, \quad f(1, 1, 0; \rho) = 0, \quad f(0, 0, 1; \rho) = 0. \tag{10.2.11}$$

We are interested in $\mathbb{E}W_n = f(0, 0, n; \rho)$. When $n = 50$, representative times are $\mathbb{E}W_{50} = 1.96$ ($\rho = 0$), $= 2.14$ ($\rho = 0.5$), $= 2.36$ ($\rho = 2.0$), $= 2.50$ ($\rho = 10$), $= 2.52$ ($\rho = \infty$).

Hudson and Kaplan (1985) studied the number of recombination events R_n^0 that occur in the history of the sample up to time W_n to ancestors of the sample having material belonging to both marginal trees. Define $f^0(a, b, c\rho)$ to be the expected number of transitions of the form $(a', b', c') \rightarrow (a' + 1, b' + 1, c' - 1)$ until reaching the state $\{(1, 1, 0), (0, 0, 1)\}$, starting from (a, b, c) . By considering the type of the first transition, we see that f^0 satisfies an equation of the form (10.2.10), with

$$g(a, b, c; \rho) = \frac{c\rho}{n(n-1) + c\rho},$$

and boundary conditions (10.2.11). The quantity we want is $\mathbb{E}R_n^0 = f^0(0, 0, n; \rho)$. When $n = 50$, representative values are $\mathbb{E}R_{50}^0 = 0.00$ ($\rho = 0$), $= 2.13$ ($\rho = 0.5$), $= 7.51$ ($\rho = 2.0$), $= 25.6$ ($\rho = 10$).

In contrast, the expected number of recombination events $\mathbb{E}R_n$ in the entire history back to the grand MRCA can be found from the random walk which makes transitions according to

$$\begin{aligned} m \rightarrow m + 1 & \quad \text{with probability } \rho/(\rho + m - 1), \quad m \geq 0 \\ m \rightarrow m - 1 & \quad \text{with probability } (m - 1)/(\rho + m - 1), \quad m \geq 1. \end{aligned}$$

R_n is the number of times the random walk makes a move of the form $m' \rightarrow m' + 1$ before reaching value 1. Standard random walk theory shows that

$$\mathbb{E}R_n = \rho \int_0^1 \frac{1 - (1 - v)^{n-1}}{v} e^{\rho v} dv. \quad (10.2.12)$$

When $n = 50$, representative times are $\mathbb{E}R_{50} = 0.00$ ($\rho = 0$), $= 2.52$ ($\rho = 0.5$), $= 16.2$ ($\rho = 2.0$), $= 24,900$ ($\rho = 10$). A comparison with the values of $\mathbb{E}R_n^0$ shows that $\mathbb{E}R_n$ and $\mathbb{E}R_n^0$ may differ dramatically.

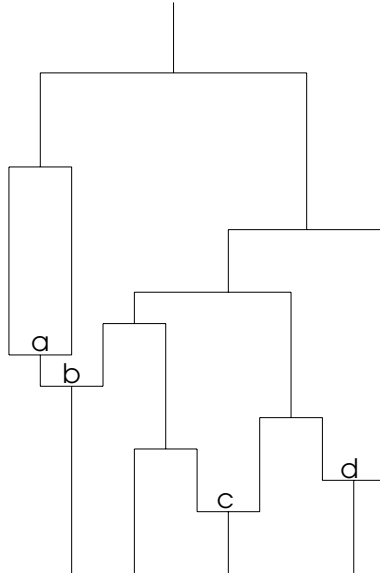
10.3 The continuous recombination model

We now consider a more general class of model in which each chromosome is represented by the unit interval $[0, 1]$. This (and the figures in this section) comes from Griffiths and Marjoram (1997). If a recombination occurs, a position Z for the break point is chosen (independently from other break points) according to a given distribution, and the recombined chromosome is formed from the lengths $[0, Z]$ and $[Z, 1]$ from the first and second parental chromosomes. Other details are as for the 2-locus model. There are several interesting potential choices for the break point distribution Z : Z is constant at 0.5, giving rise to the two-locus model studied earlier; Z is discrete, taking values $\frac{1}{m}, \dots, \frac{m-1}{m}$, giving rise to a m -locus model; and Z has a continuous distribution on $[0, 1]$, where breaks are possible at any point in $[0, 1]$; a particular choice might be the uniform distribution on $[0, 1]$.

As for the 2-locus model we are lead to the concept of ancestral graphs, but now the position at which a recombination occurs is also relevant. Figure 10.2

illustrates an ancestral graph for a sample of $n = 4$ individuals. Positions Z_1, Z_2, \dots where recombination breaks occur are labeled on the graph. The process $A_n^\rho(t)$ which records the number of ancestors of a sample of size n has identical transition rates as the corresponding process for the 2-locus model.

Fig. 10.2. Ancestral recombination graph.

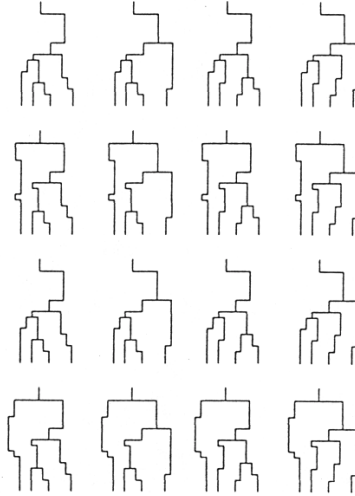


Whereas in the 2-locus model there were two ancestral trees corresponding to the ancestral graph, one for each locus, we now find that each point $x \in [0, 1]$ has an ancestral tree $\mathcal{T}(x)$ associated with it, and marginally each of these trees is described by the coalescent. To obtain $\mathcal{T}(x)$ we trace from the leaves of the ARG upward toward the MRCA. If there is a recombination vertex with label z , we take the left path if $x \leq z$, or right path if $x > z$. The MRCA in $\mathcal{T}(x)$ may occur in the graph before the grand MRCA. Figure 10.2 shows an example of $\mathcal{T}(x)$ when $x > b$ and $x < c, d$.

Since there are a finite number of recombination events in the graph, there are only a finite number of trees in $\{\mathcal{T}(x); x \in [0, 1]\}$. There are potentially 2^R if R recombination events have occurred, but some trees may be identical, or may not exist, depending on the ordering of the recombination break points. Of course (just as before) different trees share edges in the graph, and so are not independently distributed.

Figure 10.4 shows all possible trees corresponding to the ancestral graph in Figure 10.2. Trees 1 and 9 are identical; the other trees are all distinct. If $b > a$ then all trees exist as marginal trees in the graph, otherwise if $b < a$

Fig. 10.4. All possible marginal trees for the graph in Figure 10.2.



10.4 Mutation in the ARG

Mutation is superimposed on the ARG just as it was in the single locus case: Mutations arise at rate $\theta/2$ independently on different branches of the tree, and their effects are modeled by the mutation operator Γ . In the coalescent model with recombination, it often makes no sense to consider mutations that arise on lineages that are lost in the history of the sample due to recombination. Instead, we consider just those mutations which occurred on lineages having material in common with the sample. In the m -locus model, there are now m marginal trees, denoted by $\mathcal{T}_1, \dots, \mathcal{T}_m$. In we denote by $M_n^{(i)}$ the number of mutations occurring on the i th subtree back to its common ancestor, then the total number of mutations is

$$M_n = \sum_{i=1}^m M_n^{(i)}. \tag{10.4.1}$$

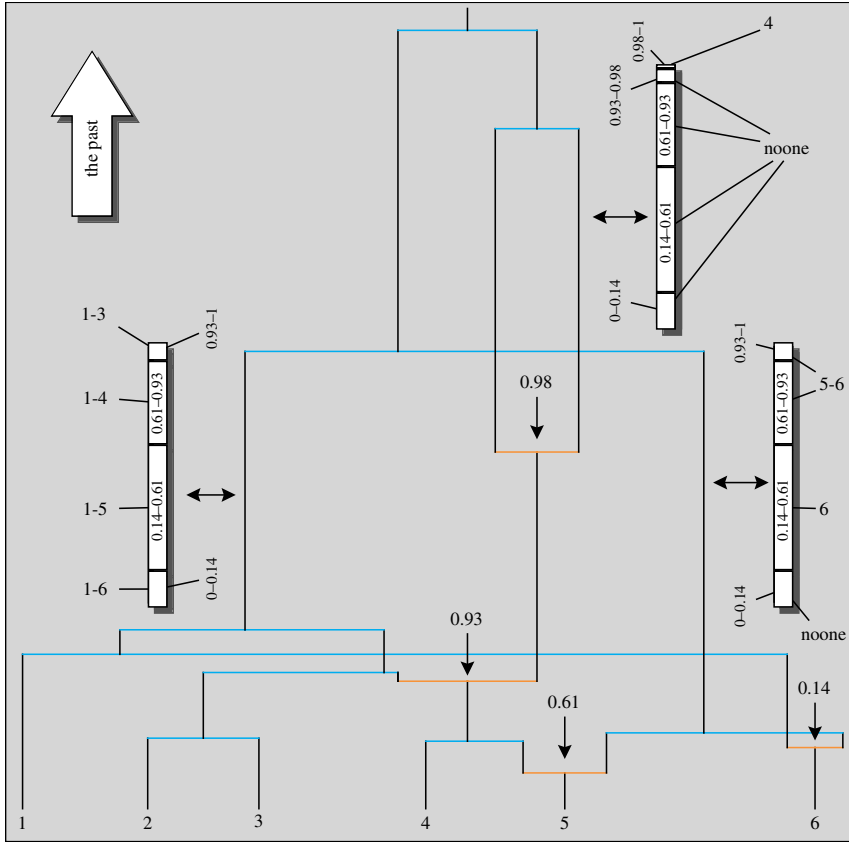
If the mutation rate at each locus is the same, then the overall mutation rate is $\Theta = m\theta$, so that

$$\mathbb{E}M_n = \sum_{i=1}^m \mathbb{E}M_n^{(i)} = \Theta \sum_{j=1}^{n-1} \frac{1}{j}. \tag{10.4.2}$$

Furthermore

$$\text{Var}(M_n) = \sum_{i=1}^m \text{Var}(M_n^{(i)}) + 2 \sum_{i=1}^m \sum_{k=i+1}^m \text{Cov}(M_n^{(i)}, M_n^{(k)}).$$

Fig. 10.5. The history of segments in an ancestral recombination graph



To evaluate the second term Σ_2 , note that conditional on the two marginal subtrees, the mutation processes on those trees are independent. Denoting the tree length at the i th locus by $L_n^{(i)}$, this leads to Hudson's (1983) observation that

$$\text{Cov}(M_n^{(i)}, M_n^{(k)}) = \frac{\theta^2}{4} \text{Cov}(L_n^{(i)}, L_n^{(k)}).$$

In Theorem 10.3 we found the covariance $F_n(\rho) \equiv F(0, 0, n; \rho)$ of the tree lengths in a two locus model with recombination parameter ρ . We can use this to find the covariances in the m -locus model in which the recombination rate ρ between any two adjacent loci is assumed to be the same. The overall recombination rate is $R = (m - 1)\rho$, and for $1 \leq i < k \leq m$, the covariance between $L_n^{(i)}$ and $L_n^{(j)}$ is given by $F_n((k - i)\rho)$. Hence

$$\Sigma_2 = \frac{\theta^2}{2} \sum_{i=1}^{m-1} \sum_{k=i+1}^m F_n((k-i)\rho) = \frac{\theta^2}{2} \sum_{k=1}^{m-1} (m-k)F_n(k\rho).$$

Combining these results, we see that

$$\begin{aligned} \text{Var}(M_n) &= m \left(\theta \sum_{j=1}^{n-1} \frac{1}{j} + \theta^2 \sum_{j=1}^{n-1} \frac{1}{j^2} \right) + \frac{\theta^2}{2} \sum_{k=1}^{m-1} (m-k)F_n(k\rho) \\ &= \Theta \sum_{j=1}^{n-1} \frac{1}{j} + \frac{\Theta}{m} \sum_{j=1}^{n-1} \frac{1}{j^2} + \frac{\Theta^2}{2m} \sum_{k=1}^{m-1} \left(1 - \frac{k}{m}\right) F_n\left(\frac{kR}{m-1}\right). \end{aligned}$$

Hudson considered the limiting case in which $m \rightarrow \infty$ while Θ and R are held fixed. This results in

$$\begin{aligned} \text{Var}(M_n) &= \Theta \sum_{j=1}^{n-1} \frac{1}{j} + \frac{\Theta^2}{2} \int_0^1 (1-w)F_n(Rw)dw \tag{10.4.3} \\ &= \Theta \sum_{j=1}^{n-1} \frac{1}{j} + \frac{1}{2} \frac{\Theta^2}{R^2} \int_0^R (R-w)F_n(w)dw. \end{aligned}$$

10.5 Simulating samples

We consider first the two-locus case. Suppose that there is an overall mutation rate of θ_A at the A locus, and θ_B at the B locus, and let $\theta = \theta_A + \theta_B$. We begin by describing the sequence of mutation, recombination, and coalescence events that occur in the history of the sample back to the MRCA.

Since mutations occur according to independent Poisson processes of rate $\theta/2$ along each lineage, we see that if there are currently m edges in the ancestral graph then the next event on the way back to the MRCA will be a mutation with probability $m\theta/(m(m-1) + m\theta + m\rho) = \theta/(m-1 + \rho + \theta)$, a recombination with probability $\rho/(m-1 + \theta + \rho)$, and a coalescence with probability $(m-1)/(m-1 + \rho + \theta)$. With these events, we may associate a random walk $\{T_k, k \geq 0\}$ which makes transitions according to

$$\begin{aligned} m \rightarrow m + 1 & \text{ with probability } \rho/(\theta + \rho + m - 1), \\ m \rightarrow m & \text{ with probability } \theta/(\theta + \rho + m - 1), \\ m \rightarrow m - 1 & \text{ with probability } (m - 1)/(\theta + \rho + m - 1), \end{aligned}$$

for $m \geq 1$. To describe a sample of size n , the process starts from $T_0 = n$, and ends at the MRCA when $T = 1$.

The effects of each mutation can be modeled in many different ways, for example allowing different combinations of infinitely-many-alleles, infinitely-many-sites, and finitely-many-sites at each locus. In the constant population size model, we can exploit the Markov chain $\{T_k, k \geq 0\}$ to provide an urn

model that can be used to simulate samples efficiently, particularly when the recombination rate ρ is not too large. First we have to generate the sequence of mutation, recombination, and coalescence events *back to the MRCA, starting at the sample*, and then superimpose the effects of each type of event starting at the MRCA and going down to the sample. Here is how this works.

Algorithm 10.1 To simulate from two-locus model.

- (i) Simulate the random walk T_k starting from n until it reaches 1 at step τ . For $k = 1, \dots, \tau$, write $U_k = T_{\tau-k+1} - T_{\tau-k}$.
- (ii) Start by generating the type of the MRCA. For example, for a stationary sample choose the type of this individual according to the stationary distribution of the mutation process. If mutation is independent at each locus this is the product of the stationary distributions of each mutation process.
- (iii) We now use the sequence U_1, U_2, \dots, U_τ (in that order) to generate the sample. For $k = 1, 2, \dots, \tau$:
 - If $U_k = -1$ then a recombination event has occurred. Choose two individuals at random without replacement from the current individuals, and recombine them. The first individual chosen contributes the A locus allele, the second the B locus allele.
 - If $U_k = 0$, a mutation has occurred. Choose an individual at random and generate a mutation. With probability θ_A/θ the mutation occurs at the A locus, in which case a transition is made according to the mutation distribution $\Gamma^A(x, \cdot)$ if the type is currently x , and similarly for the B locus.
 - If $U_k = 1$, then a coalescence has occurred. Choose an individual at random and duplicate its type.
- (iv) After τ steps of the process, the sample has size n and the distribution of the sample is just what we wanted.

It can be seen that the efficiency of this algorithm depends on the expected value of τ . When either ρ or θ is large, $\mathbb{E}\tau$ can be very large, making the simulation quite slow.

This method extends directly to simulations of samples from the general ARG. Once the locations of the recombination events have been simulated according to Algorithm 10.1, we can choose recombination break points according to any prescribed distribution on $[0,1]$. Essentially any mutation mechanism can be modeled too. For example, for the infinitely-many-sites model we can suppose that mutations occur according to a continuous distribution on $(0,1)$, and that the label of a mutation is just the position at which it occurs. In the case of variable population size this method does not work, and the ancestral recombination graph needs to be simulated first, and then mutations are superimposed from the MRCAs. Hudson (1991) is a useful reference.

10.6 Linkage disequilibrium and haplotype sharing

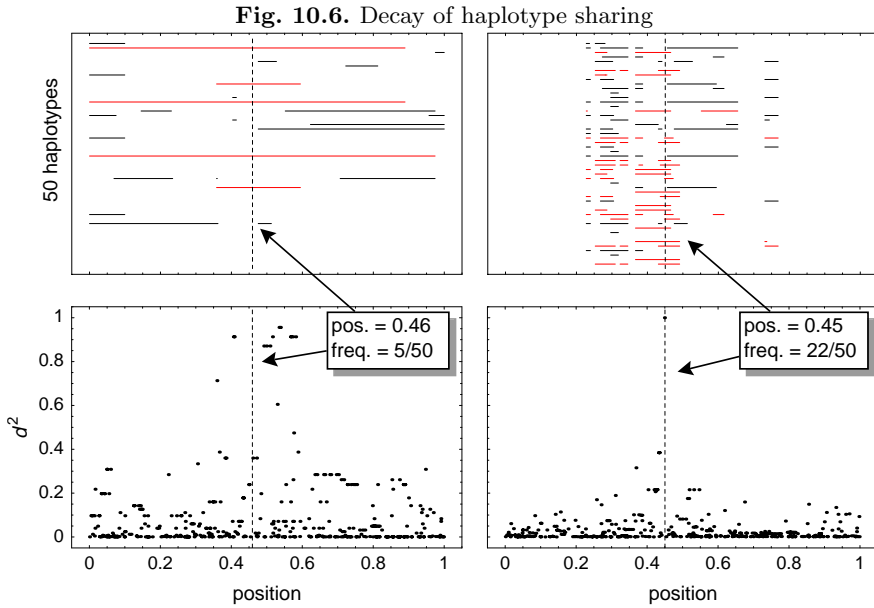
Because the genealogical trees at different linked positions in a segment are not independent of one another, neither will be the allelic states of these loci – there will be *linkage disequilibrium* (LD) between the loci. LD is usually quantified by using various measures of association between pairs of loci. Consider two such loci, each of which has two possible alleles, and denote the relative frequency of the $A_i B_j$ haplotype by $p(A_i, B_j)$, and let $p(A_i), p(B_j)$ denote the relative frequency of each allele. Among the pairwise measures of LD are

- D' , the value of $D = p(A_1, B_1) - p(A_1)p(B_1)$, normalized to have values between -1 and 1 regardless of allele frequencies;
- r^2 , the correlation in allelic state between the two loci as they occur in haplotypes;
- $d^2 = (p(B_1 | A_2) - p(B_1 | A_1))^2$, which measures the association between the alleles at (marker) locus B and the alleles at (disease) locus A .

These and other measures of LD are discussed further in Guo (1997), Hudson (2001) and Pritchard and Przeworski (2001).

Because of the history of recombination and mutation in a sample, pairwise LD is expected to be extremely variable. This is illustrated in Figure 10.6, adapted from Nordborg and Tavaré (2002). The horizontal axis, which represents chromosomal position, corresponds to roughly 100 kb. The plots illustrate the haplotype sharing and LD with respect to particular focal mutations. In the left column, a relatively low-frequency mutation (5/50=10%) was chosen as focus, and in the right column, a relatively high-frequency one (22/50=44%). The chromosomal position of these mutations are indicated by the vertical lines. The top row of plots shows the extent of haplotype sharing with respect to the MRCA of the focal mutation among the 50 haplotypes. The horizontal lines indicate segments that descend from the MRCA of the focal mutation. Light lines indicates that the current haplotype also carries the focal mutation, dark lines that it does not. Note that the light segments necessarily overlap the position of the focal mutation. For clarity, segments that do not descend from the MRCA of the focal mutation are not shown at all, and haplotypes that do not carry segments descended from the MRCA of the focal mutation are therefore invisible. The second row of plots shows the behavior of LD as measured by d^2 for different choices of markers. In each plot, the horizontal position of a dot represents the chromosomal position of the marker, and the vertical position the value of the measure (on a zero-to-one scale).

Because of interest in mapping disease susceptibility genes, the extent of LD across the human genome has been much debated. What is clear is that while there is a relationship between LD measures and distance, the inherent variability in LD makes this relationship hard to infer. In particular, it is difficult to compare studies that use different measures of pairwise LD as these measures can differ dramatically in their estimates of the range of LD.



For reviews of these issues in relation to mapping, see for example Clayton (2000), Weiss and Clark (2002), Nordborg and Tavaré (2002) and Ardlie et al. (2002).

Estimating the recombination fraction

There is a sizable literature on estimation of the scaled recombination rate ρ , among them methods that use summary statistics of the data such as Hudson (1987), Hey and Wakeley (1997), and Wakeley (1997). Griffiths and Marjoram (1996) and Fearnhead and Donnelly (2001) exploit the importance sampling approach developed in Section 6 for the infinitely-many-sites model, while Nielsen (2000) and Kuhner et al. (2000) use MCMC methods, the latter specifically for DNA sequence data. Wall (2000) has performed an extensive comparison of these approaches. One conclusion is that (reliable) estimation of pairwise recombination fractions is extremely difficult. See Fearnhead and Donnelly (2002) for another approach, and Morris et al. (2002) and the references contained therein for approaches to mapping disease genes using the coalescent.

11 ABC: Approximate Bayesian Computation

Several of the previous sections have described methods for simulating observations from a posterior distribution. One key ingredient in these methods is the likelihood function; we have until now assumed this could be computed numerically, for example using the peeling algorithm described in Section 9.4. In this section we describe some methods that can be used when likelihoods are hard or impossible to compute.

In this section, data \mathcal{D} are generated from a model \mathcal{M} determined by parameters θ . We denote the prior for θ by $\pi(\theta)$. The posterior distribution of interest is $f(\theta | \mathcal{D})$ given by

$$f(\theta | \mathcal{D}) = \mathbb{P}(\mathcal{D} | \theta)\pi(\theta)/\mathbb{P}(\mathcal{D}),$$

where $\mathbb{P}(\mathcal{D}) = \int \mathbb{P}(\mathcal{D} | \theta)\pi(\theta) d\theta$ is the normalizing constant.

11.1 Rejection methods

We have already seen examples of the rejection method for discrete data:

Algorithm 11.1

1. Generate θ from $\pi(\cdot)$
2. Accept θ with probability $h = \mathbb{P}(\mathcal{D} | \theta)$, and return to 1.

It is easy to see that accepted observations have distribution $f(\theta | \mathcal{D})$, as shown for example in Ripley (1987). As we saw in Section 7.3, the computations can often be speeded up if there is constant c such that $\mathbb{P}(\mathcal{D} | \theta) \leq c$ for all θ . h can then be replaced by h/c .

There are many variations on this theme. Of particular relevance here is the case in which the likelihood $\mathbb{P}(\mathcal{D} | \theta)$ cannot be computed explicitly. One approach is then the following:

Algorithm 11.2

1. Generate θ from $\pi(\cdot)$
2. Simulate \mathcal{D}' from model \mathcal{M} with parameter θ
3. Accept θ if $\mathcal{D}' = \mathcal{D}$, and return to 1.

The success of this approach depends on the fact that the underlying stochastic process \mathcal{M} is easy to simulate for a given set of parameters. We note also that this approach can be useful when explicit computation of the likelihood is possible but time consuming.

The practicality of algorithms like these depends crucially on the size of $\mathbb{P}(\mathcal{D})$, because the probability of accepting an observation is proportional to

$\mathbb{P}(\mathcal{D})$. In cases where the acceptance rate is too small, one might resort to approximate methods such as the following:

Algorithm 11.3

1. Generate θ from $\pi(\cdot)$
2. Simulate \mathcal{D}' from model \mathcal{M} with parameter θ
3. Calculate a measure of distance $\rho(\mathcal{D}, \mathcal{D}')$ between \mathcal{D}' and \mathcal{D}
4. Accept θ if $\rho \leq \epsilon$, and return to 1.

This approach requires selection of a suitable metric ρ as well as a choice of ϵ . As $\epsilon \rightarrow \infty$, it generates observations from the prior, and as $\epsilon \rightarrow 0$, it generates observations from the required density $f(\theta \mid \mathcal{D})$. The choice of ϵ reflects the interplay between computability and accuracy. For a given ρ and ϵ accepted observations are independent and identically distributed from $f(\theta \mid \rho(\mathcal{D}, \mathcal{D}') \leq \epsilon)$.

11.2 Inference in the fossil record

In this section, we give an application of Algorithm 11.3 to a problem concerning estimation of the time to the most recent common ancestor of primates. Our inference is based not on molecular data but on a sampling of the fossil record itself.

The problem

In Table 17 the number of primate species found as fossils in a series of stratigraphic intervals is given. Tavaré *et al.* (2002) developed a statistical method for estimating the temporal gap between the base of the stratigraphic interval in which the oldest fossil was found and the initial point of divergence of the species in the sample. The bias in the estimators and approximate confidence intervals for the parameters were found by using a parametric bootstrap approach. Estimates of the divergence time of primates (more accurately, the time of the haplorhine-strepsirrhine split) based on molecular sequence data give a time of about 90 million years. A literal interpretation of the fossil record suggests a divergence time of about 60 million years. One reason for the present studies is to reconcile these two estimates. A more detailed account of the problem is given in Soligo *et al.* (2002).

A model for speciation and sampling

We adopt the same framework as in Tavaré *et al.* (2002). We model speciation with a non-homogeneous Markov birth-and-death process. To model evolution from the last common ancestor of all living and fossil species included in the

Table 17. Data for the primate fossil record. References can be found in the supplemental material in Tavaré *et al.* (2002).

Epoch	k	T_k	Observed number of species (D_k)
Late Pleistocene	1	0.15	19
Middle Pleistocene	2	0.9	28
Early Pleistocene	3	1.8	22
Late Pliocene	4	3.6	47
Early Pliocene	5	5.3	11
Late Miocene	6	11.2	38
Middle Miocene	7	16.4	46
Early Miocene	8	23.8	36
Late Oligocene	9	28.5	4
Early Oligocene	10	33.7	20
Late Eocene	11	37.0	32
Middle Eocene	12	49.0	103
Early Eocene	13	54.8	68
Pre-Eocene	14		0

analysis, we start with two species at time 0. Species go extinct at rate λ , and so have exponential lifetimes with mean $1/\lambda$, time being measured in millions of years. A species that goes extinct at time u is replaced by an average of $m(u)$ new species. We denote by Z_t the number of species alive at time t . The expected number of species extant at time t is given by

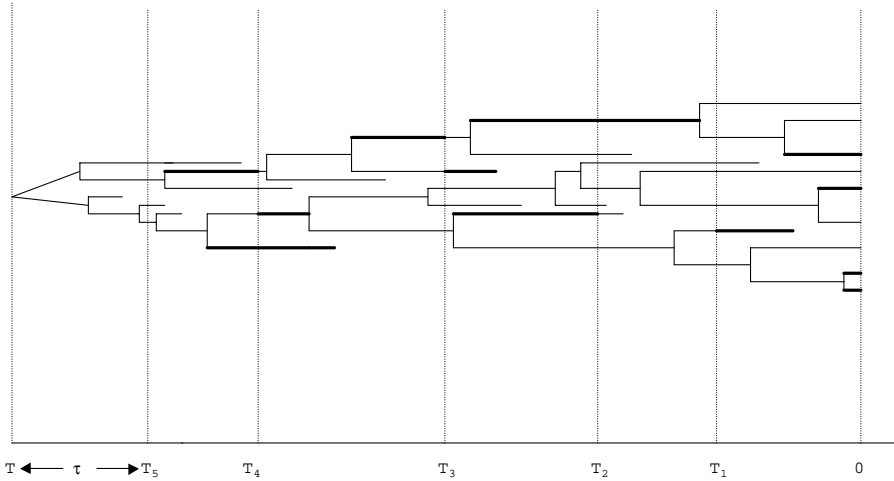
$$\mathbb{E}Z_t = 2 \exp \left\{ \lambda \int_0^t (m(u) - 1) du \right\}; \quad (11.2.1)$$

cf. Harris (1963), Chapter 5. Furthermore, if $B(s, t]$ denotes the number of species born in the interval $(s, t]$, then

$$\mathbb{E}B[s, t] = \lambda \int_s^t m(u) \mathbb{E}Z_u du, \quad s < t. \quad (11.2.2)$$

We divide time into k stratigraphic intervals, following this sequence (see Table 17 and Figure 11.1). The base of the first (youngest) stratigraphic interval is at T_1 mya and the base of the k^{th} is at T_k million years ago (mya). The earliest known fossil is found in this interval. The founding species originate at time $T := T_k + \tau$ mya, and we define a $(k+1)^{\text{st}}$ stratigraphic interval that has its base at $T_{k+1} := T$ mya and ends T_k mya. Note that no fossils have

Fig. 11.1. An illustration of the stochastic model of fossil finds. Bases of 5 stratigraphic intervals at T_1, \dots, T_5 mya are shown along the x-axis. The temporal gap between the base of the final interval and the point at which the two founding species originate is denoted by τ . Thick lines indicate species found in the fossil record. Time 0 is the present day.



been found in this interval. We wish to approximate the posterior distribution of the time τ and other parameters of the model, using as data the number of different species found in the fossil record in the first, second, \dots , k^{th} intervals. We model the number of species alive u mya by the value Z_{T-u} of the Markov branching process described earlier.

The number N_j of distinct species living in the j th stratigraphic interval having base T_j mya is the sum of those that were extant at the beginning of the interval, Z_{T-T_j} , plus those that originated in the interval, $B[T-T_j, T-T_{j-1}]$. It follows from (11.2.1) and (11.2.2) that the expected number of distinct species that can be sampled in the j th stratigraphic interval is

$$\mathbb{E}N_j = \mathbb{E}Z_{T-T_{j-1}} + \lambda \int_{T-T_j}^{T-T_{j-1}} \mathbb{E}Z_u du, \quad j = 1, \dots, k+1. \quad (11.2.3)$$

We assume that, conditional on the number of distinct species N_j that lived in the j th stratigraphic interval, the number of species D_j actually found in the fossil record in this interval is a binomial random variable with parameters N_j and α_j , $j = 1, 2, \dots, k$. Furthermore, the D_j are assumed to be conditionally independent given the N_j . The parameter α_j gives the probability of

sampling a species in the j th stratigraphic interval. A typical data set is given in Table 17.

A Bayesian approach

We write $\mathcal{D} = (D_1, \dots, D_{k+1})$ for the counts observed in the $k+1$ stratigraphic intervals, and we write θ for the vector of parameters of the process, one of which is τ , the temporal gap. The likelihood can be written in the form

$$\mathbb{P}(\mathcal{D} \mid \theta) = \mathbb{E} \prod_{j=1}^{k+1} \binom{N_j}{D_j} \alpha_j^{D_j} (1 - \alpha_j)^{N_j - D_j}, \quad (11.2.4)$$

where the expectation is over trajectories of the speciation process Z that run for time T with parameter θ , and such that both initial branches have offspring species surviving to time T . By convention the term under the expectation sign is 0 if any $D_j > N_j$.

While the acceptance probability is difficult to compute, the stochastic process itself can be simulated easily, and Algorithm 11.3 comes into play. One crucial aspect of this method is the choice of ρ in Algorithm 11.3. The counts D_1, \dots, D_{k+1} can be represented as the total number of fossils found,

$$D_+ = D_1 + \dots + D_{k+1},$$

and a vector of proportions

$$(Q_1, \dots, Q_{k+1}) := \left(\frac{D_1}{D_+}, \dots, \frac{D_{k+1}}{D_+} \right).$$

We can therefore measure the distance between \mathcal{D} and a simulated data set \mathcal{D}' by

$$\rho(\mathcal{D}, \mathcal{D}') = \left| \frac{D'_+}{D_+} - 1 \right| + \frac{1}{2} \sum_{j=1}^{k+1} |Q_j - Q'_j|. \quad (11.2.5)$$

The first term measures the relative error in the total number of fossils found in a simulated data set and the actual number, while the second term is the total variation distance between the two vectors of proportions.

Results

Tavaré *et al.* (2002) modelled the mean diversification via the logistic function, for which

$$\mathbb{E}Z_t = 2/\{\gamma + (1 - \gamma)e^{-\rho t}\}. \quad (11.2.6)$$

This form is quite flexible; for example, $\gamma = 0$ corresponds to exponential growth. They equated the expected number of species known at the present time with the observed number, and also specified the time at which the mean

diversification reached 90% of its current value. These two equations serve to determine the form of the speciation curve. They also assumed a mean species lifetime of 2.5 my (although their results were little changed by assuming a 2 my or 3 my lifetime). They modelled the sampling fractions α_j in the form

$$\alpha_j = \alpha p_j, \quad j = 1, 2, \dots, k + 1, \quad (11.2.7)$$

where the p_j are known proportions, and α is a scale parameter to be estimated from the data. The particular values of the p_j they used are given in Table 18. The average value is $\bar{p} = 0.73$.

Table 18. Sampling proportions p_j

j	1	2	3	4	5	6	7	8	9	10	11	12	13	14
p_j	1.0	1.0	1.0	1.0	0.5	0.5	1.0	0.5	0.1	0.5	1.0	1.0	1.0	0.1

Using the data from Table 17, they estimated a temporal gap of 26.7 my with an approximate 95% confidence interval of 17.2 my to 34.8 my. As the oldest known fossil primate is 54.8 my old, this is equivalent to an age of 81.5 my for the last common ancestor of living primates. The average sampling fraction $\bar{\alpha}$, defined as

$$\bar{\alpha} = \alpha \bar{p} \quad (11.2.8)$$

was estimated to be 5.7% with an upper 95% confidence limit of 7.4%.

For comparison with the earlier approach, we treat both ρ and γ as fixed parameters, so that the parameter θ is given by $\theta = (\tau, \alpha)$. The prior distributions were chosen as

$$\tau \sim U(0, 100)$$

$$\alpha \sim U(0, 0.3)$$

the notation $U(a, b)$ denotes the uniform density on (a, b) . In Tavaré *et al.* (2002), we used fixed values of $\rho = 0.2995$, $\gamma = 0.0085$. From 500 accepted observations with $\epsilon = 0.1$, we obtain the summaries in Figure 11.2 and Table 19. A median value of 27.6 my for the posterior value of the temporal gap τ is very close to that estimated in the previous analysis (Tavaré *et al.* (2002)) and is equivalent to an age of 82.4 my for the last common ancestor of living primates. The 2.5% and 97.5% points of the posterior of τ are estimated to be 15.4 my and 57.9 my, and the 95% point of the posterior for $\bar{\alpha}$ is 10%; these values are all broadly consistent with the previously published analysis. The posterior distribution of the number of present-day species serves as a

Fig. 11.2. Left panel: posterior for τ . Right panel: posterior for $\bar{\alpha}$.

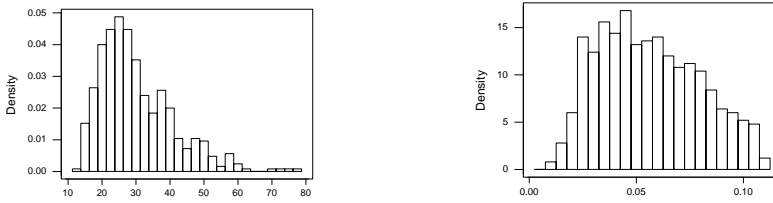


Table 19. Summary statistics for τ , $\bar{\alpha}$ and N_0 when ρ and γ are fixed.

	τ	$\bar{\alpha}(\%)$	N_0
25th percentile	22.4	3.7	180
median	27.6	5.4	253
mean	30.1	5.7	294
75th percentile	36.6	7.5	357

goodness-of-fit assessment. The observed number of extant primates, 235, is clearly a typical value under the posterior.

The analysis here can be compared to a full MCMC approach. The results are essentially indistinguishable; see Plagnol and Tavaré (2003) for further details. One advantage of approximate Bayesian approaches are their flexibility. A number of other scenarios, such as different species diversity curves and sampling schemes, can be examined quickly. For further details, see Will *et al.* (2003).

11.3 Using summary statistics

In Section 7 we found the posterior distribution conditional on a summary statistic rather than the full sequence data. The motivating idea behind this is that if the set of statistics $S = (S_1, \dots, S_p)$ is sufficient for θ , in that $\mathbb{P}(\mathcal{D} \mid S, \theta)$ is independent of θ , then $f(\theta \mid \mathcal{D}) = f(\theta \mid S)$. The normalizing constant is now $\mathbb{P}(S)$ which is typically larger than $\mathbb{P}(\mathcal{D})$, resulting in more acceptances.

In practice it is be hard, if not impossible, to identity a suitable set of sufficient statistics, and we might then resort to a more heuristic approach that uses knowledge of the particular problem at hand to suggest summary statistics that capture information about θ . With these statistics in hand,

we have the following approximate Bayesian computation scheme for data \mathcal{D} summarized by S :

Algorithm 11.4

1. Generate θ from $\pi(\cdot)$
2. Simulate \mathcal{D}' from model \mathcal{M} with parameter θ , and compute the corresponding statistics S'
3. Calculate the distance $\rho(S, S')$ between S and S'
4. Accept θ if $\rho \leq \epsilon$, and return to 1.

Examples of this algorithm approach appear frequently in the population genetics literature, including Fu and Li (1997), Weiss and von Haeseler (1998), Pritchard *et al.* (1999) and Wall (2000). Beaumont *et al.* (2002) describes a novel generalization of the rejection method in which all observations generated in steps 1 and 2 of Algorithm 11.4 are used in a local-linear regression framework to improve the simulation output. They also describe a number of other examples of this approach.

11.4 MCMC methods

There are several advantages to these rejection methods: they are usually easy to code, they generate independent observations (and so can use embarrassingly parallel computation), and they readily provide estimates of Bayes factors, which can be used for model comparison. On the other hand, for complex probability models sampling from the prior does not make good use of accepted observations, so these methods can be prohibitively slow. Here we describe an MCMC approach to problems in which the likelihood cannot be readily computed.

As we saw in Section 9, the Metropolis-Hastings Algorithm for generating observations from $f(\theta \mid \mathcal{D})$ uses output from a Markov chain. It can be described as follows:

Algorithm 11.5

1. If at θ , propose a move to θ' according to a transition kernel $q(\theta \rightarrow \theta')$
2. Calculate

$$h = \min \left(1, \frac{\mathbb{P}(\mathcal{D} \mid \theta')\pi(\theta')q(\theta' \rightarrow \theta)}{\mathbb{P}(\mathcal{D} \mid \theta)\pi(\theta)q(\theta \rightarrow \theta')} \right)$$
3. Move to θ' with probability h , else stay at θ ; go to 1.

In Marjoram *et al.* (2003) we describe an MCMC approach that is the natural analog of Algorithm 11.4, in that no likelihoods are used (or estimated) in its implementation. It is based on the following steps:

Algorithm 11.6

1. If at θ propose a move to θ' according to a transition kernel $q(\theta \rightarrow \theta')$
2. Generate \mathcal{D}' using model \mathcal{M} with parameters θ'
3. If $\rho(S', S) \leq \epsilon$, go to 4, and otherwise stay at θ and return to 1,
4. Calculate

$$h = h(\theta, \theta') = \min \left(1, \frac{\pi(\theta')q(\theta' \rightarrow \theta)}{\pi(\theta)q(\theta \rightarrow \theta')} \right)$$

5. Move to θ' with probability h , else stay at θ ; go to 1.

The stationary distribution of the chain is indeed $f(\theta \mid \rho(S', S) \leq \epsilon)$. Applications of this approach to inference about mutation rates are given in Marjoram *et al.* (2003) and Plagnol and Tavaré (2003). The method usually has to be implemented by including part of the underlying coalescent tree and the mutation process as part of the MCMC update (making it part of θ , as it were).

The method seems to allow some flexibility in studying problems where existing methods don't work well in practice, such as analyzing complex models of mutation and analyzing restriction fragment length polymorphism data. There is a need for research on implementable methods for identifying approximately sufficient statistics, and for the development of more sophisticated MCMC methods that do not use likelihoods. Such approaches will be necessary when addressing problems involving high-dimensional parameters.

11.5 The genealogy of a branching process

Thus far the genealogical processes used in these notes have been evolving in continuous time. In this section, we describe informally a method for generating the genealogical history of a sample of individuals evolving according to a discrete-time branching process.

The conventional way to describe the evolution of a Galton-Watson process is as a series of population sizes Z_0, Z_1, Z_2, \dots at times $0, 1, 2, \dots$. The number of individuals Z_{m+1} is a random sum:

$$Z_{m+1} = \sum_{j=1}^{Z_m} \xi_{mj},$$

where $\xi_{mj}, j \geq 1$ are identically distributed random variables having a distribution that may depend on m . A more detailed description of the process gives the number of families F_{mk} born into generation m that have k members, $k = 0, 1, 2, \dots$. Given Z_{m-1} , the joint distribution of $F_{mj}, j \geq 0$ is multinomial with sample size Z_{m-1} and $q_{m-1,k}, k \geq 0$; here,

$$q_{m-1,k} = \mathbb{P}(\xi_{m-1,1} = k).$$

To simulate the genealogy of a random sample of size n from generation g of the process we proceed as follows; cf. Weiss and von Haeseler (1997). Starting from Z_0 individuals, generate the family size statistics $F_{1k}, k \geq 0$. These determine Z_1 , after which the family sizes $F_{2k}, k \geq 0$ can be generated. Continuing in this way we finally generate the family sizes $F_{gk}, k \geq 0$. This done, a random subtree with n leaves can be generated backwards from generation g as follows. Randomly choose n individuals without replacement, recording which family they belong to (there being F_{gk} families of size k). Count the number A of families represented, each one corresponding to a distinct ancestor in generation $g - 1$. Next, sample A individuals from generation $g - 1$ and record which families they belong to (there now being $F_{g-1,k}$ families of size k), and so on. Iterating this scheme back through generations $g - 1, g - 2, \dots, 1$ produces a genealogical tree having the required distribution.

Versions of this scheme have been used to study the polymerase chain reaction by Weiss and von Haeseler (1997), and to estimate the time to loss of mismatch repair in a colon tumor by Tsao *et al.* (2000) and Tavaré (2004). In both examples, the effects of a mutation process are superimposed on the genealogy, thereby generating sample data. Because the simulated genealogies are relatively quick to produce, they can be used for statistical inference such as implementations of Algorithm 11.4.

Finally we note that the simulation scheme can be used in much more general settings. For example, the distribution $q_{mj}, j \geq 0$ can depend on the history of the process in generations $0, 1, \dots, m$; this covers cases of density dependent reproduction. This approach can also be applied to multitype branching processes.

12 Afterwords

This section concludes the lecture notes by giving some pointers to topics that were mentioned in the Saint Flour lectures, but have not been written up for the printed version.

12.1 The effects of selection

The previous sections have focussed on neutral genes, in which the effects of mutation could be superimposed on the underlying coalescent genealogy. When selection is acting at some loci, this separation is no longer possible and the analysis is rather more complicated.

Two basic approaches have emerged. In the first the evolution of the selected loci is modelled forward in time, and then the neutral loci are studied by coalescent methods (cf. Kaplan *et al.* (1988, 1989)). In the second a genealogical process known as the *ancestral selection graph*, the analog of the neutral coalescent, is developed by Neuhauser and Krone (1997) and Krone and Neuhauser (1997). See Neuhauser (2001), Neuhauser and Tavaré (2002) and Nordborg (2001) for reviews. Methods for simulating selected genealogies are an important current area of research; see Slatkin (2001), Slade (2000, 2001) and Fearnhead (2001) for some examples. Such simulations can be used to explore the consequences of different selection mechanisms on the pattern of variation observed in data. Methods for detecting selection in sequence data are reviewed in Kreitman (2000). Methods for inference and estimation using coalescent methods are an active area of research. For an introduction to models with spatial structure, see Nordborg (2001) for example.

12.2 The combinatorics connection

Mathematical population genetics in the guise of the Ewens Sampling Formula (3.5.3) and Poisson approximation intersect in an area of probabilistic combinatorics. This leads directly to an extremely powerful and flexible method for studying the asymptotic behavior of decomposable combinatorial structures such as permutations, polynomials over a finite field, and random mappings. The joint distribution of counts of components of different sizes can be represented as the distribution of independent random variables conditional on a weighted sum; recall (3.5.4). Consequences of this representation are exploited in Arratia and Tavaré (1994). Connections with prime factorization are outlined in the expository article of Arratia *et al.* (1997). The book of Arratia, Barbour and Tavaré (2003) provides a detailed account of the theory, which places the Ewens Sampling Formula in much the same position as the Normal distribution in the central limit theorem: informally, many decomposable combinatorial models behave asymptotically like the Ewens Sampling Formula, and the closeness of the approximation can be measured in the total variation metric. A preprint of the book can be found at

<http://www-hto.usc.edu/books/tavare/ABT/index.html>

Pitman's lecture notes, *Combinatorial Stochastic Processes*, from the 2002 Saint Flour lectures contains related material. A draft may be obtained from <http://stat-www.berkeley.edu/users/pitman/bibliog.html>

12.3 Bugs and features

Errors and typos will be reported at

<http://www-hto.usc.edu/papers/abstracts/coalescent.html>

I also intend to make a set of exercises available there.

I leave it to Søren Kierkegaard (1813-1855) to summarize why coalescents are interesting and why these notes end here:

Life can only be understood going backwards, but it must be lived going forwards.

References

1. S. Anderson, A. T. Bankier, B. G. Barrell, M. H. de Bruijn, A. R. Coulson, J. Drouin, I. C. Eperon, D. P. Nierlich, B. A. Roe, F. Sanger, P. H. Schreier, A. J. Smith, R. Staden, and I. G. Young. Sequence and organization of the human mitochondrial genome. *Nature*, 290:457–465, 1981.
2. K. G. Ardlie, L. Kruglyak, and M. Seielstad. Patterns of linkage disequilibrium in the human genome. *Nature Rev. Genet.*, 3:299–309, 2002.
3. R. Arratia, A. D. Barbour, and S. Tavaré. Random combinatorial structures and prime factorizations. *Notices of the AMS*, 44:903–910, 1997.
4. R. Arratia, A. D. Barbour, and S. Tavaré. *Logarithmic Combinatorial Structures: a Probabilistic Approach*. European Mathematical Society Publishing House, 2003. In press.
5. R. Arratia and S. Tavaré. Independent process approximations for random combinatorial structures. *Adv. Math.*, 104:90–154, 1994.
6. A. D. Barbour, L. Holst, and S. Janson. *Poisson Approximation*. Oxford University Press, Oxford, 1992.
7. M. A. Beaumont, W. Zhang, and D. J. Balding. Approximate Bayesian computation in population genetics. *Genetics*, 162:2025–2035, 2002.
8. N. G. Best, M. K. Cowles, and S. K. Vines. *CODA Manual version 0.30*. MRC Biostatistics Unit., Cambridge, UK, 1995.
9. T. A. Brown. *Genomes*. John Wiley & Sons, New York, New York, 1999.
10. P. Buneman. The recovery of trees from measures of dissimilarity. In D. G. Kendall and P. Tautu, editors, *Mathematics in the Archaeological and Historical Sciences*, pages 387–395. Edinburgh University Press, 1971.
11. C. Cannings. The latent roots of certain Markov chains arising in genetics: A new approach. I. Haploid models. *Adv. Appl. Prob.*, 6:260–290, 1974.
12. D. Clayton. Linkage disequilibrium mapping of disease susceptibility genes in human populations. *International Statistical Review*, 68:23–43, 2000.
13. J. F. Crow and M. Kimura. *An Introduction to Population Genetics Theory*. Harper and Row, New York, 1970.

14. P. Donnelly and T. G. Kurtz. The asymptotic behavior of an urn model arising in population genetics. *Stochast. Process. Applic.*, 64:1–16, 1996.
15. P. Donnelly and S. Tavaré. Coalescents and genealogical structure under neutrality. *Ann. Rev. Genet.*, 29:401–421, 1995.
16. P. Donnelly, S. Tavaré, D. J. Balding, and R. C. Griffiths. Estimating the age of the common ancestor of men from the ZFY intron. *Science*, 257:1357–1359, 1996.
17. R. L. Dorit, H. Akashi, and W. Gilbert. Estimating the age of the common ancestor of men from the ZFY intron. *Science*, 272:1361–1362, 1996.
18. S. N. Ethier and T. G. Kurtz. *Markov Processes. Characterization and Convergence*. John Wiley & Sons, Inc., New York, 1986.
19. S.N. Ethier and R.C. Griffiths. The infinitely-many-sites model as a measure valued diffusion. *Ann. Probab.*, 15:515–545, 1987.
20. S.N. Ethier and R.C. Griffiths. On the two-locus sampling distribution. *J. Math. Biol.*, 29:131–159, 1990.
21. W. J. Ewens. The sampling theory of selectively neutral alleles. *Theoret. Popn. Biol.*, 3:87–112, 1972.
22. W. J. Ewens. *Mathematical Population Genetics*. Springer-Verlag, Berlin, Heidelberg, New York, 1979.
23. W. J. Ewens. Population genetics theory - the past and the future. In S. Lessard, editor, *Mathematical and statistical developments of evolutionary theory*. Kluwer Academic Publishers, 1990.
24. W. J. Ewens and S. Tavaré. The Ewens Sampling Formula. In S. Kotz, C. B. Read, and D. L. Banks, editors, *Encyclopedia of Statistical Science*, Volume 2, pages 230–234. Wiley, New York, 1998.
25. P. Fearnhead. Perfect simulation from population genetic models with selection. *Theoret. Popul. Biol.*, 59:263–279, 2001.
26. P. Fearnhead and P. Donnelly. Estimating recombination rates from population genetic data. *Genetics*, 159:1299–1318, 2001.
27. P. Fearnhead and P. Donnelly. Approximate likelihood methods for estimating local recombination rates. *J. Royal Statist. Soc. B*, 64:657–680, 2002.
28. J. Felsenstein. The rate of loss of multiple alleles in finite haploid populations. *Theoret. Popn. Biol.*, 2:391–403, 1971.
29. J. Felsenstein. Maximum-likelihood and minimum-steps methods for estimating evolutionary trees from data on discrete characters. *Syst. Zool.*, 22:240–249, 1973.
30. J. Felsenstein. Evolutionary trees from DNA sequence data: a maximum likelihood approach. *J. Mol. Evol.*, 17:368–376, 1981.
31. J. Felsenstein, M. Kuhner, J. Yamato, and P. Beerli. Likelihoods on coalescents: a Monte Carlo sampling approach to inferring parameters from population samples of molecular data. In F. Seillier-Moisewitsch, editor, *Statistics in Molecular Biology and Genetics*, pages 163–185. Institute of Mathematical Statistics and American Mathematical Society, Hayward, California, 1999.
32. R. A. Fisher. On the dominance ratio. *Proc. Roy. Soc. Edin.*, 42:321–431, 1922.
33. G. E. Forsythe and R. A. Leibler. Matrix inversion by the Monte Carlo method. *Math. Comp.*, 26:127–129, 1950.
34. Y.-X. Fu. Estimating the age of the common ancestor of a DNA sample using the number of segregating sites. *Genetics*, 144:829–838, 1996.

35. Y.-X. Fu and W.-H. Li. Estimating the age of the common ancestor of men from the ZFY intron. *Science*, 257:1356–1357, 1996.
36. Y.-X. Fu and W.-H. Li. Estimating the age of the common ancestor of a sample of DNA sequences. *Mol. Biol. Evol.*, 14:195–199, 1997.
37. Y.-X. Fu and W.-H. Li. Coalescing into the 21st century: An overview and prospects of coalescent theory. *Theoret. Popul. Biol.*, 56:1–10, 1999.
38. W. R. Gilks, S. Richardson, and D. J. Spiegelhalter. *Markov Chain Monte Carlo in Practice*. Chapman and Hall, 1996.
39. K. Gladstien. The characteristic values and vectors for a class of stochastic matrices arising in genetics. *SIAM J. Appl. Math.*, 34:630–642, 1978.
40. R. C. Griffiths. Lines of descent in the diffusion approximation of neutral Wright-Fisher models. *Theor. Pop. Biol.*, 17:37–50, 1980.
41. R. C. Griffiths. Neutral two-locus multiple allele models with recombination. *Theoret. Popn. Biol.*, 19:169–186, 1981.
42. R. C. Griffiths. Genealogical-tree probabilities in the infinitely-many-site model. *J. Math. Biol.*, 27:667–680, 1989.
43. R. C. Griffiths. The two-locus ancestral graph. In I.V. Basawa and R.L. Taylor, editors, *Selected Proceedings of the Symposium on Applied Probability, Sheffield, 1989*, pages 100–117. Institute of Mathematical Statistics, Hayward, CA, 1991b.
44. R. C. Griffiths and P. Marjoram. Ancestral inference from samples of DNA sequences with recombination. *J. Comp. Biol.*, 3:479–502, 1996.
45. R. C. Griffiths and P. Marjoram. An ancestral recombination graph. In P. Donnelly and S. Tavaré, editors, *Progress in Population Genetics and Human Evolution*, pages 100–117. Springer Verlag, New York, 1997.
46. R. C. Griffiths and S. Tavaré. Ancestral inference in population genetics. *Statist. Sci.*, 9:307–319, 1994.
47. R. C. Griffiths and S. Tavaré. Simulating probability distributions in the coalescent. *Theor. Popn. Biol.*, 46:131–159, 1994.
48. R. C. Griffiths and S. Tavaré. Unrooted genealogical tree probabilities in the infinitely-many-sites model. *Math. Biosci.*, 127:77–98, 1995.
49. R. C. Griffiths and S. Tavaré. Monte Carlo inference methods in population genetics. *Mathl. Comput. Modelling*, 23:141–158, 1996.
50. R. C. Griffiths and S. Tavaré. Computational methods for the coalescent. In P. Donnelly and S. Tavaré, editors, *Progress in Population Genetics and Human Evolution*, pages 165–182. Springer Verlag, New York, 1997.
51. R. C. Griffiths and S. Tavaré. The age of a mutation in a general coalescent tree. *Stochastic Models*, 14:273–295, 1998.
52. R. C. Griffiths and S. Tavaré. The ages of mutations in gene trees. *Ann. Appl. Prob.*, 9:567–590, 1999.
53. R. C. Griffiths and S. Tavaré. The genealogy of a neutral mutation,. In P. J. Green, N. Hjørt, and S. Richardson, editors, *Highly Structured Stochastic Systems*. Oxford University Press,, 2003. in press.
54. S.-W. Guo. Linkage disequilibrium measures for fine-scale mapping: A comparison. *Hum. Hered.*, 47:301–314, 1997.
55. D. Gusfield. Efficient algorithms for inferring evolutionary trees. *Networks*, 21:19–28, 1991.
56. D. Gusfield. *Algorithms on Strings, Trees and Sequences. Computer Science and Computational Biology*. Cambridge University Press, Cambridge, 1997.

57. T. E. Harris. *The Theory of Branching Processes*. Springer Verlag, Berlin, 1963.
58. D. L. Hartl and E. W. Jones. *Genetics. Analysis of Genes and Genomes*. Jones and Bartlett, Sudbury, MA., Fifth edition, 2001.
59. W. K. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57:97–109, 1970.
60. J. Hey and J. Wakeley. A coalescent estimator of the population recombination fraction. *Genetics*, 145:833–846, 1997.
61. R. R. Hudson. Properties of a neutral allele model with intragenic recombination. *Theoret. Popn. Biol.*, 23:183–201, 1983.
62. R. R. Hudson. Estimating the recombination parameter of a finite population model without selection. *Genet. Res. Camb.*, 50:245–250, 1987.
63. R. R. Hudson. Gene genealogies and the coalescent process. In D. Futuyma and J. Antonovics, editors, *Oxford Surveys in Evolutionary Biology*, Volume 7, volume 7, pages 1–44. Oxford University Press, 1991.
64. R. R. Hudson. The how and why of generating gene genealogies. In N. Takahata and A. G. Clark, editors, *Mechanisms of Molecular Evolution*, pages 23–36. Sinauer, 1992.
65. R. R. Hudson. Linkage disequilibrium and recombination. In D. J. Balding, M. J. Bishop, and C. Cannings, editors, *Handbook of Statistical Genetics*, pages 309–324. John Wiley & Sons, Inc., Chichester, U.K., 2001.
66. R. R. Hudson and N. L. Kaplan. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*, 111:147–164, 1985.
67. N. L. Kaplan, T. Darden, and R. R. Hudson. The coalescent process in models with selection. *Genetics*, 120:819–829, 1988.
68. N. L. Kaplan and R. R. Hudson. The use of sample genealogies for studying a selectively neutral m -loci model with recombination. *Theoret. Popn. Biol.*, 28:382–396, 1985.
69. N. L. Kaplan, R. R. Hudson, and C. H. Langley. The “hitch-hiking” effect revisited. *Genetics*, 123:887–899, 1989.
70. S. Karlin and J. McGregor. Addendum to a paper of W. Ewens. *Theoret. Popn. Biol.*, 3:113–116, 1972.
71. S. Karlin and H. M. Taylor. *A Course in Stochastic Processes*, volume 2. Wiley, New York, NY, 1980.
72. M. Kimura and T. Ohta. The age of a neutral mutant persisting in a finite population. *Genetics*, 75:199–212, 1973.
73. J. F. C. Kingman. *Mathematics of Genetic Diversity*, volume 34 of *CBMS-NSF Regional Conference Series in Applied Mathematics*. Society for Industrial and Applied Mathematics, Philadelphia, Pennsylvania, 1980.
74. J. F. C. Kingman. The coalescent. *Stoch. Proc. Applns.*, 13:235–248, 1982.
75. J. F. C. Kingman. Exchangeability and the evolution of large populations. In G. Koch and F. Spizzichino, editors, *Exchangeability in Probability and Statistics*, pages 97–112. North-Holland Publishing Company, 1982.
76. J. F. C. Kingman. On the genealogy of large populations. *J. Appl. Prob.*, 19A:27–43, 1982.
77. J. F. C. Kingman. Origins of the coalescent: 1974–1982. *Genetics*, 156:1461–1463, 2000.
78. M. Kreitman. Methods to detect selection in populations with applications to the human. *Annu. Rev. Genomics Hum. Genet.*, 1:539–559, 2000.

79. S. M. Krone and C. Neuhauser. Ancestral processes with selection. *Theoret. Poul. Biol.*, 51:210–237, 1997.
80. M. Kuhner, J. Yamato, and J. Felsenstein. Estimating effective population size and mutation rate from sequence data using Metropolis-Hastings sampling. *Genetics*, 140:1421–1430, 1995.
81. M. Kuhner, J. Yamato, and J. Felsenstein. Maximum likelihood estimation of population growth rates based on the coalescent. *Genetics*, 149:429–434, 1998.
82. M. Kuhner, J. Yamato, and J. Felsenstein. Maximum likelihood estimation of recombination rates from population data. *Genetics*, 156:1393–1401, 2000.
83. B. Larget and D. L. Simon. Markov chain Monte Carlo algorithms for the Bayesian analysis of phylogenetic trees. *Mol. Biol. Evol.*, 16:750–759, 1999.
84. W.-H. Li and Y.-X. Fu. Coalescent theory and its applications in population genetics. In M. E. Halloran and S. Geisser, editors, *Statistics in Genetics*, pages 45–79. Springer Verlag, 1999.
85. J. S. Liu. *Monte Carlo Strategies in Scientific Computing*. Springer, New York, 2001.
86. R. Lundstrom. *Stochastic models and statistical methods for DNA sequence data*. PhD thesis, Mathematics Department, University of Utah, 1990.
87. R. S. Lundstrom, S. Tavaré, and R. H. Ward. Estimating mutation rates from molecular data using the coalescent. *Proc. Natl. Acad. Sci. USA*, 89:5961–5965, 1992.
88. R. S. Lundstrom, S. Tavaré, and R. H. Ward. Modeling the evolution of the human mitochondrial genome. *Math. Biosci.*, 112:319–335, 1992.
89. P. Marjoram, J. Molitor, V. Plagnol, and S. Tavaré. Markov chain Monte Carlo without likelihoods. *Proc. Natl. Acad. Sci. USA*, 000:000–000, 2003.
90. L. Markovtsova. *Markov chain Monte Carlo methods in population genetics*. PhD thesis, Mathematics Department, University of Southern California, 2000.
91. L. Markovtsova, P. Marjoram, and S. Tavaré. The effects of rate variation on ancestral inference in the coalescent. *Genetics*, 156:1427–1436, 2000b.
92. B. Mau, M. A. Newton, and B. Larget. Bayesian phylogenetic inference via Markov chain Monte Carlo methods. *Biometrics*, 55:1–12, 1999.
93. N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equations of state calculations by fast computing machines. *J. Chem. Phys.*, 21:1087–1092, 1953.
94. M. Möhle. Robustness results for the coalescent. *J. Appl. Prob.*, 35:438–447, 1998.
95. M. Möhle. Total variation distances and rates of convergence for ancestral coalescent processes in exchangeable population models. *Adv. Appl. Prob.*, 32:983–993, 2000.
96. M. Möhle. The coalescent in population models with time-inhomogeneous environment. *Stoch. Proc. and Applns.*, 97:199–227, 2002.
97. M. Möhle and S. Sagitov. A classification of coalescent processes for haploid exchangeable population models. *Ann. Probab.*, 29:1547–1562, 2001.
98. A. P. Morris, J. C. Whittaker, and D. J. Balding. Fine scale mapping of disease loci via shattered coalescent modelling of genealogies. *Amer. J. Hum. Genet.*, 70:686–707, 2002.
99. C. Neuhauser. Mathematical models in population genetics. In D. J. Balding, M. J. Bishop, and C. Cannings, editors, *Handbook of Statistical Genetics*, pages 153–177. John Wiley and Sons, Inc., New York, New York., 2001.

100. C. Neuhauser and S. M. Krone. The genealogy of samples in models with selection. *Genetics*, 145:519–534, 1997.
101. C. Neuhauser and S. Tavaré. The coalescent. In S. Brenner and J. Miller., editors, *Encyclopedia of Genetics*, Volume 1, pages 392–397. Academic Press, New York, 2001.
102. R. Nielsen. Estimation of population parameters and recombination rates from single nucleotide polymorphisms. *Genetics*, 154:931–942, 2000.
103. M. Nordborg. Coalescent theory. In D. J. Balding, M. J. Bishop, and C. Cannings, editors, *Handbook of Statistical Genetics*, pages 179–208. John Wiley and Sons, Inc., New York, New York., 2001.
104. N. Nordborg and S. Tavaré. Linkage disequilibrium: what history has to tell us. *Trends in Genetics*, 18:83–90, 2002.
105. N. Patil, A. J. Berno, D. A. Hinds, W. A. Barrett, J. M. Doshi, C. R. Hacker, C. R. Kautzer, D. H. Lee, C. Marjoribanks, D. P. McDonough, B. T. Nguyen, M. C. Norris, J. B. Sheehan, N. Shen, D. Stern, R. P. Stokowski, D. J. Thomas, M. O. Trulson, K. R. Vyas, K. A. Frazer, S. P. Fodor, and D. R. Cox. Blocks of limited haplotype diversity revealed by high-resolution scanning of human chromosome 21. *Science*, 294:1719–1723, 2001.
106. J. W. Pitman. Coalescents with multiple collisions. *Ann. Probab.*, 27:1870–1902, 1999.
107. V. Plagnol and S. Tavaré. Approximate Bayesian computation and MCMC. In H. Niederreiter, editor, *Monte Carlo and Quasi-Monte Carlo Methods 2002*. Springer-Verlag., 2004. In press.
108. A. Pluzhnikov. *Statistical inference in population genetics*. PhD thesis, Statistics Department, University of Chicago, 1997.
109. J. K. Pritchard and M. Przeworski. Linkage disequilibrium in humans: Models and data. *Amer. J. Hum. Genet.*, 69:1–14, 2001.
110. J. K. Pritchard, M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman. Population growth of human Y chromosomes: A study of Y chromosome microsatellites. *Mol. Biol. Evol.*, 16:1791–1798, 1999.
111. W. B. Provine. *The Origins of Theoretical Population Genetics*. University of Chicago Press, second edition, 2001.
112. B. D. Ripley. *Stochastic Simulation*. Wiley, New York, 1987.
113. S. Sagitov. The general coalescent with asynchronous mergers of ancestral lines. *J. Appl. Prob.*, 36:1116–1125, 1999.
114. I. W. Saunders, S. Tavaré, and G. A. Watterson. On the genealogy of nested subsamples from a haploid population. *Adv. Appl. Prob.*, 16:471–491, 1984.
115. S. Sawyer, D. Dykhuizen, and D. Hartl. Confidence interval for the number of selectively neutral amino acid polymorphisms. *Proc. Natl. Acad. Sci. USA*, 84:6225–6228, 1987.
116. J. Schweinsberg. Coalescents with simultaneous multiple collisions. *Electron. J. Prob.*, 5:1–50, 2000.
117. G. F. Shields, A. M. Schmeichen, B. L. Frazier, A. Redd, M. I. Vovoeda, J. K. Reed, and R. H. Ward. mtDNA sequences suggest a recent evolutionary divergence for Beringian and Northern North American populations. *Am. J. Hum. Genet.*, 53:549–562, 1993.
118. P. F. Slade. Most recent common ancestor probability distributions in gene genealogies under selection. *Theor. Popul. Biol.*, 58:291–305, 2000.
119. P. F. Slade. Simulation of ‘hitch-hiking’ genealogies. *J. Math. Biol.*, 42:41–70, 2001.

120. M. Slatkin. Simulating genealogies of selected alleles in a population of variable size. *Genet. Res.*, 78:49–57, 2001.
121. M. Slatkin and B. Rannala. Estimating the age of alleles by use of intraallelic variability. *Am. J. Hum. Genet.*, 60:447–458, 1997.
122. M. Slatkin and B. Rannala. Estimating allele age. *Annu. Rev. Genomics Hum. Genet.*, 1:225–249, 2000.
123. C. Soligo, O. Will, S. Tavaré, C. R. Marshall, and R. D. Martin. New light on the dates of primate origins and divergence. In M. J. Ravosa and M. Dagosto, editors, *Primate Origins and Adaptations*. Kluwer Academic/Plenum Publishers, New York, 2003.
124. J. Claiborne Stephens, Julie A. Schneider, Debra A. Tanguay, Julie Choi, Tara Acharya, Scott E. Stanley, Ruhong Jiang, Chad J. Messer, Anne Chew, Jin-Hua Han, Jicheng Duan, Janet L. Carr, Min Seob Lee, Beena Koshy, A. Madan Kumar, Ge Zhang, William R. Newell, Andreas Windemuth, Chuanbo Xu, Theodore S. Kalbfleisch, Sandra L. Shaner, Kevin Arnold, Vincent Schulz, Connie M. Drysdale, Krishnan Nandabalan, Richard S. Judson, Gualberto Ruanwo, and Gerald F. Vovis. Haplotype variation and linkage disequilibrium in 313 human genes. *Science*, 293:489–493, 2001.
125. M. Stephens. Times on trees and the age of an allele. *Theor. Popul. Biol.*, 57:109–119, 2000.
126. M. Stephens. Inference under the coalescent. In D. J. Balding, M. J. Bishop, and C. Cannings, editors, *Handbook of Statistical Genetics*, pages 213–238. John Wiley and Sons, Inc., New York, New York., 2001.
127. M. Stephens and P. Donnelly. Inference in molecular population genetics. *J. Roy. Statist. Soc. B*, 62:605–655, 2000.
128. F. M. Stewart. Variability in the amount of heterozygosity maintained by neutral mutations. *Theoret. Popul. Biol.*, 9:188–201, 1976.
129. C. Strobeck and K. Morgan. The effect of intragenic recombination on the number of alleles in a finite population. *Genetics*, 88:828–844, 1978.
130. F. Tajima. Evolutionary relationship of DNA sequences in finite populations. *Genetics*, 105:437–460, 1983.
131. H. Tang, D. O. Siegmund, P. Shen, P. J. Oefner, and M. W. Feldman. Frequentist estimation of coalescence times from nucleotide sequence data using a tree-based partition. *Genetics*, 161:447–459, 2002.
132. S. Tavaré. Line-of-descent and genealogical processes, and their applications in population genetics models. *Theoret. Popul. Biol.*, 26:119–164, 1984.
133. S. Tavaré. Calibrating the clock: using stochastic processes to measure the rate of evolution. In E. S. Lander and M. S. Waterman, editors, *Calculating the secrets of life*, pages 114–152. National Academy Press, Washington DC, 1993.
134. S. Tavaré. Ancestral inference from DNA sequence data. In H. G. Othmer, F. R. Adler, M. A. Lewis, and J. Dallon, editors, *Case Studies in Mathematical Modeling: Ecology, Physiology, and Cell Biology*, pages 81–96. Prentice-Hall, 1997.
135. S. Tavaré. Ancestral inference for branching processes. In P. Haccou and P. Jagers, editors, *Branching Processes in Biology: Variation, Growth, Extinction*. Cambridge University Press., 2004. In press.
136. S. Tavaré, D. J. Balding, R. C. Griffiths, and P. Donnelly. Inferring coalescence times for molecular sequence data. *Genetics*, 145:505–518, 1997.

137. S. Tavaré and W. J. Ewens. Multivariate Ewens distribution. In N. S. Johnson, S. Kotz, and N. Balakrishnan, editors, *Discrete Multivariate Distributions*, chapter 41, pages 232–246. Wiley, New York, 1997.
138. S. Tavaré, C. R. Marshall, O. Will, C. Soligo, and R. D. Martin. Using the fossil record to estimate the age of the last common ancestor of extant primates. *Nature*, 416:726–729, 2002.
139. J. L. Thorne, H. Kishino, and J. Felsenstein. Inching towards reality: an improved likelihood model of sequence evolution. *J. Mol. Evol.*, 34:3–16, 1992.
140. J. Tsao, Y. Yatabe, R. Salovaara, H. J. Järvinen, J. Mecklin, L. A. Altonen, S. Tavaré, and D. Shibata. Genetic reconstruction of individual colorectal tumor histories. *Proc. Natl. Acad. Sci. USA*, 97:1236–1241, 2000.
141. J. Wakeley. Using the variance of pairwise differences to estimate the recombination rate. *Genet. Res. Camb.*, 69:45–58, 1997.
142. J. D. Wall. A comparison of estimators of the population recombination rate. *Mol. Biol. Evol.*, 17:156–163, 2000.
143. R. H. Ward, B. L. Frazier, K. Dew, and S. Pääbo. Extensive mitochondrial diversity within a single Amerindian tribe. *Proc. Natl. Acad. Sci. USA*, 88:8720–8724, 1991.
144. M. S. Waterman. *Introduction to Computational Biology. Maps, Sequences and Genomes*. Chapman & Hall, London, 1995.
145. G. A. Watterson. Models for the logarithmic species abundance distributions. *Theoret. Popn. Biol.*, 6:217–250, 1974.
146. G. A. Watterson. On the number of segregating sites in genetical models without recombination. *Theoret. Popn. Biol.*, 7:256–276, 1975.
147. G. A. Watterson. Heterosis or neutrality? *Genetics*, 85:789–814, 1977.
148. G. A. Watterson. Reversibility and the age of an allele II. Two-allele models, with selection and mutation. *Theoret. Popul. Biol.*, 12:179–196, 1977.
149. G. A. Watterson. The homozygosity test of neutrality. *Genetics*, 88:405–417, 1978.
150. G. A. Watterson. Motoo Kimura’s use of diffusion theory in population genetics. *Theoret. Popul. Biol.*, 49:154–188, 1996.
151. G. Weiss and A. von Haeseler. Estimating the age of the common ancestor of men from the ZFY intron. *Science*, 257:1359–1360, 1996.
152. G. Weiss and A. von Haeseler. A coalescent approach to the polymerase chain reaction. *Nucleic Acids Research*, 25:3082–3087, 1997.
153. G. Weiss and A. von Haeseler. Inference of population history using a likelihood approach. *Genetics*, 149:1539–1546, 1998.
154. K. M. Weiss and A. G. Clark. Linkage disequilibrium and the mapping of complex human traits. *TIG*, 18:19–24, 2002.
155. L. S. Whitfield, J. E. Sulston, and P. N. Goodfellow. Sequence variation of the human Y chromosome. *Nature*, 378:379–380, 1995.
156. O. Will, V. Plagnol, C. Soligo, R. D. Martin, and S. Tavaré. Statistical inference in the primate fossil record: a Bayesian approach. In preparation, 2003.
157. I. J. Wilson and D. J. Balding. Genealogical inference from microsatellite data. *Genetics*, 150:499–510, 1998.
158. C. Wiuf and P. Donnelly. Conditional genealogies and the age of a neutral mutant. *Theor. Popul. Biol.*, 56:183–201, 1999.
159. S. Wright. Evolution in Mendelian populations. *Genetics*, 16:97–159, 1931.

160. S. Wright. *Evolution and the Genetics of Populations.*, volume 4, Variability within and among natural populations. University of Chicago Press, Chicago, 1978.
161. Z. Yang and B. Rannala. Bayesian phylogenetic inference using DNA sequences: A Markov chain Monte Carlo method. *Mol. Biol. Evol.*, 14:717–724, 1997.

**Ofer Zeitouni: Random Walks in Random
Environment**

Random Walks in Random Environment

Ofer Zeitouni

Department of Electrical Engineering
Department of Mathematics
Technion – Israel Institute of Technology
Haifa 32000, Israel

1	Introduction	193
1.1	Model	193
1.2	Examples	194
2	RWRE – $d=1$	195
2.1	Ergodic theorems	195
2.2	CLT for ergodic environments	209
2.3	Large deviations	213
2.4	The subexponential regime	236
2.5	Sinai’s model: non standard limit laws and aging properties	248
3	RWRE – $d > 1$	258
3.1	Ergodic Theorems	258
3.2	A Law of Large Numbers in \mathbb{Z}^d	261
3.3	CLT for walks in balanced environments	269
3.4	Large deviations for nestling walks	283
3.5	Kalikow’s condition	293
	References	308

Preface

These notes on random walks in random environments (RWRE) reflect what I hoped to cover in the 15 hours of the St Flour course on this topic, July 9–25, 2001. Of course, this turned out to be over optimistic. Departing even further from the actually delivered lectures, I have taken advantage of the year that elapsed to add some material (especially, related to multi-dimensional walks) and to correct numerous mistakes and omissions.

The manuscript consist roughly of two parts: the first deals with RWRE on \mathbb{Z} . The interest in the model began in the early 70's, and with the detailed analysis of RWRE asymptotics in the last decade, has now reached maturity (for an account of the history of the subject and many of the results through the early 90's, see [37]). I have tried to present different tools for the study of such walks, risking some repetition of results in a few cases, and deferring to the bibliographical notes a discussion of refinements and sharpening of the results. It is worthwhile to point out that RWRE's on \mathbb{Z} have already been considered in previous St Flour courses (most notably by Ledrappier [50] and by Molchanov [53]), but the emphasis in this presentation is quite different.

The second part of the notes deals with \mathbb{Z}^d . This is currently an active research area, and one hopes that much progress will be made in the next few years. My goal here was to expose the audience to some tools which have proved useful, and to point out several directions where further progress could be made. In several places, I have tried to lay the groundwork for relaxing the often made assumption of i.i.d. environment.

When preparing the notes, and taking into account the time frame of these lectures, it became clear that there were topics that had to be left out. Even the uninitiated will quickly realize that the most glaring omission is the study of RWRE's by renormalization techniques. There are three reasons for this: first, it would take too long to properly expose it. Second, these methods have not yet reached the full scope of their applicability, and in view of very active current research efforts in this direction, any account written now risks being outdated very quickly. And third, an overview of the current status of these techniques can be found in [69] and [70]. Time constraints also did not allow me to discuss random walks on Galton-Watson trees, a topic that has seen much progress in recent years.

Parts of the material presented here is based on joint work, some still unpublished, with F. Comets, A. Dembo, N. Gantert, and Y. Peres. I thank them all, both for the many hours spent together on thinking about RWRE, and for their generosity. I would also like to thank my colleagues in Haifa who suffered through a first draft of these notes in the winter of 2000. In particular, comments from D. Ioffe, H. Kaspi, E. Mayer-Wolf, A. Roitershtein and M. Zerner are gratefully acknowledged. Similarly acknowledged are useful remarks from D. Cheliotis, A. Dembo, N. Gantert, A. Guionnet, H. Kesten, D. Piau, and S.R.S. Varadhan. Comments from participants at the St Flour summer school helped improve the presentation and strengthen numerous re-

sults. I am particularly grateful to P. Bougerol who allowed me to incorporate some of his suggestions in the final text of these notes, and D. Ocone and F. Rassoul-Agha for stimulating discussions. Last but not least, I am grateful to J. Picard for the smooth and gentle running of the summer school.

A typographical comment: for aesthetic reasons, I consistently use $P_\omega^o, E_\omega^o, \mathbb{P}^o, \mathbb{E}^o$, etc., when I mean $P_\omega^0, E_\omega^0, \mathbb{P}^0, \mathbb{E}^0$.

1 Introduction

The definition of a RWRE involves two components: first, the *environment*, which is randomly chosen but kept fixed throughout the time evolution, and second, the random walk, which, given the environment, is a time homogeneous Markov chain whose transition probabilities depend on the environment. We do not attempt here a historical review of RWRE's, or in greater generality of motion in homogeneous media, except for stating that we insist on the environment being static, i.e. time independent, and that in general the random walk (conditioned on the environment) is not necessarily reversible.

1.1 Model

We begin with a general setup, that will be specialized later to the cases of interest to us. Let (V, E) denote an (infinite, oriented) graph with countable vertex set V and edges set $E = \{(v, w)\}$ (we allow, but do not require, $(v, v) \in E$). For each $v \in V$, we define its *neighborhood* N_v by

$$N_v = \{w \in V : (v, w) \in E\},$$

throughout assuming that $|N_v| < \infty$, for all $v \in V$.

For each $v \in V$, let $M_1(N_v)$ denote the collection of probability measures on V with support N_v . Formally, an element of $M_1(N_v)$, called a *transition law* at v , is a measurable function $\omega_v : V \rightarrow [0, 1]$ satisfying:

$$\begin{aligned} \text{(a)} \quad & \omega_v(w) \geq 0 \quad \forall w \in V \\ \text{(b)} \quad & \omega_v(w) = 0 \quad \forall w \notin N_v \\ \text{(c)} \quad & \sum_{w \in N_v} \omega_v(w) = 1 \end{aligned} \tag{1.1.1}$$

Note that if $v \in N_v$ then in (1.1.1c) we allow for $\omega_v(v) > 0$.

We equip $M_1(N_v)$ with the weak topology on probability measures, which makes it into a Polish space. Further, it induces a Polish structure on $\Omega = \prod_{v \in V} M_1(N_v)$. We let \mathcal{F} denote the Borel σ -algebra on Ω (which is the same as the σ -algebra generated by cylinder functions). Given a probability measure P

on (Ω, \mathcal{F}) , a *random environment* is an element ω of Ω distributed according to P .

We turn next to define the class of random walks of interest to us. For each $\omega \in \Omega$, we define the *random walk in the environment* ω as the time-homogeneous Markov chain $\{X_n\}$ taking values in V with transition probabilities

$$P_\omega(X_{n+1} = w | X_n = v) = \omega_v(w).$$

We use P_ω^v to denote the law induced on $(V^{\mathbb{N}}, \mathcal{G})$ where \mathcal{G} is the σ -algebra generated by cylinder functions and

$$P_\omega^v(X_0 = v) = 1.$$

In the sequel, we refer to $P_\omega^v(\cdot)$ as the *quenched* law of the random walk $\{X_n\}$. Note that for each $G \in \mathcal{G}$, the map

$$\omega \mapsto P_\omega^v(G)$$

is \mathcal{F} -measurable. Hence, we may define the measure $\mathbb{P}^v := P \otimes P_\omega^v$ on $(\Omega \times V^{\mathbb{N}}, \mathcal{F} \times \mathcal{G})$ from the relation

$$\mathbb{P}^v(F \times G) = \int_F P_\omega^v(G) P(d\omega), \quad F \in \mathcal{F}, G \in \mathcal{G}. \quad (1.1.2)$$

The marginal of \mathbb{P}^v on $V^{\mathbb{N}}$, denoted also \mathbb{P}^v whenever no confusion occurs, is called the *annealed law* of the random walk $\{X_n\}$; note that under \mathbb{P}^v , the random walk in random environment (RWRE) $\{X_n\}$ is *not* a Markov chain!

1.2 Examples

Throughout these notes, we only treat nearest neighbor RWRE's on \mathbb{Z}^d :

Nearest neighbor RWRE on \mathbb{Z}

Here, we take $V = \mathbb{Z}$ and $E = \cup_{z \in \mathbb{Z}} \{(z, z+1), (z, z)\}$. Then, $N_v = \{v-1, v, v+1\}$ and $M_1(N_v)$ can be identified with the three dimensional simplex; We let $\omega_z^+ := \omega_z(z+1)$, $\omega_z^- := \omega_z(z-1)$, and $\omega_z^0 := \omega_z(z)$. One defines naturally the shift θ on Ω by $(\theta\omega)_z = \omega_{z+1}$. We always make the following assumption: $(\omega, \mathcal{F}, P, \theta)$ is an ergodic system.

It is worthwhile commenting, already at this stage, that for each ω there exists a reversing measure that makes the RWRE reversible. More details are provided in Section 2.1.

Nearest neighbor RWRE on \mathbb{Z}^d

Here, $V = \mathbb{Z}^d$ and $E = \cup_{z \in \mathbb{Z}^d} \{ \cup_{y \sim z} (z, y) \cup (z, z) \}$. For each $v \in V$, N_v contains $2d + 1$ vertices, and $M_1(N_v)$ is identified with the $2d + 1$ -dimensional simplex. One may define the family of shifts $\{ \theta^e \}_{|e|_1=1}$. As in the case of $d = 1$, we always require P to be ergodic with respect to this family. We write throughout $\omega(x, e) := \omega_x(x + e)$. Unlike the case with $d = 1$, the Markov chain defined by P_ω^v is, in general, not reversible.

Bibliographical notes: the preface section contains relevant bibliography on the RWRE model in \mathbb{Z}^d , $d \geq 1$. We mention here some other models of random walks in random media that can be adapted into the general framework presented above, but that will not be considered in these notes:

- *Non nearest neighbor walks: For \mathbb{Z}^1 , see the recent thesis [7], that includes also a summary of earlier work and in particular of [43]. I am not aware of a systematic study of non nearest neighbor RWRE's on \mathbb{Z}^d , see however [79] for some results valid in that generality.*
- *Reversible random walks in random environments in \mathbb{Z}^d , $d > 1$: the prime example is the random conductance model, in which bonds on \mathbb{Z}^d carry i.i.d. conductances and modulate the transition mechanism of the walk, see [14]. Other models in the same spirit, and their surprising behavior, are described in [6] and the references therein.*
- *Random walks on Galton-Watson trees: see [15, 51, 52, 59] for recent developments.*

2 RWRE – d=1

This chapter is devoted to the study of the one-dimensional model, where sharp results are available. As a warm-up to the high dimensional case, we sometimes present different proofs of the same statement.

Our exposition progresses from ergodic properties and law of large numbers (Section 2.1), to the study of central limit theorems (Section 2.2), large deviations (Section 2.3), subexponential tail estimates (2.4), and subdiffusive behaviour and aging (Section 2.5). Each section contains a (non-exhaustive!) list pointing to the literature.

2.1 Ergodic theorems

In this section, we are interested in questions concerning transience, recurrence and laws of large numbers, in the most general nearest neighbour one-dimensional setup. Define $\rho_z = \omega_z^- / \omega_z^+$.

Assumption 2.1.1

(A1) P is stationary and ergodic.

(Note that the solution to (2.1.3) is unique due to the maximum principle, hence it is enough to verify that the function in (2.1.4) satisfies (2.1.3).)

Define $S(\omega) = \sum_{n=1}^{\infty} \rho_1 \cdots \rho_n$, $F(\omega) = \sum_{n=0}^{\infty} \rho_0^{-1} \cdots \rho_{-n}^{-1}$. Further, define the events

$$\mathcal{S}_+ = \{S(\omega) < \infty\}, \quad \mathcal{F}_+ = \{F(\omega) < \infty\}.$$

Then:

- on $T_+ := \{\mathcal{S}_+ \cap \mathcal{F}_+^c\}$, it holds that

$$\lim_{m \rightarrow \infty} [1 - \mathcal{V}_{1,m,\omega}(0)] > 0, \quad \lim_{k \rightarrow \infty} \lim_{m \rightarrow \infty} [1 - \mathcal{V}_{k,m,\omega}(0)] = 1.$$

Hence, for $\omega \in T_+$,

$$P_\omega^o(\lim_{n \rightarrow \infty} X_n = \infty) = 1.$$

- Similarly, for $\omega \in T_- := \{\mathcal{S}_+^c \cap \mathcal{F}_+\}$,

$$P_\omega^o(\lim_{n \rightarrow \infty} X_n = -\infty) = 1.$$

- Finally, if $\omega \in R := \{\mathcal{S}_+^c \cap \mathcal{F}_+^c\}$ then, for any fixed k ,

$$1 - \lim_{m \rightarrow \infty} \mathcal{V}_{k,m,\omega}(0) = \lim_{m \rightarrow \infty} \mathcal{V}_{m,k,\omega}(0) = 0,$$

and hence, for $\omega \in R$,

$$P_\omega^o(-\infty = \liminf_{n \rightarrow \infty} X_n < \limsup_{n \rightarrow \infty} X_n = \infty) = 1.$$

We observe next that both \mathcal{S}_+ and \mathcal{F}_+ are invariant events, hence $P(\mathcal{S}_+) \in \{0, 1\}$, $P(\mathcal{F}_+) \in \{0, 1\}$ by the ergodicity of P . Next, $P(\mathcal{S}_+) = 1 \Rightarrow P(\mathcal{F}_+) = 0$ by the shift-invariance of P . Thus, it is enough to prove that $P(\mathcal{S}_+) = 1$ if and only if $E_P(\log \rho_0) < 0$ and $P(\mathcal{F}_+) = 1$ if and only if $E_P(\log \rho_0) > 0$. We prove the first claim only, the second one possessing a similar proof.

Assume first $c := E_P(\log \rho_0) < 0$. Then, by the ergodic theorem, there exists an $n_0(\omega)$ with $P(n_0(\omega) < \infty) = 1$ such that $\frac{1}{n} \sum_{i=1}^n \log \rho_i \leq c/2 < 0$ for all $n > n_0(\omega)$. But then, for some $C_1(\omega) < \infty$ P -a.s.,

$$\sum_{n=1}^{\infty} \rho_1 \cdots \rho_n \leq C_1(\omega) + \sum_{k=n_0(\omega)+1}^{\infty} e^{kc/2} < \infty, \quad P\text{-a.s.}$$

implying $P(\mathcal{S}_+) = 1$. Conversely, for $\omega \in \mathcal{S}_+$, $\lim_{n \rightarrow \infty} \sum_{k=1}^n \log \rho_k = -\infty$. But, $\{Y_i = -\log \rho_i\}$ are stationary, and the claim follows from the following well known:

Lemma 2.1.5 (Kesten[40]) *For any real valued, stationary sequence $\{Y_i\}$, fix $Z_n = \sum_{i=1}^n Y_i$. Then, one has with probability 1 that the event $\{Z_n \rightarrow_{n \rightarrow \infty} \infty\}$ implies $\{\liminf_{n \rightarrow \infty} Z_n/n > 0\}$.*

Indeed, Kesten’s lemma implies that on \mathcal{S}_+ ,

$$E_P(\log \rho_0) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \log \rho_k < 0, \quad P - \text{a.s.}, \tag{2.1.6}$$

and $P(\mathcal{S}_+) = 1$ thus implies $E_P(\log \rho_0) < 0$ and completes the proof of Theorem 2.1.2. \square

Remarks: 1. P. Bougerol has kindly indicated to me the followig proof of the implication $P(\mathcal{S}_+) = 1 \Rightarrow E_P(\log \rho_0) < 0$, which bypasses the use of Kesten’s lemma: define the function $f(\omega) = \log S(\omega)$. $P(\mathcal{S}_+) = 1$ implies that $f(\omega)$ is well defined. Since $S(\omega) = \rho_1 + \rho_1 S(\theta\omega)$, it holds that $f(\omega) > \log \rho_1 + f(\theta\omega)$, and we conclude by **(A2)** that $(f(\theta\omega) - f(\omega))_+$ is P -integrable and hence $E_P[f(\omega) - f(\theta\omega)] = 0$ (this is Mañe’s lemma, apply the ergodic theorem to see it!). Using again $f(\omega) > \log \rho_1 + f(\theta\omega)$, one concludes that $0 > E_P \log \rho_1 = E_P \log \rho_0$, as claimed.

2. If P is i.i.d., Theorem 2.1.2 remains valid when the left hand side of conditions (a), (b), (c) is replaced, respectively, by

$$\begin{aligned} \text{(a')} : \quad & \sum_{n=1}^{\infty} n^{-1} P \left(\prod_{j=1}^n \rho_j > 1 \right) < \infty. \\ \text{(b')} : \quad & \sum_{n=1}^{\infty} n^{-1} P \left(\prod_{j=1}^n \rho_j < 1 \right) < \infty. \\ \text{(c')} : \quad & \sum_{n=1}^{\infty} n^{-1} P \left(\prod_{j=1}^n \rho_j < 1 \right) = \sum_{n=1}^{\infty} n^{-1} P \left(\prod_{j=1}^n \rho_j > 1 \right) = \infty. \end{aligned}$$

This is useful in particular when $E_P(\log \rho_0)$ is not well defined. See [67] for details.

Having developed transience and recurrence criteria, we turn to the law of large numbers. We first note that one cannot apply directly ergodic theorems to the sequence X_n/n : The sequence $\{X_n - X_{n-1}\}$ is not even stationary! We will exhibit two approaches to the LLN: The first is based on a hitting times decomposition. The second approach is based on the point of view of the “environment viewed from the particle”.

LLN-version I: hitting time decompositions

Introduce the following notations:

$$\bar{S} = \sum_{i=1}^{\infty} \frac{1}{\omega_{(-i)}^+} \prod_{j=0}^{i-1} \rho_{(-j)} + \frac{1}{\omega_0^+} \tag{2.1.7}$$

$$\bar{F} = \sum_{i=1}^{\infty} \frac{1}{\omega_i^-} \prod_{j=0}^{i-1} \rho_j^{-1} + \frac{1}{\omega_0^-} \tag{2.1.8}$$

Theorem 2.1.9 *Assume Assumption 2.1.1. Then,*

- (a) $E_P(\overline{S}) < \infty \Rightarrow \lim_{n \rightarrow \infty} \frac{X_n}{n} = \frac{1}{E_P(\overline{S})}, \quad \mathbb{P}^o \text{ a.s.}$
- (b) $E_P(\overline{F}) < \infty \Rightarrow \lim_{n \rightarrow \infty} \frac{X_n}{n} = -\frac{1}{E_P(\overline{F})}, \quad \mathbb{P}^o \text{ a.s.}$
- (c) $E_P(\overline{S}) = \infty \text{ and } E_P(\overline{F}) = \infty \Rightarrow \lim_{n \rightarrow \infty} \frac{X_n}{n} = 0, \quad \mathbb{P}^o \text{ a.s.}$

Remark: In the case that P is i.i.d., (a)–(c) of Theorem 2.1.9 become

- (a') $E_P(\rho_0) < 1 \Rightarrow \lim_{n \rightarrow \infty} \frac{X_n}{n} = \frac{1 - E_P(\rho_0)}{E_P\left(\frac{1}{\omega_0^+}\right)}, \quad \mathbb{P}^o \text{ a.s.}$
- (b') $E_P(\rho_0^{-1}) < 1 \Rightarrow \lim_{n \rightarrow \infty} \frac{X_n}{n} = -\frac{1 - E_P\left(\frac{1}{\rho_0}\right)}{E_P\left(\frac{1}{\omega_0^-}\right)}, \quad \mathbb{P}^o \text{ a.s.}$
- (c') $\frac{1}{E_P(\rho_0)} \leq 1 \leq E_P(\rho_0^{-1}) \Rightarrow \lim_{n \rightarrow \infty} \frac{X_n}{n} = 0, \quad \mathbb{P}^o \text{ a.s.}$

since $E_P \log \rho_0 \leq \log E_P \rho_0$ with a strict inequality whenever P is non-degenerate, it follows that one can find examples where $X_n \rightarrow \infty$ \mathbb{P}^o -a.s. but $X_n/n \rightarrow 0$, \mathbb{P}^o -a.s. This does not contradict Kesten’s lemma (Lemma 2.1.5) because $\{X_n - X_{n-1}\}$ is not in general a stationary sequence under \mathbb{P}^o .

Proof of Theorem 2.1.9

We introduce hitting times which will serve us later too. Let $T_0 = 0$, and

$$T_n = \min\{k : X_k = n\}$$

with the usual convention that the minimum over an empty set is $+\infty$. Set $\tau_0 = 0$ and

$$\tau_n = T_n - T_{n-1}, \quad n \geq 1.$$

Similarly, set

$$T_{-n} = \min\{k : X_k = -n\}$$

and

$$\tau_{-n} = T_{-n} - T_{-n+1}, \quad n \geq 1,$$

the convention being that $\tau_{\pm n} = \infty$ if $T_{\pm n} = \infty$. We have the following lemma:

Lemma 2.1.10 *If $\limsup_{n \rightarrow \infty} X_n = +\infty$, \mathbb{P}^o -a.s., then $\{\tau_i\}_{i \geq 1}$ is a stationary and ergodic sequence. If further P is strongly mixing, then $\{\tau_i\}_{i \geq 1}$ is also strongly mixing.*

Proof of Lemma 2.1.10

The stationarity of $\{\tau_i\}_{i \geq 1}$ follows from the stationarity of the environment. To see the ergodicity, let $\Xi = [0, 1]^{\mathbb{N}}$, let U_{Ξ} denote the measure on Ξ making all coordinates $\{\xi_i\}$ independent and of uniform law on $[0, 1]$, and note that $\{X_n\}$ may be constructed by writing

$$X_{n+1} = X_n + \mathbf{1}_{\{\omega_{X_n}^+ < \xi_{n+1}\}} - \mathbf{1}_{\{\xi_{n+1} \in [\omega_{X_n}^+, \omega_{X_n}^+ + \omega_{X_n}^-]\}}.$$

Suppose $A = A(\omega, \xi) = A(\tau)$ is an event, measurable w.r.t. $\mathcal{G}_n = \sigma\{\tau_i, i \geq 1\}$, which is invariant with respect to the shift $(\theta\tau)_i = \tau_{i+1}$ (we write in the sequel $\theta A = A(\theta\tau)$). We need only show that $P \otimes U_{\Xi}(A) \in \{0, 1\}$. Note however that $\theta^k A$, conditioned on $\sigma\{\omega_i, i \in \mathbb{Z}\}$, is independent of ξ_1, \dots, ξ_k . Thus, since U_{Ξ} is an i.i.d. law and hence the tail sigma-field of $\{\xi_i\}$ is trivial, it follows that $A = \theta^k A$ is, under the above conditioning, independent of $\sigma\{\xi_i, i \geq 1\}$. Thus A depends only on ω . But the shift θ on the sequence $\{\tau_i\}$ induces the usual shift θ on Ω . Thus, $\theta A = A(\theta\omega) = A(\omega)$ and hence $P(A) \in \{0, 1\}$.

To prove the strong mixing properties (which we do not actually need in the sequel), consider sets $A_1 \cdots A_k, B_1 \cdots B_j \subset \mathbb{Z}$, and let

$$\mathcal{A} = \bigcap_{i=1}^k \{\tau_i \in A_i\}, \quad \mathcal{B}^m = \bigcap_{i=1}^j \{\tau_{m+i} \in B_i\}.$$

Clearly, $\mathbb{P}^o(\mathcal{B}^m) = \mathbb{P}^o(\mathcal{B}^0)$, and thus we need to prove that whenever $\limsup_{n \rightarrow \infty} X_n = \infty$ \mathbb{P}^o -a.s., then

$$\lim_{m \rightarrow \infty} \mathbb{P}^o(\mathcal{A} \cap \mathcal{B}^m) = \mathbb{P}^o(\mathcal{A})\mathbb{P}^o(\mathcal{B}^0).$$

Toward this end, let

$$B_i^K = B_i \cap [0, K] \quad \text{and} \quad \mathcal{B}^{m,K} = \bigcap_{i=1}^j \{\tau_{m+i} \in B_i^K\}.$$

Fix $\varepsilon > 0$ and then $K = K(\varepsilon)$ large enough such that

$$\mathbb{P}^o(\mathcal{B}^m \setminus \mathcal{B}^{m,K}) = \mathbb{P}^o(\mathcal{B}^0 \setminus \mathcal{B}^{0,K}) \leq \varepsilon$$

which is possible since $\limsup_{n \rightarrow \infty} X_n = \infty$, \mathbb{P}^o -a.s. Note that, for $m > k$,

$$P_{\omega}^o(\mathcal{A} \cap \mathcal{B}^{K,m}) = P_{\omega}^o(\mathcal{A})P_{\omega}^o(\mathcal{B}^{K,m})$$

and that $P_{\omega}^o(\mathcal{A})$ is measurable with respect to $\sigma(\omega_i, i \leq k-1)$. On the other hand, since $|X_{n+1} - X_n| = 1$, on the event $\{\tau_{m+i} \leq K, 1 \leq i \leq j\}$, it holds that $X_n \geq m - K$ for $T_m \leq n \leq T_{m+j+1}$. Thus, for $m > K + k$, $P_{\omega}^o(\mathcal{B}^{K,m})$ is measurable with respect to $\sigma(\omega_i, i \geq m - K)$. It follows from the strong mixing of P that

$$\begin{aligned} \lim_{m \rightarrow \infty} \mathbb{P}^o(\mathcal{A} \cap \mathcal{B}^{K,m}) &= \lim_{m \rightarrow \infty} E_P(P_\omega^o(\mathcal{A})P_\omega^o(\mathcal{B}^{K,m})) \\ &= E_P(P_\omega^o(\mathcal{A})) \cdot E_P(P_\omega^o(\mathcal{B}^{K,0})) \\ &= \mathbb{P}^o(\mathcal{A})\mathbb{P}^o(\mathcal{B}^{K,0}). \end{aligned} \tag{2.1.11}$$

On the other hand,

$$\mathbb{P}^o(\mathcal{A} \cap \mathcal{B}^m) - \varepsilon \leq \mathbb{P}^o(\mathcal{A} \cap \mathcal{B}^{K,m}) \leq \mathbb{P}^o(\mathcal{A} \cap \mathcal{B}^m)$$

while

$$\mathbb{P}^o(\mathcal{B}^m) - \varepsilon \leq \mathbb{P}^o(\mathcal{B}^{K,m}) \leq \mathbb{P}^o(\mathcal{B}^m)$$

and one concludes from (2.1.11) that

$$\left| \lim_{m \rightarrow \infty} \mathbb{P}^o(\mathcal{A} \cap \mathcal{B}^m) - \mathbb{P}^o(\mathcal{A})\mathbb{P}^o(\mathcal{B}^0) \right| \leq \varepsilon$$

which implies the claim since ε is arbitrary. □

Remark: Note that an attempt to mimick this argument in \mathbb{Z}^d , $d \geq 1$, with T_i denoting the hitting times of hyperplanes at distance i from the origin, fails because of the extra information contained in the hitting location.

Our strategy consists now of applying the ergodic theorem to the sequence $\{\tau_i\}$. As a first step, we have the

Lemma 2.1.12 *Assume Assumption 2.1.1. Then,*

- (a) $E_{\mathbb{P}^o}(\tau_1) = E_P(\overline{S})$,
- (b) $E_{\mathbb{P}^o}(\tau_{-1}) = E_P(\overline{F})$.

Proof. We prove only (a), the proof of (b) being similar. Decompose, with $X_0 = 0$,

$$\tau_1 = \mathbf{1}_{X_1=1} + \mathbf{1}_{X_1=0}(1 + \tau'_1) + \mathbf{1}_{X_1=-1}(1 + \tau''_0 + \tau''_1). \tag{2.1.13}$$

Here, (τ'_1) is the first hitting time of 1 after time 1 (possibly infinite), $(1 + \tau''_0)$ is the first hitting time of 0 after time 1, and $1 + \tau''_0 + \tau''_1$ is the first hitting time of 1 after time $1 + \tau''_0$.

Under P_ω^o , the law of τ'_1 conditioned on the event $\{X_1 = 0\}$ is identical to the law of τ_1 , the law of τ''_0 conditioned on the event $\{X_1 = -1\}$ is $P_{\theta_{-1}\omega}^o(\tau_1 \in \cdot)$, while conditioned on the event $\{X_1 = -1\} \cap \{\tau''_0 < \infty\}$, τ''_1 also has law identical to that of τ_1 .

Consider first the case $E_{\mathbb{P}^o}(\tau_1) < \infty$. Then, both $E_\omega^o(\tau_1) < \infty$ and $E_{\theta_{-1}\omega}^o(\tau_1) < \infty$, P -a.s. Taking expectations in (2.1.13), one gets then

$$E_\omega^o(\tau_1) = 1 + (1 - \omega_0^+)E_\omega^o(\tau_1) + \omega_0^- E_{\theta_{-1}\omega}^o(\tau_1).$$

Hence,

$$E_\omega^o(\tau_1) = \frac{1}{\omega_0^+} + \rho_0 E_{\theta_{-1}\omega}^o(\tau_1).$$

Iterating this equation, we get

$$E_\omega^o(\tau_1) = \frac{1}{\omega_0^+} + \frac{\rho_0}{\omega_{(-1)}^+} + \frac{\rho_0\rho_{(-1)}}{\omega_{(-2)}^+} + \dots + \frac{\prod_{i=0}^{-(m-1)} \rho_{(-i)}}{\omega_{(-m)}^+} + \left(\prod_{i=0}^{-m+1} \rho_{(-i)} \right) E_{\theta^{-m}\omega}^o(\tau_1). \quad (2.1.14)$$

Omitting the last term, taking expectations on both sides, and then taking $m \rightarrow \infty$ using dominated convergence, we get

$$E_{\mathbb{P}^o}(\tau_1) \geq E_P(\bar{S}). \quad (2.1.15)$$

To see the reverse inequality, note that by (2.1.13), for any $M < \infty$,

$$E_\omega^o(\tau_1 \mathbf{1}_{\tau_1 < M}) \leq 1 + (1 - \omega_0^+)E_\omega^o(\tau_1 \mathbf{1}_{\tau_1 < M}) + \omega_0^- E_{\theta^{-1}\omega}^o(\tau_1 \mathbf{1}_{\tau_1 < M}).$$

Iterating, we get that

$$E_\omega^o(\tau_1 \mathbf{1}_{\tau_1 < M}) \leq \bar{S} + M \prod_{i=0}^{-m+1} \rho_{(-i)}.$$

Taking expectations, we get that

$$E_{\mathbb{P}^o}(\tau_1 \mathbf{1}_{\tau_1 < M}) \leq E_P(\bar{S}) + M E_P \left(\prod_{i=0}^{-m+1} \rho_{(-i)} \right).$$

Assuming $E_P(\bar{S}) < \infty$ and hence $E_P \left(\prod_{i=0}^{-m+1} \rho_{(-i)} \right) \xrightarrow{m \rightarrow \infty} 0$, we get that

$$E_{\mathbb{P}^o}(\tau_1 \mathbf{1}_{\tau_1 < M}) \leq E_P(\bar{S}).$$

Taking $M \rightarrow \infty$ and using monotone convergence we conclude, using also (2.1.15), that

$$E_{\mathbb{P}^o}(\tau_1 \mathbf{1}_{\tau_1 < \infty}) = E_P(\bar{S}),$$

completing the proof that $E_{\mathbb{P}^o}(\tau_1) < \infty \Rightarrow E_{\mathbb{P}^o}(\tau_1) = E_P(\bar{S})$.

It thus remains to show that $E_{\mathbb{P}^o}(\tau_1) = \infty \Rightarrow E_P(\bar{S}) = \infty$. Note next that if $E_P(\log \rho_0) \leq 0$, we have by Theorem 2.1.2 that

$$E_{\mathbb{P}^o}(\tau_1 \mathbf{1}_{\tau_1 < \infty}) = E_{\mathbb{P}^o}(\tau_1)$$

hence $E_{\mathbb{P}^o}(\tau_1) = \infty$ implies $E_P(\bar{S}) = \infty$. On the other hand, if $E_P(\log \rho_0) > 0$ then $\prod_{i=1}^0 \rho_{(-j)} \rightarrow_{j \rightarrow \infty} \infty$, P -a.s. by the ergodic theorem and hence also $E_P(\bar{S}) = \infty$. This concludes the proof of Lemma 2.1.12. \square

Remark: In fact, a similar proof shows that in the uniformly elliptic case, $E_\omega^o(\tau_1) = \bar{S}$, for every environment ω .

An application of Lemmas 2.1.10 and 2.1.12 yields that in case (a)

$$\frac{T_n}{n} = \frac{\sum_{i=1}^n \tau_i}{n} \xrightarrow{n \rightarrow \infty} E_{\mathbb{P}^o}(\tau_1) < \infty, \quad \mathbb{P}^o\text{-a.s.} \tag{2.1.16}$$

On the other hand, we have the following:

Lemma 2.1.17 *Assume $T_n/n \rightarrow \alpha$, for some constant $\alpha < \infty$. Then,*

$$\frac{X_n}{n} \xrightarrow{n \rightarrow \infty} \frac{1}{\alpha}.$$

Proof of Lemma 2.1.17

Let k_n be the unique (random) integers such that

$$T_{k_n} \leq n < T_{k_n+1}.$$

Note that $X_n < k_n + 1$ while $X_n \geq k_n - (n - T_{k_n})$. Hence,

$$\frac{k_n}{n} - \left(1 - \frac{T_{k_n}}{n}\right) \leq \frac{X_n}{n} \leq \frac{k_n + 1}{n}.$$

But, $\lim_{n \rightarrow \infty} k_n/n = \lim_{n \rightarrow \infty} n/T_n$ (due to the existence of the second limit and the definition of k_n). Thus,

$$\frac{1}{\alpha} \geq \limsup_{n \rightarrow \infty} \frac{X_n}{n} \geq \liminf_{n \rightarrow \infty} \frac{X_n}{n} \geq \frac{1}{\alpha}. \quad \square$$

Lemma 2.1.17 and (2.1.16) complete the proof of Theorem 2.1.9 in case (a). Case (b) is similar, while case (c) is a minor modification of the above argument and is left out. □

Bibliographical notes: The proof of Theorem 2.1.2 is essentially from [67], except that the use of Kesten’s lemma is borrowed from [1]. See also [50] for an “ergodic” approach. The rest of the section is an adaptation of the argument in [67], which requires a strongly mixing assumption. The proof of ergodicity in Lemma 2.1.10 was suggested to me by P. Bougerol. F. Rassoul-Agha has kindly shown me a different proof of this fact.

Transience and recurrence results for non nearest-neighbour RWRE on \mathbb{Z} , in terms of certain Lyapunov exponents of products of random matrices, are developed in [43], see also [50] and [7]. This is further developed in [3], [39], where transience and recurrence criteria for RWRE on graphs of the form $\mathbb{Z} \times G$, G finite, are derived.

LLN-version II: auxiliary Markov chains

We use the evaluation of the LLN as an excuse for introducing the machinery of the “environment viewed from the particle”. The first step consists of introducing an auxiliary Markov chain.

Starting from the RWRE X_n , define $\bar{\omega}(n) = \theta^{X_n} \omega$. The sequence $\{\bar{\omega}_n\}$ is a process with paths in $\Omega^{\mathbb{N}}$. What is maybe more useful is that it is in fact a Markov process. More precisely:

Lemma 2.1.18 *The process $\{\bar{\omega}(n)\}$ is a Markov process under either P_ω^0 or \mathbb{P}^0 , with state space Ω and transition kernel*

$$M(\omega, d\omega') = \omega_0^+ \delta_{\theta\omega=\omega'} + \omega_0^- \delta_{\theta^{-1}\omega=\omega'} + \omega_0^0 \delta_{\omega=\omega'}.$$

Proof. For bounded functions $f_i : \Omega \rightarrow \mathbb{R}$,

$$\begin{aligned} E_\omega^o \left(\prod_{i=1}^n f_i(\bar{\omega}(i)) \right) &= E_\omega^o \left(\prod_{i=1}^n f_i(\theta^{X_i} \omega) \right) \\ &= E_\omega^o \left(\prod_{i=1}^{n-1} f_i(\theta^{X_i} \omega) E_\omega^{X_{n-1}} (f_n(\theta^{X_n} \omega)) \right) \\ &= E_\omega^o \left(\prod_{i=1}^{n-1} f_i(\theta^{X_i} \omega) \left[\omega_{X_{n-1}}^+ f_n(\theta \cdot \theta^{X_{n-1}} \omega) \right. \right. \\ &\quad \left. \left. + \omega_{X_{n-1}}^- f_n(\theta^{-1} \cdot \theta^{X_{n-1}} \omega) + \omega_{X_{n-1}}^0 f_n(\theta^{X_{n-1}} \omega) \right] \right) \\ &= E_\omega^o \left(\prod_{i=1}^{n-1} f_i(\theta^{X_i} \omega) M f_n(\bar{\omega}(n-1)) \right) \end{aligned} \tag{2.1.19}$$

where

$$M f(\omega) = \int f(\omega') M(\omega, d\omega'),$$

which proves the Markov property of $\{\bar{\omega}(n)\}$ under P_ω^o . Integrating both sides of (2.1.19) with respect to P yields the Markov property under \mathbb{P}^o . \square

Our next step is to construct an invariant measure for the transition kernel M . In most of this section we will assume that $E_P \log \rho_0 < 0$, implying, by Theorem 2.1.2, that $T_1 < \infty$, \mathbb{P}^o -a.s. Whenever $E_{\mathbb{P}^o}(T_1) < \infty$, define the measures

$$Q(B) = E_{\mathbb{P}^o} \left(\sum_{i=0}^{T_1-1} \mathbf{1}_{\{\bar{\omega}(i) \in B\}} \right), \quad \bar{Q}(B) = \frac{Q(B)}{Q(\Omega)} = \frac{Q(B)}{E_{\mathbb{P}^o} T_1}.$$

Using Lemma 2.1.12, one checks that under Assumption 2.1.1 and if $E_P(\bar{S}) < \infty$ then $E_{\mathbb{P}^o} T_1 < \infty$, and $\bar{Q}(\cdot)$ in this case is a probability measure.

Lemma 2.1.20 *Assume Assumption 2.1.1 and $E_P(\bar{S}) < \infty$. Then, $Q(\cdot)$ is invariant under the Markov kernel M , that is*

$$Q(B) = \iint \mathbf{1}_{\omega' \in B} M(\omega, d\omega') Q(d\omega).$$

Proof. We have

$$\begin{aligned}
 & \iint \mathbf{1}_{\omega' \in B} M(\omega, d\omega') Q(d\omega) \\
 &= \sum_{k=0}^{\infty} E_{\mathbb{P}^o} \left(T_1 > k; \mathbf{1}_{\bar{\omega}(k+1) \in B} \right) \\
 &= \sum_{k=0}^{\infty} E_{\mathbb{P}^o} \left(T_1 = k + 1; \mathbf{1}_{\bar{\omega}(k+1) \in B} \right) + \sum_{k=0}^{\infty} E_{\mathbb{P}^o} \left(T_1 > k + 1; \mathbf{1}_{\bar{\omega}(k+1) \in B} \right) \\
 &= \mathbb{P}^o(T_1 < \infty; \bar{\omega}(T_1) \in B) + \sum_{k=1}^{\infty} \mathbb{P}^o(T_1 > k; \bar{\omega}(k) \in B).
 \end{aligned}$$

But $\mathbb{P}^o(T_1 < \infty) = 1$ while $\mathbb{P}^o(\bar{\omega}(T_1) \in B) = P(\theta\omega \in B) = P(\omega \in B)$, hence

$$= \sum_{k=0}^{\infty} \mathbb{P}^o(T_1 > k; \bar{\omega}(k) \in B) = Q(B). \quad \square$$

Define next

$$\Lambda(\omega) = \frac{1}{\omega_0^+} \left[1 + \sum_{i=1}^{\infty} \prod_{j=1}^i \rho_j \right].$$

It is not hard to check, by the shift invariance of P , that the condition $E_P(\Lambda(\omega)) < \infty$ is equivalent to $E_P(\bar{S}) < \infty$, c.f. Section 2.1. We next claim the

Lemma 2.1.21 *Under the assumptions of Lemma 2.1.20, it holds that*

$$\frac{dQ}{dP} = \Lambda(\omega).$$

Proof. Note first that by Jensen’s inequality, $E_P(\Lambda) < \infty$ implies that $E_P(\log \rho_0) < 0$ and hence $X_n \rightarrow_{n \rightarrow \infty} \infty$, \mathbb{P}^o -a.s., by Theorem 2.1.2. Let $f : \Omega \rightarrow \mathbb{R}$ be measurable. Then,

$$\int f dQ = E_{\mathbb{P}^o} \left(\sum_{i=0}^{T_1-1} f(\bar{\omega}_i) \right) = E_{\mathbb{P}^o} \left(\sum_{i \leq 0} f(\theta^i \omega) N_i \right)$$

where $N_i = \{\#k \in [0, T_1) : X_k = i\}$ (note the difference in the role the index i plays in the two sums!). Using the shift invariance of P , we get

$$\begin{aligned}
 \int f dQ &= \sum_{i \leq 0} E_P \left(f(\theta^i \omega) E_{\omega}^o N_i \right) \\
 &= \sum_{i \leq 0} E_P \left(f(\omega) E_{\theta^{-i}\omega}^o N_i \right) = E_P \left(f(\omega) \left(\sum_{i \leq 0} E_{\theta^{-i}\omega}^o N_i \right) \right).
 \end{aligned}$$

Hence,

$$\frac{dQ}{dP} = \sum_{i \leq 0} E_{\theta^{-i}\omega}^o N_i, \tag{2.1.22}$$

and the right hand side converges, P -a.s.

In order to prove both the convergence in (2.1.22) and the lemma, we turn to evaluate $E_{\omega}^o N_i$. Define, for $i \leq 0$,

$$\begin{aligned} \eta_{i,0} &= \min\{k \leq T_1 : X_k = i\} \\ \theta_{i,0} &= \min\{\eta_{i,0} < k \leq T_1 : X_{k-1} = i, X_k = i - 1\} \end{aligned}$$

and, for $j \geq 1$,

$$\begin{aligned} \eta_{i,j} &= \min\{\theta_{i,j-1} < k \leq T_1 : X_k = i\} \\ \theta_{i,j} &= \min\{\eta_{i,j} < k \leq T_1 : X_{k-1} = i, X_k = i - 1\} \end{aligned}$$

(with the usual convention that the minimum over an empty set is $+\infty$). We refer to the time interval $(\theta_{i,j-1}, \eta_{i,j})$ as the j -th excursion from $i - 1$ to i . For any $j \geq 0$, any $i \leq 0$, define

$$\begin{aligned} U_{i,j} &= \{\#\ell \geq 0 : \theta_{i+1,j} < \theta_{i,\ell} < \eta_{i+1,j+1}\} \\ Z_{i,j} &= \{\#k \geq 0 : X_{k-1} = i, X_k = i, \theta_{i+1,j} < k < \eta_{i+1,j+1}\}. \end{aligned}$$

Note that $U_{i,j}$ is the number of steps from i to $i - 1$ during the $j + 1$ -th excursion from i to $i + 1$, whereas $Z_{i,j}$ is the number of steps from i to i during the same excursion. The Markov property implies that

$$\begin{aligned} &P_{\omega}^o\left(U_{i,\ell} = k_{\ell}, Z_{i,\ell} = m_{\ell}, \ell = 1, \dots, L \mid \{U_{i',j}\}_{i' > i}, \eta_{i+1,L+1} < \infty\right) \\ &= \prod_{\ell=1}^L \left[\left(\frac{\omega_i^-}{\omega_i^- + \omega_i^+}\right)^{k_{\ell}} \left(\frac{\omega_i^+}{\omega_i^- + \omega_i^+}\right) \left(\frac{\omega_i^0}{\omega_i^0 + \omega_i^+}\right)^{m_{\ell}} \left(\frac{\omega_i^+}{\omega_i^0 + \omega_i^+}\right) \right]. \end{aligned} \tag{2.1.23}$$

Defining $U_i = \sum_j U_{i,j}$, $Z_i = \sum_j Z_{i,j}$, and noting that $\mathbb{P}^o(\{U_i < \infty\} \cap \{Z_i < \infty\}) = 1$ because $X_n \rightarrow \infty$, \mathbb{P}^o -a.s., (2.1.23) implies that $\{U_i\}$ is under P_{ω}^o an (inhomogeneous) branching process with geometric offspring distribution of parameter $\frac{\omega_i^-}{\omega_i^- + \omega_i^+}$. Further,

$$\begin{aligned} E_{\omega}^o(U_i | U_{i+1}, \dots, U_0) &= \rho_i U_{i+1} \\ E_{\omega}^o(Z_i | U_{i+1}, \dots, U_0) &= \frac{\omega_i^0}{\omega_i^+} U_{i+1} \end{aligned} \tag{2.1.24}$$

and using the relation $N_i = U_i + U_{i+1} + Z_i$, \mathbb{P}^o -a.s., we get

$$E_{\omega}^o(N_i | U_{i+1}, \dots, U_0) = E_{\omega}^o\left(U_i + U_{i+1} + Z_i | U_{i+1}, \dots, U_0\right) = \frac{1}{\omega_i^+} E_{\omega}^o U_{i+1}.$$

Iterating (2.1.24), one gets

$$E_\omega^\circ N_i = \frac{1}{\omega_i^+} \rho_0 \cdots \rho_{i+1}.$$

Hence, using (2.1.22), and the assumption,

$$\frac{dQ}{dP} = \frac{1}{\omega_0^+} \left[1 + \sum_{i=1}^\infty \prod_{j=1}^i \rho_j \right] < \infty, P\text{-a.s.}$$

which completes the proof of the Lemma. □

Remark: Note that $dQ/dP > 0$, P -a.s., and hence under the assumption $E_P(S) < \infty$ it holds that $Q \sim P$. This fact is true in greater generality, see the discussion in [69] and in Section 3.3 below.

Corollary 2.1.25 *Under the law induced by $\bar{Q} \otimes P_\omega^\circ$, the sequence $\{\bar{\omega}(n)\}$ is stationary and ergodic.*

Proof. The stationarity follows from the stationarity of \bar{Q} . Let $\bar{\theta}$ denote the shift on $\bar{\Omega} = \Omega^\mathbb{N}$, that is, for $\bar{\omega} \in \bar{\Omega}$, $\bar{\theta}\bar{\omega}(n) = \bar{\omega}(n + 1)$. Denote by \bar{P}_ω the law of the sequence $\{\bar{\omega}(n)\}$ with $\bar{\omega}(0) = \omega$, that is, for any measurable sets $B_i \subset \Omega$,

$$\begin{aligned} \bar{P}_\omega(\bar{\omega}(i) \in B_i, i = 1, \dots, \ell) \\ = \int_{B_1} \cdots \int_{B_\ell} M(\omega, d\omega^1) M(\omega^1, d\omega^2) \cdots M(\omega^{\ell-1}, d\omega^\ell) \end{aligned}$$

and set $\bar{Q} = \bar{Q} \otimes \bar{P}_\omega$ (as usual, we also use \bar{Q} to denote the corresponding marginal induced on $\bar{\Omega}$).

We need to show that for any invariant A , that is $A \in \bar{\Omega}$ such that $\bar{\theta}A = A$, $\bar{Q}(A) \in \{0, 1\}$. Set $\varphi(\omega) = \bar{P}_\omega(A)$, we claim that $\{\varphi(\bar{\omega}(n))\}$ is a martingale with respect to the filtration $\mathcal{G}_n = \sigma(\bar{\omega}(0), \dots, \bar{\omega}(n))$: indeed,

$$\varphi(\bar{\omega}(n)) = \bar{P}_{\bar{\omega}(n)}(A) = E_{\bar{Q}}(\mathbf{1}_{\theta^n A} | \mathcal{G}_n) = E_{\bar{Q}}(\mathbf{1}_A | \mathcal{G}_n),$$

where the second equality is due to the Markov property and the third due to the invariance of A . Hence, by the martingale convergence theorem,

$$\varphi(\bar{\omega}(n)) \xrightarrow[n \rightarrow \infty]{} \mathbf{1}_A, \quad \bar{Q}\text{-a.s.} \tag{2.1.26}$$

Further, $Q(\varphi(\omega) \notin \{0, 1\}) = 0$ because otherwise there exists an interval $[a, b]$ with $\{0\}, \{1\} \notin [a, b]$ and $Q(\varphi(\omega) \in [a, b]) > 0$, while

$$\frac{1}{n} \sum_0^{n-1} \mathbf{1}_{\{\varphi(\bar{\omega}(n)) \in [a, b]\}} \rightarrow E_{\bar{Q}}(\mathbf{1}_{\{\varphi(\bar{\omega}(0)) \in [a, b]\}} | \mathcal{J}), \tag{2.1.27}$$

where \mathbb{J} is the invariant σ -field.

Taking expectations in (2.1.27) and using (2.1.26), one concludes that

$$0 = \overline{\mathbb{Q}}\left(\varphi(\overline{\omega}(0)) \in [a, b]\right) = \overline{\mathbb{Q}}\left(\varphi(\omega) \in [a, b]\right),$$

a contradiction. Thus for some measurable $B \subset \Omega$, $\varphi(\omega) = \mathbf{1}_B, \overline{\mathbb{Q}} - \text{a.s.}$ Further, the Markov property and invariance of A yield that $M\mathbf{1}_B = \mathbf{1}_B, \overline{\mathbb{Q}}\text{-a.s.}$ and hence $P\text{-a.s.}$ But then,

$$\mathbf{1}_B = M\mathbf{1}_B \geq \omega_0^+ \mathbf{1}_{\theta B}, P\text{-a.s.}$$

Since $E A(\omega) < \infty$ implies $P(\omega_0^+ = 0) = 0$, it follows that $\mathbf{1}_B \geq \mathbf{1}_{\theta B}, P\text{-a.s.}$, and then $E_P(\mathbf{1}_B) = E_P(\mathbf{1}_{\theta B})$ implies that $\mathbf{1}_B = \mathbf{1}_{\theta B}, P\text{-a.s.}$ But then, by ergodicity of P , $P(B) \in \{0, 1\}$, and hence $\overline{\mathbb{Q}}(B) \in \{0, 1\}$. Since $\overline{\mathbb{Q}}(A) = E_{\overline{\mathbb{Q}}}\varphi(\omega) = \overline{\mathbb{Q}}(B)$, the conclusion follows. \square

We are now ready to give the:

Proof of Theorem 2.1.9 - Environment version We begin with case (a), noting that the proof of case (b) is identical by the transformation $\omega_i \mapsto \hat{\omega}_{-i}$, where $\hat{\omega}_i^+ = \omega_i^-, \hat{\omega}_i^- = \omega_i^+$. Set $d(x, \omega) = E_\omega^x(X_1 - x)$. Then

$$\begin{aligned} X_n &= \sum_{i=1}^n (X_i - X_{i-1}) = \sum_{i=1}^n \left(X_i - X_{i-1} - d(X_{i-1}, \omega) \right) + \sum_{i=1}^n d(X_{i-1}, \omega) \\ &:= M_n + \sum_{i=1}^n d(X_{i-1}, \omega). \end{aligned} \tag{2.1.28}$$

But, under P_ω^o , M_n is a martingale, with $|M_{n+1} - M_n| \leq 2$; Hence, with $\mathcal{G}_n = \sigma(M_1, \dots, M_n)$,

$$\begin{aligned} E_\omega^o(e^{\lambda M_n}) &= E_\omega^o\left(e^{\lambda M_{n-1}} E_\omega^o(e^{\lambda(M_n - M_{n-1})} | \mathcal{G}_{n-1})\right) \\ &\leq E_\omega^o\left(e^{\lambda M_{n-1}} e^{2\lambda^2}\right) \end{aligned}$$

and hence, iterating, $E_\omega^o(e^{\lambda M_n}) \leq e^{2n\lambda^2}$ (this is a version of Azuma’s inequality, see [19, Corollary 2.4.7]). Chebycheff’s inequality then implies

$$\frac{M_n}{n} \rightarrow 0, \quad \mathbb{P}^o\text{-a.s.}$$

(and even with exponential rate). Next, note that

$$\sum_{i=1}^n d(X_{i-1}, \omega) = \sum_{i=1}^n d(0, \overline{\omega}(i-1)).$$

The ergodicity of $\{\overline{\omega}(i)\}$ under $\overline{\mathbb{Q}} \otimes P_\omega^o$ implies that

$$\frac{1}{n} \sum_{i=1}^n d(0, \bar{\omega}(i-1)) \longrightarrow E_{\bar{Q}}(d(0, \bar{\omega}(0))), \bar{Q} \otimes P_{\omega}^o\text{-a.s.} \quad (2.1.29)$$

But,

$$\begin{aligned} & E_{\bar{Q}}(d(0, \bar{\omega}(0))) \\ &= \frac{E_P \left[\Lambda(\omega)(\omega_0^+ - \omega_0^-) \right]}{E_P(\Lambda(\omega))} \\ &= \frac{1 + E_P \left(\omega_1^- \left[\frac{1}{\omega_1^+} + \sum_{i=2}^{\infty} \prod_{j=2}^i \rho_j \right] - \omega_0^- \left[\frac{1}{\omega_0^+} + \sum_{i=1}^{\infty} \prod_{j=1}^i \rho_j \right] \right)}{E_P(\Lambda(\omega))} \\ &= \frac{1}{E_P(\Lambda(\omega))} = \frac{1}{E_P(\bar{S}(\omega))}. \end{aligned}$$

Finally, since $E_P(\Lambda(\omega)) < \infty$, (2.1.29) holds also \mathbb{P}^o -a.s., completing the proof of the theorem in cases (a),(b).

Case (c) is handled by appealing to Lemma 2.1.12. Suppose $\limsup X_n = +\infty, \mathbb{P}^o - \text{a.s.}$ Then, $\tau_1 < \infty, \mathbb{P}^o\text{-a.s.}$ Define $\tau_i^K = \min(\tau_i, K)$. Note that under P_{ω}^o , the random variables $\{\tau_i^K\}$ are independent and bounded, and hence, with $G_n^K = n^{-1} \sum_{i=1}^n \tau_i^K$, we have

$$|G_n^K - E_{\omega}^o G_n^K| \rightarrow_{n \rightarrow \infty} 0, \quad P_{\omega}^o - \text{a.s.}$$

But $f(\omega) := E_{\omega}^o \tau_1^K$ is a bounded, measurable, local function on Ω , and $E_{\omega}^o G_n^K = n^{-1} \sum_{i=1}^n f(\theta^i \omega)$. Hence, by the ergodic theorem, $E_{\omega}^o G_n^K \rightarrow_{n \rightarrow \infty} E_{\mathbb{P}^o} \tau_1^K, P - \text{a.s.}$ Since, by Lemma 2.1.12 we have $E_{\mathbb{P}^o} \tau_1^K \rightarrow_{K \rightarrow \infty} \infty$, we conclude that

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \tau_i \geq \lim_{K \rightarrow \infty} E_{\mathbb{P}^o} \tau_1^K = \infty, \mathbb{P}^o - \text{a.s.}$$

This immediately implies $\limsup_{n \rightarrow \infty} X_n/n \leq 0, \mathbb{P}^o - \text{a.s.}$ The reverse inequality is proved by considering the sequence $\{\tau_{-i}\}$, yielding part (c) of the Theorem. □

Remark: Exactly as in Lemma 2.1.17, it is not hard to check that under Assumption 2.1.1, it holds that

$$\lim_{n \rightarrow \infty} \frac{T_n}{n} = E_P(\bar{S}), \quad \mathbb{P}^o - \text{a.s.} \quad (2.1.30)$$

Bibliographical notes: The construction presented here goes back at least to [45]. Our presentation is heavily influenced by [1] and [69].

2.2 CLT for ergodic environments

In this section, we continue to look at the environment from the point of view of the particle. Our main goal is to prove the following:

Theorem 2.2.1 *Assume 2.1.1. Further, assume that for some $\varepsilon > 0$,*

$$E_Q(\overline{S}^{2+\varepsilon}(\omega) + \overline{S}(\theta^{-1}\omega)^{2+\varepsilon}) < \infty, \tag{2.2.2}$$

and that

$$\sum_{n \geq 1} \sqrt{E_P\left(E_P\left(v_P \overline{S}(\omega) - 1 \mid \sigma(\omega_i, i \leq -n)\right)^2\right)} < \infty, \tag{2.2.3}$$

where $v_P := 1/E_P(\overline{S}(\omega))$. Then, with

$$\sigma_{P,1}^2 := v_P^2 E_Q\left(\omega_0^+(\overline{S}(\omega) - 1)^2 + \omega_0^-(\overline{S}(\theta^{-1}\omega) + 1)^2 + \omega_0^0\right),$$

and

$$\sigma_{P,2}^2 := E_P(v_P \overline{S}(\omega) - 1)^2 + 2 \sum_{n=1}^{\infty} E_P\left((v_P \overline{S}(\omega) - 1)(v_P \overline{S}(\theta^n \omega) - 1)\right),$$

we have that

$$\mathbb{P}^o\left(\frac{X_n - nv_P}{\sigma_P \sqrt{n}} > x\right) \rightarrow_{n \rightarrow \infty} \Phi(-x),$$

where

$$\Phi(x) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{\theta^2}{2}} d\theta,$$

and $\sigma_P^2 = \sigma_{P,1}^2 + v_P \sigma_{P,2}^2$.

Proof. The basic idea in the proof is to construct an appropriate martingale, and then use the Martingale CLT and the CLT for stationary ergodic sequences. We thus begin with recalling the version of these CLT's most useful to us.

Lemma 2.2.4 ([26], pg. 417) *Suppose $(Z_n, \mathcal{F}_n)_{n \geq 0}$ is a martingale difference sequence, and let $V_n = \sum_{1 \leq k \leq n} E(Z_k^2 | \mathcal{F}_{k-1})$. Assume that*

- (a) $\frac{V_n}{n} \rightarrow_{n \rightarrow \infty} \sigma^2$, in probability.
- (b) $\frac{1}{n} \sum_{m \leq n} E\left(Z_m^2 \mathbf{1}_{\{|Z_m| > \varepsilon \sqrt{n}\}}\right) \rightarrow_{n \rightarrow \infty} 0$.

Then, $\sum_{i=1}^n Z_i / \sigma \sqrt{n}$ converges in distribution to a standard Gaussian random variable.

Lemma 2.2.5 ([26], p. 419) *Suppose $\{Z_n\}_{n \in \mathbb{Z}}$ is a stationary, zero mean, ergodic sequence, and set $\mathcal{F}_n = \sigma(Z_i, i \leq n)$. Assume that*

$$\sum_{n \geq 0} \sqrt{E(E(Z_0 | \mathcal{F}_{-n}))^2} < \infty. \tag{2.2.6}$$

Then, $\left\{ \sum_{i=1}^{nt} Z_i / \sigma \sqrt{n} \right\}_{t \in [0,1]}$ converges in distribution to a standard Brownian motion, where

$$\sigma^2 = EZ_0^2 + 2 \sum_{n=1}^{\infty} E(Z_0 Z_n).$$

We next recall that by Theorem 2.1.9,

$$\frac{X_n}{n} \rightarrow v_P, \quad \mathbb{P}^o\text{-a.s.},$$

where $v_P := 1/E_P(\bar{S})$. One is tempted to use the martingale M_n appearing in the environment proof of Theorem 2.1.9 (see (2.1.28)), however this strategy is not so successful because of the difficulties associated with separating the fluctuations in M_n and $\sum_{i=1}^n d(X_{i-1}, \omega)$. Instead, write

$$f(x, n, \omega) = x - v_P n + h(x, \omega), \quad x \in \mathbb{Z}.$$

We want to make $f(X_n, n, \omega)$ into a martingale w.r.t. $\mathcal{F}_n := \sigma(X_1, \dots, X_n)$ and the law P_ω^o . This is automatic if we can ensure that

$$E_\omega^{X_n} f(X_{n+1}, n+1, \omega) = f(X_n, n, \omega), \quad P_\omega^o\text{-a.s.} \tag{2.2.7}$$

Developing this equality and defining $\Delta(x, \omega) = h(x+1, \omega) - h(x, \omega)$, we get that (2.2.7) holds true if a bounded solution to the equation

$$\Delta(x, \omega) = - \left[\frac{\omega_x^+ - \omega_x^- - v_P}{\omega_x^+} \right] + \frac{\omega_x^-}{\omega_x^+} \Delta(x-1, \omega)$$

exists. One may verify that $\Delta(x, \omega) = -1 + v_P \bar{S}(\theta^x \omega)$ is such a solution.

Fixing $h(0, \omega) = 0$, and defining $\bar{M}_0 = 0$ and $\bar{M}_n = f(X_n, n, \omega)$, one concludes that \bar{M}_n is a martingale, and further

$$\begin{aligned} E_\omega^o \left((\bar{M}_{k+1} - \bar{M}_k)^2 | \mathcal{F}_k \right) &= \omega_{X_k}^+ v_P^2 (\bar{S}(\theta^{X_k} \omega) - 1)^2 + \omega_{X_k}^- v_P^2 (\bar{S}(\theta^{X_k-1} \omega) + 1)^2 + \omega_{X_k}^0 v_P^2 \\ &= v_P^2 \left[\bar{\omega}(k)_0^+ (\bar{S}(\bar{\omega}_k) - 1)^2 + \bar{\omega}(k)_0^- (\bar{S}(\theta^{-1} \bar{\omega}_k) + 1)^2 + \bar{\omega}(k)_0^0 \right]. \end{aligned}$$

Hence,

$$\frac{V_n}{n} = \frac{1}{n} \sum_{k=1}^n E_\omega^o \left((\bar{M}_{k+1} - \bar{M}_k)^2 | \mathcal{F}_k \right) \xrightarrow[n \rightarrow \infty]{} \sigma_{P,1}^2, \quad \mathbb{P}^o\text{-a.s.},$$

using the machinery developed in Section 2.1. The integrability condition (2.2.2) is enough to apply the Martingale CLT (Lemma 2.2.4), and one concludes that for any $\delta > 0$,

$$P \left(\left| P_\omega^o \left(\frac{\overline{M}_n}{\sigma_{P,1}\sqrt{n}} \geq x \right) - \Phi(-x) \right| > \delta \right) \rightarrow_{n \rightarrow \infty} 0. \tag{2.2.8}$$

Note that since both $P_\omega^o(\overline{M}_n \geq x\sigma_{P,1}\sqrt{n})$ and $\Phi(x)$ are monotone in x , and that $\Phi(\cdot)$ is continuous, the convergence in (2.2.8) actually is uniform on \mathbb{R} . Further, note that

$$h(X_n, \omega) = \sum_{j=1}^{X_n-1} \Delta(j, \omega) = \sum_{j=1}^{nv_P} \Delta(j, \omega) + R_n := Z_n + R_n.$$

Note that, for every $\delta > 0$ and some $\delta_n \rightarrow 0$,

$$\begin{aligned} \mathbb{P}^o \left(\frac{|R_n|}{\sqrt{n}} \geq \delta \right) &\leq \mathbb{P}^o \left(|X_n - nv_P| \geq \delta_n n \right) \\ + P \left(\max_{j_-, j_+ \in (-n\delta_n, n\delta_n)} \left| \sum_{i=j_-}^{j_+} \frac{\Delta(i, \omega)}{\sqrt{n}} \right| \geq \delta \right) &:= P_{1,n}(\delta_n) + P_{2,n}(\delta, \delta_n) \xrightarrow[n \rightarrow \infty]{} 0, \end{aligned} \tag{2.2.9}$$

where the convergence of the first term is due (choosing an appropriate $\delta_n \rightarrow_{n \rightarrow \infty} 0$ slowly enough) to Theorem 2.1.9 and that of the second one due to $E_P \Delta(i, \omega) = 0$ and the stationary invariance principle (Lemma 2.2.5), which can be applied, for any $\delta_n \rightarrow_{n \rightarrow \infty} 0$, due to (2.2.3).

Another application of Lemma 2.2.5 yields that

$$\lim_{n \rightarrow \infty} P(Z_n \geq z\sqrt{nv_P}\sigma_{P,2}) = \Phi(-z). \tag{2.2.10}$$

Writing $X_n - nv_P = \overline{M}_n - Z_n - R_n$, and using that $R_n/\sqrt{n} \rightarrow_{n \rightarrow \infty} 0$ in \mathbb{P}^o -probability, one concludes that

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}^o \left(\frac{X_n - nv_P}{\sqrt{n}} > x \right) &= \lim_{n \rightarrow \infty} E_P(P_\omega^o(\overline{M}_n/\sqrt{n} > x + Z_n/\sqrt{n})) \\ &= \lim_{n \rightarrow \infty} E_P \left(\Phi \left(-\frac{x + Z_n/\sqrt{n}}{\sigma_{P,1}} \right) \right), \end{aligned} \tag{2.2.11}$$

where the second equality is due to the uniform convergence in (2.2.8). Combining (2.2.11) with (2.2.10) yields the claim. \square

Remark: The alert reader will have noted that under assumptions (2.1.1) and (2.2.2), and a mild mixing assumption on P which ensures that for any $\delta > 0$ and $\delta_n \rightarrow 0$, $P_{2,n}(\delta, \delta_n) \rightarrow_{n \rightarrow \infty} 0$, c.f. (2.2.9),

$$P_\omega^o \left(\frac{X_n - v_P n - Z_n}{\sqrt{n}\sigma_{P,1}} > x \right) \rightarrow_{n \rightarrow \infty} \Phi(-x).$$

That is, using a random centering one also has a *quenched* CLT.

Exercise 2.2.12 Check that the integrability conditions (2.2.2) and (2.2.3) allow for the application of Lemmas 2.2.4 and 2.2.5 in the course of the proof of Theorem 2.2.1.

Exercise 2.2.13 Check that in the case of P being a product measure, the assumption (2.2.2) in Theorem 2.2.1 can be dropped.

Bibliographical notes: The presentation here follows the ideas of [45], as developed in [53]. The latter provides an explicit derivation of the CLT in case $P(\omega_0^0 = 0) = 1$, but it seems that in his derivation only the quenched CLT is derived and the random centering then is missing. A different approach to the CLT is presented in [1], using the hitting times $\{\tau_i\}$; It is well suited to yield the quenched CLT, and under strong assumptions on P which ensure that the random quenched centering vanishes P -a.s., also the annealed CLT. Note however that the case of P being a product measure is not covered in the hypotheses of [1]. See [7] for some further discussion and extensions.

There are situations where limit laws which are not of the CLT type can be exhibited. The proof of such results uses hitting time decompositions, and techniques as discussed in Section 2.4. We refer to Section 2.5 and its bibliographical notes for an example of such a situation and additional information.

2.3 Large deviations

Having settled the issue of the LLN, the next logical step (even if not following the historical development) is the evaluation of the probabilities of large deviations. As already noted in the evaluation of the CLT in Section 2.2, there can be serious differences between quenched and annealed probabilities of deviations. In order to address this, we make the following definitions; throughout this section, \mathcal{X} denotes a completely regular topological space.

Definition 2.3.1 A function $I : \mathcal{X} \rightarrow [0, \infty]$ is a rate function if it is lower semicontinuous. It is a good rate function if its level sets are compact.

Definition 2.3.2 A sequence of \mathcal{X} valued random variables $\{Z_n\}$ satisfies the quenched Large Deviations Principle (LDP) with speed n and deterministic rate function I if for any Borel set A ,

$$-I(A^\circ) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^\circ(Z_n \in A) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^\circ(Z_n \in A) \leq -I(\bar{A})$$

P -a.s. (2.3.3)

where A° denotes the interior of A , \bar{A} the closure of A , and for any Borel set F ,

$$I(F) = \inf_{x \in F} I(x) . \tag{2.3.4}$$

Definition 2.3.5 A sequence of \mathcal{X} valued random variables $\{Z_n\}$ satisfies the annealed LDP with speed n and rate function I if, for any Borel set A ,

$$-I(A^o) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(Z_n \in A) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(Z_n \in A) \leq -I(\bar{A}). \tag{2.3.6}$$

Finally, we note the

Definition 2.3.7 A LDP is called weak if the upper bound in (2.3.3) or (2.3.6), holds only with \bar{A} compact.

For background on the LDP we refer to [19]. It is well known, c.f. [19, Lemma 4.1.4] that if the LDP holds then the rate function is uniquely defined. The following easy lemma is intuitively clear: annealed deviation probabilities allow for atypical fluctuations of the environment and hence are not smaller than corresponding quenched deviation probabilities:

Lemma 2.3.8 Let $\{A_n\}$ be a sequence of events, subsets of $\Omega \times \mathbb{Z}^N$. Then,

$$c := \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(A_n) \geq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(A_n), \quad P - a.s. \tag{2.3.9}$$

Further,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(A_n) \geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(A_n), \quad P - a.s. \tag{2.3.10}$$

In particular, if a sequence of \mathcal{X} valued random variables $\{Z_n\}$ satisfies annealed and quenched LDP's with rate functions $I_a(\cdot), I_q(\cdot)$, respectively, then,

$$I_a(x) \leq I_q(x), \forall x \in \mathcal{X}.$$

Proof. Assume first $c < 0$. Fix $\delta > 0$ and let $B_n^\delta = \{\omega : P_\omega^o(A_n) \geq \exp((c + \delta)n)\}$. Then, by the definition of c , see (2.3.9), and Markov's bound, for n large enough,

$$P(B_n^\delta) \leq e^{-\delta n/2}.$$

Hence, $\omega \in B_n^\delta$ occurs only finitely many times, P -a.s., implying that for P -almost all ω there exists an $n_0(\omega)$ such that for all $n \geq n_0(\omega)$, $P_\omega^o(A_n) < \exp((c + \delta)n)$. Hence,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(A_n) \leq c + \delta, \quad P - a.s.$$

(2.3.9) follows by the arbitrariness of $\delta > 0$. Next, set $\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(A_n) := c_1 \leq c$. Define $\{n_k\}$ such that

$$\lim_{k \rightarrow \infty} \frac{1}{n_k} \log \mathbb{P}^o(A_{n_k}) = c_1.$$

Apply now the first part of the lemma to conclude that

$$c_1 \geq \limsup_{k \rightarrow \infty} \frac{1}{n_k} \log P_\omega^o(A_{n_k}) \geq \liminf_{k \rightarrow \infty} \frac{1}{n_k} \log P_\omega^o(A_{n_k}) \geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(A_n)$$

$P - \text{a.s.}$

The case $c = 0$ is the same, except that (2.3.9) is trivial. This completes the proof. \square

Quenched LDP’s

The LDP in the quenched setting makes use in its proof of the hitting times $\{\tau_i\}$. Introduce, for any $\lambda \in \mathbb{R}$,

$$\begin{aligned} \varphi(\lambda, \omega) &= E_\omega^o(e^{\lambda\tau_1} \mathbf{1}_{\{\tau_1 < \infty\}}), \quad f(\lambda, \omega) = \log \varphi(\lambda, \omega) \\ G(\lambda, P, u) &= \lambda u - E_P(f(\lambda, \omega)). \end{aligned}$$

We need throughout the following modification of Assumption 2.1.1.

Assumption 2.3.11

- (B1) P is stationary and ergodic,
- (B2) There exists an $\varepsilon > 0$ such that $P(\omega_0^+ \notin (0, \varepsilon))P(\omega_0^- \notin (0, \varepsilon)) = 1$,
- (B3) $P(\omega_0^+ + \omega_0^- > 0) = 1$, $P(\omega_0^0 > 0, \omega_0^+ \omega_0^- = 0) = 0$, and $P(\omega_0^+ = 0)P(\omega_0^- = 0) = 0$.

Note that we allow for the possibility of having one sided transitions (e.g., moves to the right only) of the RWRE. This allows one to deal with the case where “random nodes” are present.

Define

$$\begin{aligned} \rho_{\min} &:= \inf[\rho : P(\rho_0 < \rho) > 0], \\ \rho_{\max} &:= \sup[\rho : P(\rho_0 > \rho) > 0], \\ \omega_{\max}^0 &:= \sup[\alpha : P(\omega_0^0 > \alpha) > 0]. \end{aligned}$$

With P_N denoting the restriction of P to the first N coordinates $\{\omega_i\}_{i=0}^{N-1}$, we say that P is locally equivalent to the product of its marginals if for any N finite, $P_N \sim \otimes^N P_1$.

Finally, we say that a measure P is *extremal* if it is locally equivalent to the product of its marginals and in addition it satisfies the following condition:

- (C5) Either $\rho_{\min} \leq 1$ and $\rho_{\max} \geq 1$, or if $\rho_{\min} > 1$ then for all $\delta > 0$, $P(\rho_0 < \rho_{\min} + \delta, \omega_0^0 > \omega_{\max}^0 - \delta) > 0$, or if $\rho_{\max} < 1$ then for all $\delta > 0$, $P(\rho_0 > \rho_{\max} - \delta, \omega_0^0 > \omega_{\max}^0 - \delta) > 0$.

Note that **(C5)**, which is used only in the proof of the annealed LDP, can be read off the support of P_0 and represents an assumption concerning the inclusion of “extremal environments” in the support of P . The introduction of this assumption is not essential and can be avoided at the cost of a slightly more cumbersome proof, see the remarks at the end of this chapter.

For a fixed $\varepsilon > 0$, we denote by $M_1^{e,\varepsilon}$ the set of probability measures satisfying Assumption 2.3.11 with parameter ε in **(B2)**. Define also the maps $F : \Omega \mapsto \Omega$ by $(F\omega)_k^+ = \omega_k^-$, $(F\omega)_k^- = \omega_k^+$, and $(\text{Inv } \omega)_k = (F\omega)_{-k}$. We now have:

Theorem 2.3.12 *Assume Assumption 2.3.11.*

a) *The random variables $\{T_n/n\}$ satisfy the weak quenched LDP with speed n and convex rate function*

$$I_P^{\tau,q}(u) = \sup_{\lambda \in \mathbb{R}} G(\lambda, P, u).$$

b) *Assume further that $E_P \log \rho_0 \leq 0$. Then, the random variables X_n/n satisfy the quenched LDP with speed n and good convex rate function*

$$I_P^q(v) = \begin{cases} v I_P^{\tau,q}\left(\frac{1}{v}\right) & , 0 < v \leq 1 \\ |v| \left(I_P^{\tau,q}\left(\frac{1}{|v|}\right) - E_P(\log \rho_0) \right) & , -1 \leq v < 0 \end{cases}$$

and

$$I_P^q(0) = \lim_{v \downarrow 0} v I_P^{\tau,q}\left(\frac{1}{v}\right).$$

c) *Finally, if $E_P \log \rho_0 > 0$, define $P^{\text{Inv}} := P \circ \text{Inv}^{-1}$. Then, $E_{P^{\text{Inv}}}(\log \rho_0) < 0$, and the LDP for (X_n/n) holds with good convex rate function*

$$I_P^q(v) = I_{P^{\text{Inv}}}^q(-v).$$

Proof. It should come as no surprise that we begin with the LDP for T_n/n . We divide the proof of Theorem 2.3.12 into the following steps:

Step I: $E_P \log \rho_0 \leq 0$, quenched LDP for T_n/n with convex rate function $I_P^{\tau,q}(\cdot)$:

$$(I.1) \quad \text{upper bound, lower tail: } P_\omega^o(T_n \leq nu)$$

$$(I.2) \quad \text{upper bound, upper tail: } P_\omega^o(T_n \geq nu)$$

$$(I.3) \quad \text{lower bound}$$

Step II: $E_P \log \rho_0 > 0$, quenched LDP for T_n/n with convex rate function $I_P^{\tau,q}(\cdot) + E_P(\log \rho_0)$,

Step III: quenched LDP for X_n/n with convex rate function $I_P^q(\cdot)$.

As a preliminary step we have the following technical lemma, whose proof is deferred:

Lemma 2.3.13 *Assume $P \in M_1^{e,\varepsilon}$ and $E_P(\log \rho_0) \leq 0$; Then*

(a) *The convex function $I_P^{r,q}(\cdot) : \mathbb{R} \mapsto [0, \infty]$ is nonincreasing on $[1, E_P(\overline{S})]$, nondecreasing on $[E_P(\overline{S}), \infty)$. Further, if $E_P(\overline{S}) < \infty$ then $I_P^{r,q}(E_P(\overline{S})) = 0$.*

(b) *For any $1 < u < E_P(\overline{S})$, there exists a unique $\lambda_0 = \lambda_0(u, P)$ such that $\lambda_0 < 0$ and*

$$u = \int \frac{d}{d\lambda} \log \varphi(\lambda, \omega) \Big|_{\lambda=\lambda_0} P(d\omega). \tag{2.3.14}$$

Further,

$$\inf_{P \in M_1^{e,\varepsilon}} \lambda_0(u, P) > -\infty. \tag{2.3.15}$$

(c) *There is a deterministic $\lambda_{\text{crit}} := \lambda_{\text{crit}}(P) \in [0, \infty]$ such that*

$$\varphi(\lambda, \omega) \begin{cases} < \infty, & \lambda < \lambda_{\text{crit}}, & P - \text{a.s.} \\ = \infty, & \lambda > \lambda_{\text{crit}}, & P - \text{a.s.} \end{cases}$$

with $\lambda_{\text{crit}} < \infty$ if $P(\omega_0^+ \omega_0^- = 0) = 0$. In the latter case, $E_\omega^o(e^{\lambda_{\text{crit}} \tau_1}) < e^{-\lambda_{\text{crit}}/\varepsilon}$, P -a.s., and with

$$u_{\text{crit}} = \begin{cases} \infty, & E_P \left[\frac{E_\omega^o(\tau_1 e^{\lambda_{\text{crit}} \tau_1})}{E_\omega^o(e^{\lambda_{\text{crit}} \tau_1})} \right] = \infty \\ E_P \left(\frac{d}{d\lambda} \log \varphi(\lambda, \omega) \Big|_{\lambda=\lambda_{\text{crit}}} \right), & E_P \left[\frac{E_\omega^o(\tau_1 e^{\lambda_{\text{crit}} \tau_1})}{E_\omega^o(e^{\lambda_{\text{crit}} \tau_1})} \right] < \infty, \end{cases}$$

and $E_P(\overline{S}) \leq u < u_{\text{crit}}$, there exists a unique $\lambda_0 := \lambda_0(u, P)$ such that $\lambda_0 \geq 0$ and (2.3.14) holds.

(d) *Assume P is extremal and further assume that $\rho_{\text{max}} < 1$. Then,*

$$\lambda_{\text{crit}} = \bar{\lambda} := -\log \left(\omega_{\text{max}}^0 + \frac{2(1 - \omega_{\text{max}}^0)\sqrt{\rho_{\text{max}}}}{1 + \rho_{\text{max}}} \right).$$

Further, define

$$\hat{\omega}^+ = (1 - \omega_{\text{max}}^0)/(1 + \rho_{\text{max}}), \hat{\omega}^- = \rho_{\text{max}}\hat{\omega}^+, \hat{\omega}^0 = \omega_{\text{max}}^0,$$

and let $\tilde{\omega}^{\text{min}}$ denote the deterministic environment with $\tilde{\omega}_k^{\text{min}} = \hat{\omega}$. Then, for any $\lambda \leq \lambda_{\text{crit}}$, and any ω such that $\omega_i^0 \leq \omega_{\text{max}}^0$, $\rho_i \leq \rho_{\text{max}}$,

$$\varphi(\lambda, \omega) \leq \varphi(\lambda, \tilde{\omega}^{\text{min}}) < \infty. \tag{2.3.16}$$

Step I.1: Obviously, it is enough to deal with $u \leq E_P(\overline{S})$. Indeed, for $u > E_P(\overline{S})$ we have by (2.1.30) that $\mathbb{P}^o(T_n \leq nu) \xrightarrow[n \rightarrow \infty]{} 1$, and there is nothing to prove. Next, by Chebycheff's inequality, for all $\lambda \leq 0$,

$$\begin{aligned}
 P_\omega^o\left(\frac{T_n}{n} \leq u\right) &\leq e^{-\lambda nu} E_\omega^o\left(e^{\lambda \sum_{i=1}^n \tau_i}\right) = e^{-\lambda nu} \prod_{i=1}^n E_{\theta^i \omega}^o\left(e^{\lambda \tau_1}\right) \\
 &= e^{-\lambda nu} \prod_{i=1}^n \varphi(\lambda, \theta^i \omega), \quad P - \text{ a.s.}
 \end{aligned}
 \tag{2.3.17}$$

where the first equality is due to the Markov property and the second due to $\tau_i < \infty$, \mathbb{P}^o - a.s. (the null set in (2.3.17) does not depend on λ). An application of the ergodic theorem yields that

$$\frac{1}{n} \log \prod_1^n \varphi(\lambda, \theta^i \omega) \longrightarrow E_P\left(f(\lambda, \omega)\right), \quad P - \text{ a.s.}$$

first for all λ rational and then for all λ by monotonicity. Thus,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(\frac{T_n}{n} \leq u\right) \leq -\sup_{\lambda \leq 0} G(\lambda, P, u), \quad P - \text{ a.s.}$$

Note that if $E_P(\bar{S}) = \infty$ then clearly $E_P[\log E_\omega^o(e^{\lambda \tau_1})] = \infty$ by Jensen’s inequality for $\lambda > 0$, and then $\sup_{\lambda \leq 0} G(\lambda, P, u) = I_P^{r,q}(u)$. If $E_P(\bar{S}) < \infty$ then, because $u < E_P(\bar{S})$, it holds that for any $\lambda > 0$,

$$\lambda u - E_P f(\lambda, \omega) \leq \lambda E_P(\bar{S}) - E_P f(\lambda, \omega) \leq 0,$$

where Jensen’s inequality was used in the last step. Since $G(0, P, u) = 0$, it follows that also in this case $\sup_{\lambda \leq 0} G(\lambda, P, u) = I_P^{r,q}(u)$. Hence,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(\frac{T_n}{n} \leq u\right) \leq -I_P^{r,q}(u) = -\inf_{w \leq u} I_P^{r,q}(w),$$

where the last inequality is due to part a) of Lemma 2.3.13, completing Step I.1.

Step I.2: is similar, using this time $\lambda \geq 0$.

Step I.3: The proof of the lower bound is based on a change of measure argument. We present it here in full detail for $u < u_{\text{crit}}$. Fix $\lambda_0 = \lambda_0(u, P)$ as in Lemma 2.3.13, and set a probability measure $Q_{\omega,n}^o$ such that

$$\frac{dQ_{\omega,n}^o}{dP_\omega^o} = \frac{1}{Z_{n,\omega}} \exp\left(\lambda_0 T_n\right), \quad Z_{n,\omega} = \mathbb{E}_\omega^o\left(\exp\left(\lambda_0 T_n\right)\right),$$

and let $\bar{Q}_{\omega,n}^o$ denote the induced law on $\{\tau_1, \dots, \tau_n\}$. Due to the Markov property, $\bar{Q}_{\omega,n}^o$ is a product measure, whose first n marginals do not depend on n , hence we will write \bar{Q}_ω^o instead of $\bar{Q}_{\omega,n}^o$ when integrating over events depending only on $\{\tau_i\}_{i < n}$. But, for any $\delta > 0$,

$$\begin{aligned}
 &P_\omega^o\left(\frac{T_n}{n} \in (u - \delta, u + \delta)\right) \\
 &\geq \exp\left(-nu\lambda_0 - n\delta|\lambda_0| + \sum_{i=1}^n \log \varphi\left(\lambda_0, \theta^i \omega\right)\right) \overline{Q}_\omega^o\left(\left|\frac{T_n}{n} - u\right| \leq \delta\right).
 \end{aligned} \tag{2.3.18}$$

By the ergodic theorem and the fact that $u < u_{\text{crit}}$, it holds that

$$E_{\overline{Q}_\omega^o}(T_n/n) \xrightarrow{n \rightarrow \infty} E_P\left(E_{\overline{Q}_\omega^o}(\tau_1)\right) = u, P - \text{ a.s.} \tag{2.3.19}$$

where we used again (2.3.14). On the other hand, again because $\lambda_0 < \lambda_{\text{crit}}$ it holds that there exists an $\eta > 0$ such that

$$E_P\left(E_{\overline{Q}_\omega^o}\left(e^{\eta\tau_1}\right)\right) < \infty,$$

implying that

$$\overline{Q}_\omega^o\left(\left|\frac{T_n}{n} - u\right| \geq \delta\right) \xrightarrow{n \rightarrow \infty} 0, P - \text{ a.s.} \tag{2.3.20}$$

Combining (2.3.20) with (2.3.18), we get

$$\begin{aligned}
 &\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(\frac{T_n}{n} \in (u - \delta, u + \delta)\right) \\
 &\geq -u\lambda_0 - \delta|\lambda_0| + \liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \log \varphi(\lambda_0, \theta^i \omega) \\
 &= -u\lambda_0 - \delta|\lambda_0| + E_P(\log \varphi(\lambda_0, \omega)) \\
 &= -G(\lambda_0, P, u) - \delta|\lambda_0| = -I_P^{T,q}(u) - \delta|\lambda_0|, P - \text{ a.s.}
 \end{aligned}$$

where the first equality is due to the ergodic theorem and the last one to Lemma 2.3.13. This completes Step I.3 when $u < u_{\text{crit}}$, since $\delta > 0$ is arbitrary. For $u > u_{\text{crit}}$, the proof is similar, except that one needs to truncate the variables $\{\tau_i\}$, we refer to [12, Theorem 4] for details. Step I is complete, except for the:

Proof of Lemma 2.3.13

We consider in what follows only the case $P(\omega_0^+ \omega_0^- = 0) = 0$, the modifications in the case where random nodes are allowed are left to the reader.

a) The convexity of $I_P^{T,q}(\cdot)$ is immediate from its definition as a supremum of affine functions.

As in the course of the proof of Step I, recall that

$$\sup_{\lambda \in \mathbb{R}} G(\lambda, u, P) = \begin{cases} \sup_{\lambda \leq 0} G(\lambda, u, P), & u < E_P(\overline{S}) \\ \sup_{\lambda \geq 0} G(\lambda, u, P), & u > E_P(\overline{S}) \\ 0, & u = E_P(\overline{S}). \end{cases}$$

The stated monotonicity properties are then immediate.

b)+c) Recall the path decomposition (2.1.13). Exponentiating and taking expectations using $\tau_1 < \infty$, \mathbb{P}^o - a.s., we have that if $\varphi(\lambda, \omega) < \infty$ then

$$\varphi(\lambda, \omega) = \omega_0^+ e^\lambda + \omega_0^0 e^\lambda \varphi(\lambda, \omega) + \omega_0^- e^\lambda \varphi(\lambda, \omega) \varphi(\lambda, \theta^{-1}\omega). \quad (2.3.21)$$

Thus $\varphi(\lambda, \omega) < \infty$ implies $\varphi(\lambda, \theta^{-1}\omega) < \infty$, yielding that $\mathbf{1}_{\varphi(\lambda, \omega) < \infty}$ is constant P - a.s., and hence for all λ rational, $P(\varphi(\lambda, \omega) < \infty) \in \{0, 1\}$. This, and the monotonicity of $\varphi(\lambda, \omega)$ in λ , immediately yields the existence of a deterministic λ_{crit} . (We note in passing that (2.3.21) gives, by iterating, a representation of $\varphi(\lambda, \omega)$ as a continued fraction, but we do not need this now.) We also conclude from (2.3.21) that for $\lambda < \lambda_{\text{crit}}$ it holds that $\varphi(\lambda, \omega) \leq e^{-\lambda}/\varepsilon$, P -a.s., which implies by monotone convergence that $\varphi(\lambda_{\text{crit}}, \omega) < \infty$, P -a.s.

Next, for $\lambda < 0$ we have that

$$g(\lambda) := \int \frac{E_\omega^o(\tau_1 e^{\lambda\tau_1})}{E_\omega^o(e^{\lambda\tau_1})} P(d\omega) = \int \frac{d}{d\lambda} \log \varphi(\lambda, \omega) P(d\omega).$$

Further, $g(0) = E_{\mathbb{P}^o}(\tau_1)$, whereas $g(\cdot) \geq 1$ is strictly monotone increasing, satisfying $g(\lambda) \xrightarrow{\lambda \rightarrow -\infty} 1$. This implies (2.3.14). Finally, to see (2.3.15), note that

$$\begin{aligned} 1 \leq \frac{E_\omega^o(\tau_1 e^{\lambda\tau_1})}{E_\omega^o(e^{\lambda\tau_1})} &\leq \frac{\omega_0^+ e^\lambda + E_\omega(\tau_1 e^{\lambda\tau_1} \mathbf{1}_{\tau_1 \geq 2})}{\omega_0^+ e^\lambda} \\ &\leq 1 + \frac{c e^{3\lambda/2}}{\omega_0^+ e^\lambda} \leq 1 + \frac{c}{\varepsilon} e^{\lambda/2} \xrightarrow{\lambda \rightarrow -\infty} 1, \end{aligned}$$

where c is a constant independent of ω or λ . Hence, $g(\lambda) \xrightarrow{\lambda \rightarrow -\infty} 1$ uniformly in $M_1^{e, \varepsilon}$.

d) Assume that P is extremal. The first inequality in (2.3.16) follows from a simple coupling argument: let $\bar{\varphi}(\lambda) := E_{\bar{\omega}_{\min}}^o[e^{\lambda\tau_1}]$. By the recursions (2.3.21), it holds that if $\bar{\varphi}(\lambda) < \infty$ then as long as $\lambda \leq \bar{\lambda}$ it holds that

$$\bar{\varphi}(\lambda) = \frac{(1 - \omega_{\max}^0 e^\lambda) - \sqrt{(1 - \omega_{\max}^0 e^\lambda)^2 - 4\hat{\omega}^+ \hat{\omega}^- e^{2\lambda}}}{2\hat{\omega}^- e^\lambda}.$$

Thus, we have to show that if $\lambda > \lambda_{\text{crit}}$ then $E_\omega^o(e^{\lambda\tau_1}) = \infty$, P -a.s. Since $E_{\bar{\omega}_{\min}}^o(e^{\lambda\tau_1}) = \infty$, we may find an M large enough such that $E_{\bar{\omega}_{\min}}^o(e^{\lambda\tau_1} \mathbf{1}_{\tau_1 < M}) > 1/\varepsilon + 1$. Since the last expression is local, i.e. depends only on $\{\omega_i\}_{i=0}^{-M+1}$, it follows (from the assumption of local equivalence to the product of marginals) that with P positive probability, $E_\omega^o(e^{\lambda\tau_1}) > 1/\varepsilon$, and hence by part (c) actually $E_\omega^o(e^{\lambda\tau_1}) = \infty$ with P positive probability, and hence with P probability 1. \square

Remark: Before proceeding, we note that a direct consequence of Lemma 2.3.13 is that if $E_P(\log \rho_0) \leq 0$, then for $u < E_P(\bar{S})$,

$$I_P^{\tau,q}(u) = \lambda_0 u - E_P(f(\lambda_0, \omega)) = \sup_{\lambda \in \mathbb{R}} G(\lambda, P, u) > G(0, P, u) = 0$$

since the function $G(\cdot, P, u)$ is strictly concave.

Step II: Recall the transformation $\text{Inv} : \Omega \mapsto \Omega$ and the law $P^{\text{Inv}} = P \circ \text{Inv}^{-1}$. Proving the LDP for T_n/n when $E_P(\log \rho_0) > 0$ is the same, by space reversal, as proving the quenched LDP for T_{-n}/n under the law P^{Inv} on the environment. Note that in this case, $E_{P^{\text{Inv}}}(\log \rho_0) < 0$, and further, $P \in M_1^{e,\varepsilon}$ implies that $P^{\text{Inv}} \in M_1^{e,\varepsilon}$. Thus, Step II will be completed if we can prove a quenched LDP for T_{-n}/n for $P \in M_1^{e,\varepsilon}$ satisfying $E_P \log \rho_0 < 0$. We turn to this task now.

Note that if $P(\omega_0^- = 0) > 0$ then $P_\omega^o(T_{-n} < \infty) = 0$ for some $n = n(\omega)$ large enough, and the LDP for T_{-n}/n is trivial. We thus assume throughout that $\omega_0^- \geq \varepsilon$, P -a.s. As a first step in the derivation of the LDP, we compute logarithmic moment generating functions. Define, for any $\lambda \in \mathbb{R}$,

$$\bar{\varphi}(\lambda, \omega) = E_\omega^o(e^{\lambda\tau_{-1}} \mathbf{1}_{\{\tau_{-1} < \infty\}}), \quad \bar{f}(\lambda, \omega) = \log \bar{\varphi}(\lambda, \omega).$$

Lemma 2.3.22 *Assume $P \in M_1^{e,\varepsilon}$ and further assume that $\min(\omega_0^+, \omega_0^-) > \varepsilon$, P -a.s. Then,*

$$E_P(\bar{f}(\lambda, \omega)) = E_P(f(\lambda, \omega)) + E_P \log \rho_0. \tag{2.3.23}$$

Proof of Lemma 2.3.22:

Define the map $I_n : \Omega \mapsto \Omega$ by

$$(I_n \omega)_k = \begin{cases} \omega_k, & k \notin [0, n] \\ (F\omega)_{n-k}, & k \in [0, n]. \end{cases}$$

Introduce

$$\varphi_n(\lambda, \omega) = E_\omega^o(e^{\lambda\tau_1}; \tau_1 < T_{-(n+1)}), \quad \bar{\varphi}_n(\lambda, \omega) = E_\omega^o(e^{\lambda\tau_{-1}}; \tau_{-1} < T_{(n+1)}).$$

We will show below that

$$G_n(\lambda, \omega) := \varphi_n(\lambda, \theta^n \omega) \bar{\varphi}_{n-1}(\lambda, \omega) = \varphi_{n-1}(\lambda, \theta^n \omega) \bar{\varphi}_n(\lambda, \omega) := F_n(\lambda, \omega). \tag{2.3.24}$$

Because $\min(\omega_0^+, \omega_0^-) > \varepsilon$, P -a.s., the function $\log \varphi_n(\lambda, \omega)$ and $\log \bar{\varphi}_n(\lambda, \omega)$ are P -integrable for each n . Taking logarithms in (2.3.24), we find that $E_P(\log \varphi_n(\lambda, \omega)) - E_P(\log \bar{\varphi}_n(\lambda, \omega))$ does not depend on n . On the other hand, both terms are monotone in n , hence by monotone convergence either both sides of (2.3.24) are $+\infty$ or both are finite, in which case

$$E_P(\log \varphi(\lambda, \omega)) - E_P(\log \bar{\varphi}(\lambda, \omega)) = E_P \left(\log \left(\frac{\varphi_0(\lambda, \omega)}{\bar{\varphi}_0(\lambda, \omega)} \right) \right) = -E_P(\log \rho_0),$$

yielding (2.3.23).

We thus turn to the proof of (2.3.24). It is straight forward to check, by space inversion, that $F_n(\lambda, I_n\omega) = G_n(\lambda, \omega)$. Thus, the proof of (2.3.24) will be complete once we show that $F_n(\lambda, I_n\omega) = F_n(\lambda, \omega)$. Toward this end, note that by the Markov property,

$$\begin{aligned} \bar{\varphi}_n(\lambda, \omega) &= E_\omega^o(e^{\lambda T_{-1}}; T_{-1} < T_{n+1}) \\ &= E_\omega^o(e^{\lambda T_{-1}}; T_{-1} < T_n) \\ &\quad + E_\omega^o(e^{\lambda T_n}; T_n < T_{-1})E_\omega^n(e^{\lambda T_{-1}}; T_{-1} < T_{n+1}). \end{aligned}$$

Hence, defining

$$\begin{aligned} B_n(\lambda, \omega) &:= E_\omega^o(e^{\lambda T_{-1}}; T_{-1} < T_n), \\ C_n(\lambda, \omega) &:= E_\omega^o(e^{\lambda T_n}; T_n < T_{-1}), \end{aligned}$$

one has, using again space reversal and the Markov property in the second equality,

$$\begin{aligned} F_n(\lambda, \omega) &= E_\omega^n(e^{\lambda T_{n+1}}; T_{n+1} < T_0)E_\omega^o(e^{\lambda T_{-1}}; T_{-1} < T_n) \\ &\quad + E_\omega^n(e^{\lambda T_{n+1}}; T_{n+1} < T_0)E_\omega^n(e^{\lambda T_{-1}}; T_{-1} < T_{n+1})E_\omega^o(e^{\lambda T_n}; T_n < T_{-1}) \\ &\quad = B_n(\lambda, \omega)B_n(\lambda, I_n\omega) \\ &\quad + E_\omega^n(e^{\lambda T_{n+1}}; T_{n+1} < T_0)E_\omega^n(e^{\lambda T_0}; T_0 < T_{n+1})E_\omega^o(e^{\lambda T_{-1}}; T_{-1} < T_{n+1}) \\ &\quad \quad E_\omega^o(e^{\lambda T_n}; T_n < T_{-1}) \\ &\quad = B_n(\lambda, \omega)B_n(\lambda, I_n\omega) + C_n(\lambda, \omega)C_n(\lambda, I_n\omega)F_n(\lambda, \omega), \end{aligned}$$

implying the invariance of $F_n(\lambda, \omega)$ under the action of I_n on Ω , except possibly at λ where $C_n(\lambda, \omega)C_n(\lambda, I_n\omega) = 1$. The latter λ is then handled by continuity. This completes the proof of Lemma 2.3.22 \square

Step II now is completed by following the same route as in the proof of Step I, using Lemma 2.3.22 to transfer the analytic results of Lemma 2.3.13 to this setup. The details, which are straightforward and are given in [12], are omitted here. \square

Remarks: 1. Note that the conclusion of Lemma 2.3.22 extends immediately, by the ergodic decomposition, to stationary measures $P \in M_1^{s, \varepsilon}$.

2. Lemma 2.3.22 is the key to the large deviations principle, and deserves some discussion. First, by taking $\lambda \uparrow 0$, one sees that if $E_P(\log \rho_0) \leq 0$ then $E_P[\log P_\omega^o(\tau_{-1} < \infty)] = E_P(\log \rho_0)$. Next, let $\bar{\tau}_{-1}, \bar{\tau}_{-2}, \bar{\tau}_{-3}, \dots, \bar{\tau}_{-N}$ have the distribution of $\tau_{-1}, \tau_{-2}, \tau_{-3}, \dots, \tau_{-N}$ under P_ω^o conditioned on $T_{-N} < \infty$. In fact the law of $\{\bar{\tau}_{-i}\}_{i=1}^N$ does not depend on N . This can be seen by a discrete h -transform: the distributions of $X_0^{T_{-N}} := (X_0, \dots, X_{T_{-N}})$ under P_ω^o , conditioned on $T_{-N} < \infty$, $N = 1, 2, \dots$ form a consistent family whose extension is again a Markov chain. To see this, let $\tilde{P}_{\omega, N}^o := P_\omega^o(\cdot | T_{-N} < \infty)$, restricted to $X_0^{T_{-N}}$. Denoting $x_1^n := (x_1, \dots, x_n)$, compute (with $x_i > -N$),

$$\begin{aligned}
 \tilde{P}_{\omega,N}^o(X_{n+1} = x_n + 1 | X_1^n = x_1^n) &= \frac{\tilde{P}_{\omega,N}^o(X_{n+1} = x_n + 1, X_1^n = x_1^n)}{\tilde{P}_{\omega,N}^o(X_1^n = x_1^n)} \\
 &= \frac{P_{\omega}^o(X_{n+1} = x_n + 1, X_1^n = x_1^n, T_{-N} < \infty)}{P_{\omega}^o(X_1^n = x_1^n, T_{-N} < \infty)} \\
 &= \frac{P_{\omega}^o(X_{n+1} = x_n + 1, X_1^n = x_1^n) P_{\theta^{x_n+1}\omega}^o(T_{-N-x_n-1} < \infty)}{P_{\omega}^o(X_1^n = x_1^n) P_{\theta^{x_n}\omega}^o(T_{-N-x_n} < \infty)} \\
 &= P_{\omega}^o(X_{n+1} = x_n + 1 | X_1^n = x_1^n) P_{\theta^{x_n+1}\omega}(T_{-1} < \infty) \\
 &= \omega_{x_n}^+ P_{\theta^{x_n+1}\omega}^o(T_{-1} < \infty),
 \end{aligned}$$

where we used the Markov property in the third and in the fourth equality. The last term depends neither on N nor on x_1^{n-1} . Therefore, the extension of $(\tilde{P}_{\omega,N})_{N \geq 1}$ is the distribution of the Markov chain with transition probabilities $\tilde{\omega}_i^+ = \omega_i^+ P_{\theta^{i+1}\omega}(T_{-1} < \infty)$, $\tilde{\omega}_i^0 = \omega_i^0$, $i \in \mathbb{Z}$. In particular, $\bar{\tau}_{-1}, \bar{\tau}_{-2}, \bar{\tau}_{-3}, \dots$ are independent under P_{ω}^o and, with a slight abuse of notations, form a stationary sequence under \mathbb{P}^o . Note now that if we set

$$\bar{\phi}(\lambda, \omega) := E_{\omega}^o(e^{\lambda \bar{\tau}_{-1}}) = \frac{\bar{\varphi}(\lambda, \omega)}{P_{\omega}^o(T_{-1} < \infty)} \tag{2.3.25}$$

then Lemma 2.3.22 tells us that $E_P \bar{\phi}(\lambda, \omega) = E_P \varphi(\lambda, \omega)$. In particular, $\mathbb{E}_P(\bar{\tau}_1) = \mathbb{E}_P(\tau_1) = E_P(\bar{S})$ if $E_P(\log \rho_0) \leq 0$ and, repeating the arguments leading to the LDP of T_n/n , we find that the sequence of random variables T_{-n}/n , *conditioned on* $T_{-n} < \infty$, satisfy a quenched LDP under P_{ω}^o with the same rate function as T_n/n !

Step III: By space reversal, it is enough to prove the result for $E_P(\log \rho_0) \leq 0$. Further, as in Step II, it will be enough to consider the case where $\min(\omega_0^+, \omega_0^-) \geq \varepsilon$, P -a.s. Since $I_P^{r,q}(\cdot)$ is convex, and since $x \mapsto xf(1/x)$ is convex if $f(\cdot)$ is convex, it follows that $I_P^q(\cdot)$ is convex on $(0, 1]$ and on $[-1, 0)$ separately. If $\lambda_{\text{crit}}(P) = 0$ then $I_P^q(0) = 0$ and the convexity on $[-1, 1]$ follows. In the general case, note that I_P^q is continuous at 0, and $(I_P^q)'(0^-) = -(I_P^q)'(0^+) + E_P(\log \rho_0)$. Note that for $\lambda \leq \lambda_{\text{crit}}$, by the Markov property,

$$E_{\omega}^o(e^{\lambda T_M} \mathbf{1}_{\tau_{-1} < \tau_M}) = E_{\omega}^o(e^{\lambda \tau_{-1}} \mathbf{1}_{\tau_{-1} < \tau_M}) \varphi(\lambda, \theta^{-1}\omega) E_{\omega}^o(e^{\lambda T_M}),$$

and hence,

$$1 \geq E_{\omega}^o(e^{\lambda \tau_{-1}} \mathbf{1}_{\tau_{-1} < \tau_M}) \varphi(\lambda, \theta^{-1}\omega),$$

leading (by taking $M \rightarrow \infty$) to the conclusion that

$$\varphi(\lambda, \omega) \bar{\varphi}(\lambda, \theta^{-1}\omega) \leq 1.$$

Taking logarithms and P -expectations, we conclude that

$$E_P f(\lambda, \omega) + E_P \bar{f}(\lambda, \omega) \leq 0.$$

Combined with Lemma 2.3.22, we deduce that for all $\lambda \leq \lambda_{\text{crit}}$, $2E_P(f(\lambda, \omega)) \leq -E_P(\log \rho_0)$. Hence,

$$(I_P^q)'(0^+) = -E_P(\log E_\omega(e^{\lambda_{\text{crit}}\tau_1})) = -E_P(f(\lambda_{\text{crit}}, \omega)) \geq \frac{1}{2}E_P(\log \rho_0),$$

implying that $(I_P^q)'(0^-) \leq (I_P^q)'(0^+)$, and hence that $I_P^q(\cdot)$ is convex on $[-1, 1]$.

Using the monotonicity of $I_P^{\tau, q}(\cdot)$, c.f. the remark following the proof of Step I, it follows easily that $I_P^q(\cdot)$ is non increasing on $[-1, v_P]$ and non decreasing on $[v_P, 1]$.

Let $v > v_P$. We have

$$P_\omega^\circ\left(\frac{X_n}{n} \geq v\right) \leq P_\omega^\circ(T_{\lfloor nv \rfloor} \leq n) = P_\omega^\circ\left(\frac{T_{\lfloor nv \rfloor}}{\lfloor nv \rfloor} \leq \frac{n}{\lfloor nv \rfloor}\right).$$

Step I and the monotonicity of $I_P^{\tau, q}(\cdot)$ now imply

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^\circ\left(\frac{X_n}{n} \geq v\right) \leq -vI_P^{\tau, q}\left(\frac{1}{v}\right),$$

which yields the required upper bound by the monotonicity of $I_P^q(\cdot)$. The same argument applies to yield the desired upper bound on $P_\omega^\circ\left(\frac{X_n}{n} \leq v\right)$ for $v < 0$, by considering the hitting times $T_{\lceil nv \rceil}$.

In the same way, for any $0 < \eta < \delta/2$,

$$P_\omega^\circ\left((v + \delta) \geq \frac{X_n}{n} \geq (v - \delta)\right) \geq P_\omega^\circ\left((1 - \eta)n \leq T_{\lceil nv \rceil} \leq n\right),$$

hence, from Step I it follows that for $v \geq 0$,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^\circ\left(\frac{X_n}{n} \in (v - \delta, v + \delta)\right) \geq -vI_P^{\tau, q}\left(\frac{1 - \eta}{v}\right), \quad P - \text{a.s.},$$

and the lower bound is obtained by letting $\eta \rightarrow 0$. The same argument also yields the lower bounds for $v < 0$, using this time the function $I_P^{\tau, q}(\cdot)$.

Next, we turn to evaluate an upper bound on $P_\omega^\circ(X_n/n \leq v)$, $0 \leq v < v_P$, with $v_P \geq 0$. Starting with $v = 0$, let $\eta, \delta > 0$, with $\delta < v_P$. Then,

$$\begin{aligned} P_\omega^\circ(X_n \leq 0) &\leq P_\omega^\circ(T_{\lceil n\delta \rceil} \geq n) + P_\omega^\circ(T_{\lceil n\delta \rceil} < n, \frac{X_n}{n} \leq 0) && (2.3.26) \\ &\leq P_\omega^\circ(T_{\lceil n\delta \rceil} \geq n) + \sum_{1/\eta \leq k, l; (k+l)\eta \leq 1/\delta} P_\omega^\circ\left(\frac{T_{\lceil n\delta \rceil}}{n\delta} \in [k\eta, (k+1)\eta)\right) \times \\ &\quad P_{\theta_{\lceil n\delta \rceil}^\omega}\left(\frac{T_{\lceil n\delta \rceil} - \lceil n\delta \rceil}{n\delta} \in [l\eta, (l+1)\eta)\right) \sum_{-2n\delta\eta \leq m - n(1 - (k+l)\delta\eta) \leq 0} P_\omega^\circ(X_m \leq 0), \end{aligned}$$

by the strong Markov property. Define the random variable

$$a = \limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{m: -2n\delta\eta \leq m-n \leq 0} \log P_\omega^o(X_m \leq 0) ,$$

and note, using the inequality

$$P_\omega^o(X_n \leq 0) \geq P_\omega^o(X_m \leq 0) \inf_{i \leq 0} P_{\theta^i \omega}^o[X_{n-m} = -(n-m)]$$

with a worst-environment estimate, that

$$a - C\delta\eta \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o[X_n \leq 0] \leq a \tag{2.3.27}$$

with $C = -2 \log \varepsilon > 0$. The first two probabilities in the right-hand side of (2.3.26) will be estimated using Step I. By convexity, the rate functions $I_P^{\tau,q}$ and $I_P^{-\tau,q} := I_P^{\tau,q} - E_P(\log \rho_0)$ are continuous, so that the oscillation

$$w(\delta; \eta) = \max\{|I_P^{\tau,q}(u) - I_P^{\tau,q}(u')| + |I_P^{-\tau,q}(u) - I_P^{-\tau,q}(u')|; u, u' \in [1, 1/\delta], |u - u'| \leq \eta\}$$

tends to 0 with η , for all fixed δ . From the proof of Step II, it is not difficult to see that the third term in the right-hand side of (2.3.26) can be estimated similarly (it does not cause problems to consider $P_{\theta^{[n\delta]}\omega}^o$ instead of P_ω^o):

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_{\theta^{[n\delta]}\omega}^o\left(\frac{T_{-[n\delta]}}{n\delta} \in [l\eta, (l+1)\eta]\right) \leq -\delta (I_P^{-\tau,q}(l\eta) - w(\delta; \eta)) \quad P\text{-a.s.}$$

Finally, we get from (2.3.27) and (2.3.26)

$$a \leq C\delta\eta + \max\{-I_P^q(\delta), \max_{1/\eta \leq k, l; (k+l)\eta \leq 1/\delta} [-\delta\eta(kI_P^q(1/k\eta) + lI_P^q(-1/l\eta)) + 2\delta w(\delta; \eta) + (1 - (k+l+2)\delta\eta)a']\} .$$

By convexity and since $\delta \leq v_P$, it holds $kI_P^q(1/k\eta) + lI_P^q(-1/l\eta) \geq (k+l)I_P^q(0) \geq (k+l)I_P^q(\delta)$, and therefore $a' := a + I_P^q(\delta)$ is such that

$$a' \leq C\delta\eta + \left(\max_{1/\eta \leq k, l; (k+l)\eta \leq 1/\delta} [2\delta w(\delta; \eta) + 2\delta\eta I_P^q(\delta) + (1 - (k+l+2)\delta\eta)a'] \right)^+ .$$

Computing the maximum for positive a' , we derive that $2a' \leq C\eta + 2(w(\delta; \eta) + \eta I_P^q(\delta))$. Letting now $\eta \rightarrow 0$ and $\delta \rightarrow 0$, we conclude that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(X_n \leq 0) \leq -I_P^q(0) , \quad P\text{-a.s.} \tag{2.3.28}$$

In fact, the same proof actually shows that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(\exists \ell \geq n : X_\ell \leq 0) \leq -I_P^q(0) , \quad P\text{-a.s.} \tag{2.3.29}$$

For an arbitrary $v \in [0, v_P)$, we write

$$\begin{aligned}
 P_\omega^o\left(\frac{X_n}{n} \leq v\right) &\leq P_\omega^o(\exists \ell \geq n : X_\ell \leq nv) \leq P_\omega^o(T_{[nv]} \geq n) \\
 + \sum_{k: v/\eta \leq k < 1/\eta} P_\omega^o\left(\frac{T_{[nv]}}{n} \in [k\epsilon, (k+1)\epsilon)\right) &P_{\theta^{[nv]}\omega}^o\left(\exists \ell \geq n - n(k+1)\eta : X_\ell \leq 0\right)
 \end{aligned} \tag{2.3.30}$$

where the two first probabilities in the right-hand side can be estimated using Step I, and concerning the last one we note that from (2.3.29) one has that

$$\frac{1}{n - n(k+1)\eta} \log P_{\theta^{[nv]}\omega}^o\left(\exists \ell \geq n - n(k+1)\eta : X_\ell \leq 0\right) \xrightarrow{n \rightarrow \infty} -I_P^q(0),$$

in probability, and hence a.s. along a random subsequence. Therefore,

$$\begin{aligned}
 &\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(\exists \ell \geq n : X_\ell \leq nv\right) \\
 &\leq \limsup_{\eta \rightarrow 0} \left(-I_P^q(v) \vee \max_{v/\eta \leq k \leq 1/\eta} [-k\eta I_P^q(v/k\eta) - (1 - k\eta)I_P^q(0)] \right) \\
 &= -I_P^q(v),
 \end{aligned} \tag{2.3.31}$$

by convexity. But, due to Kingman’s sub-additive ergodic theorem, the left hand side of the last expression converges P -a.s., resulting with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(X_n \leq nv\right) \leq -I_P^q(v), \quad P - a.s..$$

The upper bound for general subsets of $[0, 1]$ follows by noting the convexity of $I_P^q(\cdot)$. □

Remarks 1. If P is extremal, a simpler proof of (2.3.28) can be given. Indeed, note that $I_P^q(0) = \lambda_{\text{crit}}(0)$, and, by extremality,

$$\begin{aligned}
 P_\omega^o(X_n \leq 0) &\leq P_\omega(X_m \leq 0, \text{ some } m \geq 0) \leq P_{\tilde{\omega}^{\min}}(X_m \leq 0, \text{ some } m \geq 0) \\
 &\leq \sum_{m=n}^\infty P_{\tilde{\omega}^{\min}}(X_m \leq 0).
 \end{aligned}$$

A simple computation reveals that $P_{\tilde{\omega}^{\min}}(X_m \leq 0) \leq C_\lambda e^{-\lambda m}$ for any $\lambda < \lambda_{\text{crit}}$, yielding that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(X_n \leq 0) \leq -I_P^q(0), \quad P - a.s. \tag{2.3.32}$$

2. A lot of information is available concerning the shape of the rate function $I_P^q(\cdot)$, and in particular concerning the existence of pieces where the rate function is not strictly convex. We refer to the discussion in [12] for details.

Annealed LDP’s

The LDP in the annealed setting also makes use of the hitting times T_n and T_{-n} . For technical reasons, we need to make stronger hypotheses on the environment. To state these, define the empirical process

$$R_n(\omega) := \frac{1}{n} \sum_{i=0}^{n-1} \delta_{\theta^i \omega} .$$

R_n takes values in the space $M_1(\Omega)$ of probability measures on Ω , which we equip with the topology of weak convergence. We also need to introduce the specific relative entropy

$$h(\cdot|P) : M_1(\Omega) \mapsto [0, \infty], \quad h(Q|P) := \begin{cases} \lim_{N \rightarrow \infty} \frac{1}{N} H(Q_N|P_N), & Q \text{ stationary} \\ \infty, & \text{otherwise,} \end{cases}$$

where Q_N, P_N denote the restriction of Q, P to the first N coordinates $\{\omega_i\}_{i=0}^{N-1}$ and $H(\cdot|\cdot)$ denotes the relative entropy:

$$H(\mu|\nu) = \begin{cases} \int \log \left(\frac{d\mu}{d\nu}(x) \right) \mu(dx), & \mu \ll \nu \\ \infty, & \text{otherwise} \end{cases} .$$

Assumption 2.3.33

- (C1) P is stationary and ergodic
- (C2) There exists an $\varepsilon > 0$ such that $\min(\omega_0^+, \omega_0^-) > \varepsilon$, P - a.s.,
- (C3) $\{R_n\}$ satisfies under P the process level LDP in $M_1(\Omega)$ with good rate function $h(\cdot|P)$,
- (C4) P is locally equivalent to the product of its marginals and, for any stationary measure $\eta \in M_1(\Omega)$ there is a sequence $\{\eta^n\}$ of stationary, ergodic measures with $\eta^n \xrightarrow{n \rightarrow \infty} \eta$ weakly and $h(\eta^n|P) \rightarrow h(\eta|P)$.
- (C5) P is extremal.

We note that product measures and Markov processes with bounded transition kernels satisfy (C1)–(C4) of Assumption 2.3.33, see [27, Lemma 4.8] and [23]. Define now

$$I_P^{\tau, \alpha}(u) = \inf_{\eta \in M_1^{\varepsilon, \varepsilon}(\Omega)} [I_\eta^{\tau, q}(u) + h(\eta|P)] , \quad I_P^\alpha(v) = \inf_{\eta \in M_1^{\varepsilon, \varepsilon}(\Omega)} [I_\eta^q(v) + |v|h(\eta|P)] .$$

We now have the annealed analog of Theorem 2.3.12:

Theorem 2.3.34 *Assume Assumption 2.3.33. Then, the random variables $\{T_n/n\}$ satisfy the weak annealed LDP with speed n and rate function $I_P^{\tau, \alpha}(\cdot)$. Further, the random variables X_n/n satisfy the annealed LDP with speed n and good rate function $I_P^\alpha(\cdot)$.*

Proof. Throughout, $M_1^{s,\varepsilon}$ denotes the set of stationary probability measures $\eta \in M_1(\Omega)$ satisfying $\text{supp } \eta_0 \subset \text{supp } P_0$. If $E_P(\log \rho_0) \leq 0$ then $\lambda_{\text{crit}} = \lambda_{\text{crit}}(P)$ is as in Lemma 2.3.13, whereas if $E_P(\log \rho_0) > 0$ then $\lambda_{\text{crit}} = \lambda_{\text{crit}}(P^{\text{Inv}})$.

Let $M_1^{s,\varepsilon,P} = \{\mu \in M_1^{s,\varepsilon} : \text{supp } \mu_0 \subset \text{supp } P_0\}$. The following lemma, whose proof is deferred, is key to the transfer of quenched LDP's to annealed LDP's:

Lemma 2.3.35 *Assume P satisfies Assumption 2.3.33. Then, the function $(\mu, \lambda) \mapsto \int f(\lambda, \omega)\mu(d\omega)$ is continuous on $M_1^{s,\varepsilon,P} \times (-\infty, \lambda_{\text{crit}}]$.*

Steps I.1 + I.2: weak annealed LDP upper bound for T_n/n : We have, for $\lambda \leq 0$,

$$\begin{aligned} \mathbb{P}^o(T_n/n \leq u) &\leq e^{-\lambda nu} \mathbb{E}^o \left(\exp \left(\lambda \sum_{j=1}^n \tau_j \right) \mathbf{1}_{\tau_j < \infty, j=1, \dots, n} \right) \\ &= e^{-\lambda nu} E_P \left(\prod_{j=1}^n E_\omega^o \left(e^{\lambda \tau_j} \mathbf{1}_{\tau_j < \infty} \right) \right) = e^{-\lambda nu} E_P \left(\exp \left(\sum_{j=0}^{n-1} f(\lambda, \theta^j \omega) \right) \right) \\ &= e^{-\lambda nu} E_P \left(\exp \left(n \int f(\lambda, \omega) R_n(d\omega) \right) \right). \end{aligned} \tag{2.3.36}$$

By Assumption 2.3.33, $\{R_n\}$ satisfies a LDP with rate function $h(\cdot|P)$. Lemma 2.3.35 ensures that we can apply Varadhan's lemma (see [19, Lemma 4.3.6]) to get

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log E_P \left(\exp \left(n \int f(\lambda, \omega) R_n(d\omega) \right) \right) \\ \leq \sup_{\eta \in M_1^{s,\varepsilon}} \left[\int f(\lambda, \omega) \eta(d\omega) - h(\eta|P) \right]. \end{aligned} \tag{2.3.37}$$

Going back to (2.3.36), this yields the upper bound

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(T_n/n \leq u) &\leq \inf_{\lambda \leq 0} \sup_{\eta \in M_1^{s,\varepsilon}} \left[\int f(\lambda, \omega) \eta(d\omega) - h(\eta|P) - \lambda u \right] \\ &= - \sup_{\lambda \leq 0} \inf_{\eta \in M_1^{s,\varepsilon}} [G(\lambda, \eta, u) + h(\eta|P)]. \end{aligned} \tag{2.3.38}$$

Since $\eta \rightarrow - \int f(\lambda, \omega) \eta(d\omega) + h(\eta|P)$ is lower semi-continuous (for $\lambda \leq 0$) and $M_1(\Omega)$ is compact, the infimum in (2.3.38) is achieved for each λ , on measures in $M_1^{s,\varepsilon}$, for otherwise $h(\eta|P) = \infty$. Further, by (2.3.15), the supremum over λ can be taken over a compact set (recall that $\infty > u > 1$). By the Minimax theorem (see [64, Theorem 4.2] for this version), the min-max is equal to the max-min in (2.3.38). Further, since taking first the supremum in λ in the right

hand side of (2.3.38) yields a lower semicontinuous function, an achieving $\bar{\eta}$ exists, and then, due to compactness, there exists actually an achieving pair $\bar{\lambda}, \bar{\eta}$. We will show below that the infimum may be taken over stationary, ergodic measures only, that is

$$\inf_{\eta \in M_1^{s,\varepsilon}} \sup_{\lambda \leq 0} (G(\lambda, \eta, u) + h(\eta|P)) = \inf_{\eta \in M_1^{e,\varepsilon}} \sup_{\lambda \leq 0} (G(\lambda, \eta, u) + h(\eta|P)). \tag{2.3.39}$$

Then,

$$\begin{aligned} \text{R.H.S. of (2.3.38)} &= - \inf_{\eta \in M_1^{e,\varepsilon}} \sup_{\lambda \leq 0} (G(\lambda, \eta, u) + h(\eta|P)) \\ &= - \inf_{\eta \in M_1^{e,\varepsilon}} \inf_{w \leq u} (I_\eta^{\tau,q}(w) + h(\eta|P)). \end{aligned} \tag{2.3.40}$$

The second equality in (2.3.40) is obtained as follows: set $M_u = \{\eta \in M_1^{e,\varepsilon} : E_\eta(E_\omega^o(\tau_1 | \tau_1 < \infty)) > u\}$, $M_u^- = \{\eta \in M_1^{e,\varepsilon} : E_\eta(E_\omega^o(\tau_1 | \tau_1 < \infty)) \leq u\}$. For $\eta \in M_u$,

$$\inf_{w \leq u} I_\eta^{\tau,q}(w) = I_\eta^{\tau,q}(u) = \sup_{\lambda \in \mathbb{R}} G(\lambda, \eta, u) = \sup_{\lambda \leq 0} G(\lambda, \eta, u).$$

Further, recall that $I_\eta^{\tau,q}(\cdot)$ is convex with minimum value $\max(0, E_\eta(\log \rho_0))$ achieved at $E_\eta(E_\omega^o(\tau_1 | \tau_1 < \infty))$. Then, for $\eta \in M_u^-$,

$$\inf_{w \leq u} I_\eta^{\tau,q}(w) = \max(0, E_\eta(\log \rho_0))$$

whereas Jensen’s inequality implies that for such η ,

$$\sup_{\lambda \leq 0} G(\lambda, \eta, u) = G(0, \eta, u) = \max(0, E_\eta(\log \rho_0)),$$

completing the proof of (2.3.40). Hence,

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(T_n/n \leq u) &\leq - \inf_{w \leq u} \inf_{\eta \in M_1^{e,\varepsilon}} (I_\eta^{\tau,q}(w) + h(\eta|P)) \\ &= - \inf_{w \leq u} I_P^{\tau,q}(w). \end{aligned} \tag{2.3.41}$$

Turning to the proof of (2.3.39), we have, due to **(C4)** in Assumption 2.3.33, a sequence of stationary, ergodic measures with $\eta^n \rightarrow \bar{\eta}$ and $h(\eta^n|P) \rightarrow h(\bar{\eta}|P)$. Let λ_n be the maximizers in (2.3.39) corresponding to η^n . We have

$$\begin{aligned} \inf_{\eta \in M_1^{e,\varepsilon}} \sup_{\lambda \leq 0} \left(\left[\lambda u - \int f(\lambda, \omega) \eta(d\omega) \right] + h(\eta|P) \right) &\leq \left[\lambda_n u - \int f(\lambda_n, \omega) \eta^n(d\omega) \right] \\ &\quad + h(\eta^n|P). \end{aligned} \tag{2.3.42}$$

W.l.o.g. we can assume, by taking a subsequence, that $\lambda_n \rightarrow \lambda^* \leq 0$. Using the joint continuity in Lemma 2.3.35, we have, for $\varepsilon' > 0$ and $n \geq N_0(\varepsilon')$,

$$\begin{aligned} & \lambda_n u - \int f(\lambda_n, \omega) \eta^n(d\omega) + h(\eta^n|P) \\ & \leq \left[\lambda^* u - \int f(\lambda^*, \omega) \bar{\eta}(d\omega) \right] + h(\bar{\eta}|P) + \epsilon' \\ & \leq \inf_{\eta \in M_1^{s, \epsilon}} \sup_{\lambda \leq 0} \left(\left[\lambda u - \int f(\lambda, \omega) \eta(d\omega) \right] + h(\eta|P) \right) + \epsilon'. \end{aligned}$$

But this shows the equality in (2.3.39), since the reverse inequality there is trivial.

The upper bound for the upper tail, that is for $\frac{1}{n} \log P[\infty > \frac{1}{n} \sum_{j=1}^n \tau_j \geq u]$, where $1 < u < \infty$, is achieved similarly. We detail the argument since there is a small gap in the proof presented in [12]. First, exactly as in (2.3.38), one has

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(T_n/n \geq u) & \leq \inf_{0 \leq \lambda \leq \lambda_{\text{crit}}} \sup_{\eta \in M_1^{s, \epsilon}} [-G(\lambda, \eta, u) - h(\eta|P)] \\ & = - \sup_{0 \leq \lambda \leq \lambda_{\text{crit}}} \inf_{\eta \in M_1^{s, \epsilon}} [G(\lambda, \eta, u) + h(\eta|P)]. \end{aligned} \tag{2.3.43}$$

One may now apply the min-max theorem to deduce that the right hand side of (2.3.43) equals

$$\inf_{\eta \in M_1^{s, \epsilon}} \sup_{0 \leq \lambda \leq \lambda_{\text{crit}}} [G(\lambda, \eta, u) + h(\eta|P)] = \inf_{\eta \in M_1^{s, \epsilon}} \sup_{0 \leq \lambda \leq \lambda_{\text{crit}}} [G(\lambda, \eta, u) + h(\eta|P)],$$

where the second equality is proved by the same argument as in (2.3.39). Here a new difficulty arises: the supremum is taken over $\lambda \in [0, \lambda_{\text{crit}}(P)]$, but in general $\lambda_{\text{crit}}(\eta) \geq \lambda_{\text{crit}}(P)$ and hence the identification of the last expression with a variational problem involving $I_\eta^{T, q}(\cdot)$ is not immediate. To bypass this obstacle, we note, first by replacing η with $(1 - n^{-1})\eta + n^{-1}P$ and then using again **(C4)** to approximate with an ergodic measure, that the last expression equals

$$\inf_{\{\eta \in M_1^{e, \epsilon}, \lambda_{\text{crit}}(\eta) = \lambda_{\text{crit}}(P)\}} \sup_{0 \leq \lambda \leq \lambda_{\text{crit}}} [G(\lambda, \eta, u) + h(\eta|P)].$$

From here, one proceeds as in the case of the lower tail, concluding that

$$\begin{aligned} & \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(T_n/n \geq u) \\ & \leq - \inf_{\{\eta \in M_1^{e, \epsilon}, \lambda_{\text{crit}}(\eta) = \lambda_{\text{crit}}(P)\}} \sup_{0 \leq \lambda \leq \lambda_{\text{crit}}} [G(\lambda, \eta, u) + h(\eta|P)] \\ & = - \inf_{\{\eta \in M_1^{e, \epsilon}, \lambda_{\text{crit}}(\eta) = \lambda_{\text{crit}}(P)\}} \inf_{w \geq u} I_\eta^{T, q}(w) \leq - \inf_{\eta \in M_1^{e, \epsilon}} \inf_{w \geq u} I_\eta^{T, q}(w). \end{aligned}$$

This will then complete the proof of the (weak) upper bound, as soon as we prove the convexity of $I_P^{T, a}$. But, the function

$$\begin{aligned} & \sup_{\lambda \in \mathbb{R}} \inf_{\eta \in M_1^{s,\varepsilon}} [G(\lambda, \eta, u) + h(\eta|P)] \\ &= \sup_{\lambda \in \mathbb{R}} \left[\lambda u + \inf_{\eta \in M_1^{s,\varepsilon}} \left(- \int f(\lambda, \omega) \eta(d\omega) + h(\eta|P) \right) \right], \end{aligned} \tag{2.3.44}$$

being a supremum over affine functions in u , is clearly convex in u , while one shows, exactly as in (2.3.39), that

$$\inf_{\eta \in M_1^{s,\varepsilon}} \sup_{\lambda \in \mathbb{R}} [G(\lambda, \eta, u) + h(\eta|P)] = \inf_{\eta \in M_1^{e,\varepsilon}} \sup_{\lambda \in \mathbb{R}} [G(\lambda, \eta, u) + h(\eta|P)] \tag{2.3.45}$$

and therefore

$$\inf_{\eta \in M_1^{s,\varepsilon}} \sup_{\lambda \in \mathbb{R}} [G(\lambda, \eta, u) + h(\eta|P)] = \inf_{\eta \in M_1^{e,\varepsilon}} [I_\eta^{\tau,q}(u) + h(\eta|P)] = I_P^{\tau,a}(u).$$

Recalling that, as we saw above, supremum and infimum in (2.3.44) can be exchanged, this completes the proof of the upper bounds for the annealed LDP's for T_n/n .

Step I.3: Annealed lower bounds for T_n/n : We will use the following standard argument.

Lemma 2.3.46 *Let P be a probability distribution, (\mathcal{F}_n) be an increasing sequence of σ -fields and A_n be \mathcal{F}_n -measurable sets, $n = 1, 2, 3, \dots$. Let (Q_n) be a sequence of probability distributions such that $Q_n[A_n] \rightarrow 1$ and*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} H(Q_n|P) \Big|_{\mathcal{F}_n} \leq h$$

where $H(\cdot|P) \Big|_{\mathcal{F}_n}$ denotes the relative entropy w.r.t. P on the σ -field \mathcal{F}_n and h is a positive number. Then we have

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P[A_n] \geq -h.$$

Proof of Lemma 2.3.46. From the basic entropy inequality ([22], p. 423),

$$Q_n[A_n] \leq \frac{\log 2 + H(Q_n|P) \Big|_{\mathcal{F}_n}}{\log(1 + 1/P[A_n])}, \quad A_n \in \mathcal{F}_n,$$

we have $-Q_n[A_n] \log P[A_n] \leq \log 2 + H(Q_n|P) \Big|_{\mathcal{F}_n}$. Dividing by n and taking limits we obtain the desired result. □

We prove the lower bound for the lower tail only, the upper tail being handled by the same truncation as in the quenched case, see [12] for details. For $\eta \in M_1^{e,\varepsilon}$ satisfying $E_\eta(\log \rho_0) \leq 0$, define \overline{Q}_ω^o as in Step I.3 of Theorem 2.3.12, and let $\overline{Q}_\eta^o = \eta(d\omega) \otimes \overline{Q}_\omega^o$. Let $A_n = \{|n^{-1}T_n - u| < \delta\}$. We know already

that $\overline{Q}_\omega^o[A_n^c] \xrightarrow{n \rightarrow \infty} 0$, η - a.s., and this implies $\overline{Q}_\eta^o[A_n^c] \xrightarrow{n \rightarrow \infty} 0$. Let $\mathcal{F}_n := \sigma(\{\tau_i\}_{i=1}^n, \{\omega_j\}_{j=-\infty}^n)$, $\mathcal{F}_n^\omega = \sigma(\{\omega_j\}_{j=-\infty}^n)$. Note that

$$\overline{Q}_\eta^o|_{\mathcal{F}_n} = \eta|_{\mathcal{F}_n^\omega}(d\omega) \otimes \overline{Q}_\omega^o|_{\mathcal{F}_n}.$$

Hence,

$$H(\overline{Q}_\eta^o|\mathbb{P}^o)|_{\mathcal{F}_n} = H(\eta|P)|_{\mathcal{F}_n^\omega} + \int H(\overline{Q}_\omega^o|P_\omega^o)|_{\mathcal{F}_n} \eta(d\omega). \tag{2.3.47}$$

Considering the second term in (2.3.47), we have

$$\begin{aligned} & \frac{1}{n} \int H(\overline{Q}_\omega^o|P_\omega^o)|_{\mathcal{F}_n} \eta(d\omega) \\ &= -\frac{1}{n} \int \log Z_{n,\omega} \eta(d\omega) + \lambda_0(u, \eta) \int \frac{T_n}{n} d\overline{Q}_\omega^o \eta(d\omega) \\ &= -\frac{1}{n} \int \sum_{j=1}^n \log \varphi(\lambda_0(u, \eta), \theta^{j-1}\omega) \eta(d\omega) + \lambda_0(u) \int \frac{T_n}{n} d\overline{Q}_\omega^o \eta(d\omega) \end{aligned}$$

and we see, as in the proof of the lower bound of Theorem 2.3.12, that

$$\frac{1}{n} \int H(\overline{Q}_\omega^o|P_\omega^o)|_{\mathcal{F}_n} \eta(d\omega) \xrightarrow{n \rightarrow \infty} \lambda_0(u, \eta)u - E_\eta f(\lambda_0(u, \eta), \omega) \leq I_\eta^{\tau, q}(u).$$

Considering the first term in (2.3.47), we know that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} H(\eta|P)|_{\mathcal{F}_n^\omega} = h(\eta|P).$$

Hence,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} H(\overline{Q}_\eta^o|\mathbb{P}^o)|_{\mathcal{F}_n} \leq I_\eta^{\tau, q}(u) + h(\eta|P),$$

and we can now apply Lemma 2.3.46 to conclude that for any $\eta \in M_1^{e, \varepsilon}$ satisfying $E_\eta(\log \rho_0) \leq 0$ one has,

$$\liminf_{n \rightarrow \infty} E_P(A_n) \geq - (I_\eta^{\tau, q}(u) + h(\eta|P)).$$

As in the quenched case, one handles $\eta \in M_1^{e, \varepsilon}$ satisfying $E_\eta(\log \rho_0) > 0$ by repeating the above argument with the required (obvious) modifications, replacing \overline{Q}_ω^o by $\overline{Q}_\omega^o(\cdot|T_n < \infty)$. This completes the proof of Step I. \square

Proof of Lemma 2.3.35: For $\kappa > 1$, decompose $\varphi(\lambda, \omega)$ as follows:

$$\begin{aligned} E_\omega^o(e^{\lambda\tau_1} \mathbf{1}_{\tau_1 < \infty}) &= E_\omega^o(e^{\lambda\tau_1}; \tau_1 < \kappa) + E_\omega^o(e^{\lambda\tau_1}; \infty > \tau_1 \geq \kappa) \\ &:= \varphi_1^\kappa(\lambda, \omega) + \varphi_2^\kappa(\lambda, \omega), \end{aligned} \tag{2.3.48}$$

where $(\lambda, \omega) \rightarrow \log \varphi_1^\kappa(\lambda, \omega)$ is bounded and continuous. We also have

$$0 \leq \log \left(1 + \frac{\varphi_2^\kappa(\lambda, \omega)}{\varphi_1^\kappa(\lambda, \omega)} \right) \leq \log \left(1 + \frac{\varphi_2^\kappa(\lambda_{\text{crit}}, \omega)}{\varepsilon e^\lambda} \right).$$

Hence, the required continuity of the function $(\mu, \lambda) \rightarrow \int f(\lambda, \omega)\mu(d\omega)$ will follow from (2.3.48) as soon as we show that for any fixed constant $C_1 < 1$,

$$\lim_{\kappa \rightarrow \infty} \sup_{\mu \in M_1^{s, \varepsilon, P}} \int \log \left(1 + \frac{\varphi_2^\kappa(\lambda_{\text{crit}}, \omega)}{C_1} \right) \mu(d\omega) = 0. \tag{2.3.49}$$

If $\rho_{\min} < 1$ and $\rho_{\max} > 1$ then one can easily check, by a coupling argument using **(C4)**, that $\lambda_{\text{crit}} = 0$ (for a detailed proof see [12, Lemma 4]). Then, for each $\epsilon' > 0$ there exists a $\kappa_\mu = \kappa(\epsilon', \mu)$ large enough such that,

$$E_\mu \left(\log \left(1 + \frac{P_\omega^o(\infty > \tau_1 > \kappa_\mu)}{P_\omega^o(\tau_1 < \infty)} \right) \right) < \epsilon'.$$

Further, in this situation, for stationary, ergodic μ ,

$$\int f(0, \omega)\mu(d\omega) = \left(- \int \log \rho_0(\omega)\mu(d\omega) \right) \wedge 0. \tag{2.3.50}$$

In particular, $\mu \mapsto \int f(0, \omega)\mu(d\omega)$, being linear, is uniformly continuous on the compact set $M_1^{s, \varepsilon}$. Therefore, using (2.3.48), one sees that for each such μ one can construct a neighborhood B_μ of μ such that, for each $\nu \in B_\mu \cap M_1^{s, \varepsilon}$,

$$E_\nu \left(\log \left(1 + \frac{P_\omega^o(\infty > \tau_1 > \kappa_\mu + 1)}{P_\omega^o(\tau_1 < \infty)} \right) \right) < \epsilon'.$$

By compactness, it follows that there exists an $\kappa = \kappa(\epsilon')$ large enough such that, for all $\mu \in M_1^{s, \varepsilon}$,

$$E_\mu \left(\log \left(1 + \frac{P_\omega^o(\infty > \tau_1 > \kappa)}{P_\omega^o(\tau_1 < \infty)} \right) \right) < \epsilon'.$$

Using the inequality $\log(1 + cx) \leq c \log(1 + x)$, valid for $x \geq 0, c \geq 1$, one finds that for κ large enough,

$$\sup_{\mu \in M_1^{s, \varepsilon}} \int \log \left(1 + \frac{\varphi_2^\kappa(0, \omega)}{C_1} \right) \mu(d\omega) \leq \epsilon'/C_1,$$

proving (2.3.49) under the condition $\rho_{\min} < 1, \rho_{\max} > 1$.

We next handle the case $\rho_{\max} < 1$. We now complete the proof of Lemma 2.3.35 in the case $\rho_{\min} > 1$. We have $f(\lambda, \omega) \geq \lambda + \log \omega_0^+ \geq \lambda + \log \varepsilon$. We show that $(\lambda, \omega) \mapsto \varphi(\lambda, \omega)$ is continuous as long as $\omega_i \leq \omega^{\max}, \rho_i \leq \rho_{\max}$ and $\lambda \leq \lambda_{\text{crit}}$, which is enough to complete the proof. Write, for $\lambda \leq \lambda_{\text{crit}}$,

$$E_\omega(e^{\lambda\tau_1} \mathbf{1}_{\tau_1 < \infty}) = E_\omega(e^{\lambda\tau_1}; \tau_1 < \kappa) + E_\omega(e^{\lambda\tau_1}; \infty > \tau_1 \geq \kappa) \tag{2.3.51}$$

and observe that the first term in the right hand side of (2.3.51) is continuous as a function of ω and the second term goes to 0 for $\kappa \rightarrow \infty$, uniformly in ω . More precisely, due to (2.3.16), for all ω considered here,

$$E_\omega[e^{\lambda\tau_1}; \infty > \tau_1 \geq \kappa] \leq E_{\tilde{\omega}_{\min}}(e^{\lambda_{\text{crit}}\tau_1}; \tau_1 \geq \kappa) \xrightarrow{\kappa \rightarrow \infty} 0. \tag{2.3.52}$$

Finally, in the case $\rho_{\min} > 1$, the conclusion follows from the duality formula (2.3.23) and Remark 1 that follows its proof, by reducing the claim to the case $\rho_{\max} < 1$. □

Step II: The proof is identical to Step I, and is omitted.

Step III: The proof of all statements, except for the convexity of I_P^a , and the upper bound on $\mathbb{P}^o(X_n \leq nv)$, follow the argument in the quenched case. The latter proofs can be found in [12]. □

Remarks: 1. We note that under the conditions of Theorem 2.3.34, if $E_P \log \rho_0 \leq 0$ then both $I_P^a(v) \neq 0$ and $I_P^q(v) \neq 0$ for $v \notin [0, v_P]$. Indeed, since $h(\eta|P) \neq 0$ unless $\eta = P$, it holds that $I_P^a(v) = 0$ only if $I_P^q(v) = 0$. If $E_P \log \rho_0 = 0$, $v_P = 0$ and then for any $v \neq 0$, $I_P^q(v) = |v|I_P^{T,q}(1/|v|) > 0$ by the remark following the proof of Lemma 2.3.13. On the other hand, if $E_P \log \rho_0 < 0$, the same argument applies for $v > v_P$ while for $v < 0$ we have that $I_P^q(v) \geq -|v|E_P \log \rho_0 > 0$.

2. The condition **(C5)** can be avoided altogether. This is not hard to see if one is interested only in the LDP for T_n/n . Indeed, **(C5)** was used mainly in describing a worst case environment in the course of the proof of Lemma 2.3.35, see also part (d) of Lemma 2.3.13. When it is dropped, the following lemma, whose proof we provide below, replaces Lemma 2.3.35 when deriving the annealed LDP for T_n/n :

Lemma 2.3.53 *Assume P satisfies Assumption 2.3.33 except for **(C5)**. Then, $\lambda_{\text{crit}}(P)$ depends only on $\text{supp}(P_0)$, and the map $(\mu, \lambda) \mapsto E_\mu(f(\lambda, \omega))$ is continuous on $M_1^{s,\varepsilon,P} \times (-\infty, 0] \cup [0, \lambda_{\text{crit}})$.*

Given Lemma 2.3.53, we omit **(C5)** and replace **(C4)** in Assumption 2.3.33 by

(C4') *P is locally equivalent to the product of its marginals and, for any stationary measure $\eta \in M_1(\Omega)$ with $h(\eta|P) < \infty$ there is a sequence $\{\eta^n\}$ of stationary, ergodic measures, locally equivalent to the product of P 's marginals, with $\text{supp}((\eta^n)_0) = \text{supp}(P_0)$, $\eta^n \xrightarrow{n \rightarrow \infty} \eta$ weakly and $h(\eta^n|P) \rightarrow h(\eta|P)$.*

One now checks (we omit the details) that all approximations carried out in the proof of the upper bound of the upper tail of T_n/n can still be done, yielding the annealed LDP for T_n/n . To transfer this LDP to an annealed LDP for X_n/n does require a new argument, we refer to [16] for details.

We conclude our discussion of large deviation principles with the:

Proof of Lemma 2.3.53: Set $\Xi = \text{supp}(P_0)$ and define $\bar{\lambda} := \inf_{\omega \in \Xi} \lambda_{\text{crit}}(\omega)$ where

$$\lambda_{\text{crit}}(\omega) := \sup\{\lambda \in \mathbb{R} : E_\omega^o(e^{\lambda\tau_1} \mathbf{1}_{\{\tau_1 < \infty\}})\}.$$

By definition, $\lambda_{\text{crit}}(P) \geq \bar{\lambda}$. On the other hand, if $\lambda > \bar{\lambda}$ then there exists a $\bar{\omega} \in \Xi^{\mathbb{Z}}$ with $E_{\bar{\omega}}^o(e^{\lambda\tau_1} \mathbf{1}_{\{\tau_1 < \infty\}}) = \infty$. Fix $K = e^{-\lambda}/\varepsilon$, and using monotone convergence, choose an M large enough such that

$$\varphi_{M,\bar{\omega}}(\lambda) := E_{\bar{\omega}}^o(e^{\lambda\tau_1} \mathbf{1}_{\{\tau_1 < M\}}) > K + 1.$$

Since $\varphi_{M,\bar{\omega}}$ depends only on $\{\omega_i, i \in (-M, 0)\}$, it holds that with positive P -probability, $E_{\bar{\omega}}^o(e^{\lambda\tau_1} \mathbf{1}_{\{\tau_1 < M\}}) \geq K + 1$. But, if $\lambda \leq \lambda_{\text{crit}}(P)$ it follows from the recursions (2.3.21) that $\varphi(\lambda, \omega) < K$, P -a.s., a contradiction unless $\bar{\lambda} = \lambda_{\text{crit}}(P)$. In particular, $\lambda_{\text{crit}}(P)$ depends only on $\text{supp}(P_0)$. Note that the characterization of $\lambda_{\text{crit}}(P)$ as $\bar{\lambda}$ implies that for any $\mu \in M_1^{s,\varepsilon,P}$ it holds that $\lambda_{\text{crit}}(\mu) \geq \lambda_{\text{crit}}(P)$.

Next, as in the course of the proof of Lemma 2.3.35, see (2.3.49), it is enough to show that for any $\lambda < \lambda_{\text{crit}}(P)$,

$$\lim_{\kappa \rightarrow \infty} \sup_{\mu \in M_1^{s,\varepsilon,P}} \int \varphi_2^\kappa(\lambda, \omega) \mu(d\omega) = 0. \tag{2.3.54}$$

But, since $\varphi(\lambda, \omega) \leq e^{-\lambda}/\varepsilon$ μ -a.s. for all $\mu \in M_1^{s,\varepsilon,P}$ (use again the recursions (2.3.21) and that $\lambda_{\text{crit}}(\mu) \geq \lambda_{\text{crit}}(P)$), it holds that

$$\int \varphi_2^\kappa(\lambda, \omega) \mu(d\omega) \leq e^{(\lambda - \lambda_{\text{crit}})M} \frac{e^{-\lambda}}{\varepsilon},$$

yielding immediately (2.3.54). □

Bibliographical notes: The first quenched LDP result is due to Greven and Den Hollander, [34], who proved it for i.i.d. environments using the method of the environment viewed from the particle. Our derivation here follows the hitting times approach developed in [12], except that the proof of Lemma 2.3.22 follows the article [58]. Extensions of the LDP’s in this chapter to more general models allowing for (non geometric) holding times is presented in [16], where the derivation avoids completely coupling arguments and thus bypasses altogether the need for **(C5)** in deriving the annealed LDP for X_n/n .

The “process level LDP” for R_n was first proved in [23] in the context of Markov chains with law P satisfying appropriate regularity conditions. It was extended to various ergodic situation in [55] and [56], see also [11]. We refer to [27] and [19, Chapter 6] for further information. Our presentation of the annealed LDP follows here [12], where additional information on the shape of the rate functions etc. can be found. Note that [12] treats the case $P(\omega_0^0 = 0) = 1$. In the exposition here, we corrected and simplified some of the arguments in [12], following [16], where a RWRE with general (i.e., not necessarily geometric) holding times is considered. Finally, a completely different approach to the derivation of the LDP, both annealed and quenched, is described in [79].

2.4 The subexponential regime

We saw in Section 2.3 that, at least when P satisfies Assumption 2.3.33 and $E_P \log \rho_0 \leq 0$, we have that for any δ small enough, any $v \in [0, v_P]$,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) \\ = \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) = 0. \end{aligned} \tag{2.4.1}$$

Our goal in this section is to obtain more precise information on the rate of convergence in (2.4.1). Surprisingly, it turns out that it is better to consider first the annealed case.

Throughout this section, we impose the following assumption on the law P . Together with **(C4)** there, it implies Assumption 2.3.33.

Assumption 2.4.2

- (D1)** *There exists an $\varepsilon > 0$ such that $\min(\omega_0^+, \omega_0^-) > \varepsilon$, P -a.s.*
- (D2)** $\rho_{\min} < 1$, $\rho_{\max} > 1$, and $E_P \log \rho_0 \leq 0$.
- (D3)** P is α -mixing with $\alpha(n) = \exp(-n(\log n)^{1+\eta})$ for some $\eta > 0$; that is, for any ℓ -separated measurable bounded by 1 functions f_1, f_2 ,

$$E_P \left(f_1(\omega)(f_2(\omega) - E_P f_2(\omega)) \right) \leq \alpha(\ell).$$

(functions f_i are ℓ separated if f_i is measurable on $\sigma(\omega_j, j \in I_i)$ with I_i intervals satisfying $\text{dist}(I_i, I_k) > \ell$ for any $i \neq k$).

It is known that **(D3)** implies **(C1)** and **(C3)** of Assumption 2.3.33, see [10]. In particular, letting $\bar{R}_k := k^{-1} \sum_{i=0}^{k-1} \log \rho_i$, it implies that \bar{R}_k satisfies the LDP with good rate function $J(\cdot)$. We add the following assumption on $J(\cdot)$:

- (D4)** $J(0) > 0$.

Condition **(D4)** implies that $E_P(\log \rho_0) < 0$. Define next $s := \min_{y \geq 0} \frac{1}{y} J(y)$. Note that the condition $E_P(\bar{S}) < \infty$ and the existence of a LDP for \bar{R}_k with good rate function $J(\cdot)$ are enough to imply, by Varadhan’s lemma, that $0 \geq \sup_y (y - J(y))$, and in particular that $s \geq 1$. (In the case where P is a product measure, we can identify s as satisfying $E_P(\rho_0^s) = 1$, and then $E_P(\bar{S}) < \infty$, which is equivalent to $E_P(\rho_0) < 1$, implies that $s > 1$.)

Annealed subexponential estimates

Theorem 2.4.3 *Assume P satisfies Assumption 2.4.2, and $v_P > 0$. Then, for any $v \in (0, v_P)$ and any $\delta > 0$ small enough,*

$$\lim_{n \rightarrow \infty} \frac{\log \mathbb{P}^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right)}{\log n} = 1 - s.$$

Proof. We begin by proving the lower bound. Fix $0 < v - \delta < v - 4\eta < v < v_P$; let

$$L_k = \max\left\{n \geq T_k : (k - X_n)\right\}$$

denote the largest excursion of $\{X_n\}$ to the left of k after hitting it. Observe that the event $\{n^{-1}X_n \in (v - \delta, v + \delta)\}$ contains the event

$$A := \left\{ \frac{(v - 4\eta)n}{v_P} < T_{(v-2\eta)n} < n, L_{(v-2\eta)n} < \frac{\eta n}{2}, T_{vn} > n \right\}, \tag{2.4.4}$$

namely, the RWRE hits $(v - 2\eta)n$ at about the expected time, from which point its longest excursion to the left is less than $\eta n/2$, but the RWRE does not arrive at position vn by time n .

Next, note that by (2.1.4),

$$P_\omega^\circ\left(L_{(v-2\eta)n} \geq \eta n/2\right) \leq \sum_{i=0}^\infty \prod_{j=-(\eta n/2-1)}^i \rho_{(v-2\eta)n+j}. \tag{2.4.5}$$

Hence, using the LDP for \bar{R}_k , we have for all n large enough

$$\begin{aligned} \mathbb{P}^\circ\left(L_{(v-2\eta)n} \geq \eta n/2\right) &\leq \sum_{i=0}^\infty E\left(e^{(\eta n/2+i)\bar{R}_{\eta n/2+i}}\right) \\ &\leq e^{-\eta n \sup_y (y - J(y))/4} \leq e^{-\delta_1 n} \end{aligned} \tag{2.4.6}$$

for some $\delta_1 > 0$. Thus, for all n large enough,

$$\begin{aligned} \mathbb{P}^\circ(A) &\geq \mathbb{P}^\circ\left(\frac{(v - 4\eta)n}{v_P} < T_{(v-2\eta)n} < n, T_{vn} > n\right) - e^{-\delta_1 n} \\ &\geq \mathbb{E}^\circ\left(P_\omega^\circ\left(\frac{(v - 4\eta)n}{v_P} < T_{(v-2\eta)n} < n\right)\right. \\ &\quad \left.P_\omega^{(v-\eta)n}\left(T_{vn} > \frac{4\eta n}{v_P}, L_{(v-\eta)n} < \eta n/2\right)\right) - e^{-\delta_1 n} \\ &\geq B \cdot C - \alpha(\eta n/2) - 2e^{-\delta_1 n}, \end{aligned}$$

where

$$\begin{aligned} B &= \mathbb{P}^\circ\left(\frac{(v - 4\eta)n}{v_P} < T_{(v-2\eta)n} < n\right) \\ C &= \mathbb{P}^\circ\left(T_{\eta n} > \frac{4\eta n}{v_P}\right). \end{aligned}$$

and $\alpha(\cdot)$ is as in **(D3)**.

Next, note that $B \rightarrow_{n \rightarrow \infty} 1$ by (2.1.16). We will prove below that for any $\delta' > 0$,

$$C \geq n^{1-s-2\delta'} \tag{2.4.7}$$

and this implies that for all n large, $\mathbb{P}^o(A) \geq n^{1-s-4\delta'}$, which yields the required lower bound (recall δ' is arbitrary!) as soon as we prove (2.4.7).

Turning to the proof of (2.4.7), fix y such that $\frac{J(y)}{y} \leq s + \frac{\delta'}{4}$, $K = \lceil n^{\frac{\delta'}{4}} \rceil$, $k = \lceil \frac{1}{y} \log n \rceil$, and set $\overline{m}_K = \lceil \eta n / K \rceil$. Now, using **(D3)**,

$$\begin{aligned} P \left(\bigcap_{j=1}^{\overline{m}_K} \{ \overline{R}_k(\theta^{jK} \omega) \leq y \} \right) &\leq (P(\overline{R}_k(\omega) \leq y))^{\overline{m}_K} + \overline{m}_K \alpha(K - k) \\ &= (1 - P(\overline{R}_k(\omega) > y))^{\overline{m}_K} + \overline{m}_K \alpha(K - k) \\ &\leq \left(1 - e^{-k(J(y) + \frac{\delta' y}{4})} \right)^{\overline{m}_K} + \overline{m}_K \alpha(K - k) \leq 1 - n^{1-s-\delta'}, \end{aligned}$$

for all n large enough. Hence,

$$P\left(\exists j \in \{1, \dots, \overline{m}_K\} : \overline{R}_k(\theta^{jK} \omega) > y\right) \geq n^{1-s-\delta'}. \tag{2.4.8}$$

On the other hand, let ω and $j \leq \overline{m}_K$ be such that $\overline{R}_k(\theta^{jK} \omega) > y$. Then, using (2.1.6) in the second inequality, for such ω ,

$$\begin{aligned} P_\omega^o \left(T_{\eta n} > \frac{4\eta n}{v_P} \right) &\geq P_\omega^{jK} \left(T_k > \frac{4\eta n}{v_P} \right) \geq (1 - e^{-ky})^{\frac{4\eta n}{v_P}} \\ &\geq \left(1 - \frac{1}{n} \right)^{\frac{6\eta n}{v_P}} \geq e^{-\frac{8\eta}{v_P}}. \end{aligned} \tag{2.4.9}$$

Combining (2.4.8) and (2.4.9), we conclude that

$$C \geq n^{1-s-\delta'} \cdot e^{-\frac{8\eta}{v_P}},$$

as claimed.

We next turn to the proof of the upper bounds. We may and will assume that $s > 1$, for otherwise there is nothing to prove. We first note that, for some $\delta'' := \delta''(\delta) > 0$,

$$\begin{aligned} \mathbb{P}^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) &\leq \mathbb{P}^o \left(\frac{X_n}{n} < v + \delta \right) \\ &\leq \mathbb{P}^o(T_{n(v+2\delta)} > n) + \mathbb{P}^o(L_0 > n\delta) \\ &\leq \mathbb{P}^o(T_{n(v+2\delta)} > n) + e^{-\delta'' n} \end{aligned} \tag{2.4.10}$$

where the stationarity of P was used in the second inequality, and (2.4.6) in the third. Thus, the required upper bound follows once we show that for any $v < v_P$, any $\delta' > 0$,

$$\mathbb{P}^o(T_{nv} > n) \leq n^{1-s+\delta'} \tag{2.4.11}$$

for all n large enough.

Set $a := \sup_y (y - J(y))$. Because $s > 1$ and $J(0) > 0$, it holds that $a < 0$. Fix $A > -s/a$, and set $k = k(n) = A \log n$. Next, define the process $\{Y_n\}$ in $\mathbb{Z}^{\mathbb{N}}$ and the hitting times $\tilde{T}_{ik} = \min(n \geq 0 : Y_n = ik), i = 0, 1, \dots$ such that the only change between the processes $\{X_n\}$ and $\{Y_n\}$ is that the process $\{Y_n\}_{n \geq \tilde{T}_{ik}}$ is reflected at position $(i - 1)k$ (with a slight abuse of notations, we continue to use $P_\omega^\circ, \mathbb{P}^\circ$ to denote the law of $\{Y_n\}$ as well as that of $\{X_n\}$). Set $m_k = \lceil vn/k \rceil + 1$, and $\tilde{\tau}_k^{(i)} = \tilde{T}_{ik} - \tilde{T}_{(i-1)k}, i = 1, \dots, m_k$. Note that the $\tilde{\tau}_k^{(i)}$ are identically distributed, each stochastically dominated by T_k . Hence, $\mathbb{E}^\circ \tilde{T}_{ik} \leq \mathbb{E}^\circ T_{ik}$. On the other hand, fixing $\lambda \in (1/s, 1)$, we will see below (cf. Lemma 2.4.16) that $\mathbb{E}^\circ(T_k^{1/\lambda}) \leq ck^{1/\lambda}$ for some $c := c(\lambda)$, yielding, by Hölder's inequality, that

$$\mathbb{E}^\circ T_k \leq \mathbb{E}^\circ(\tilde{T}_k) + \mathbb{P}^\circ(L_0 \geq k)^{1-\lambda} \mathbb{E}^\circ(T_k^{1/\lambda})^\lambda \leq \mathbb{E}^\circ(\tilde{T}_k) + ck \mathbb{P}^\circ(L_0 \geq k)^{1-\lambda}.$$

Thus, using (2.4.6) and the fact that $\mathbb{E}^\circ(T_k)/k = v_P$, we conclude that $\lim_{k \rightarrow \infty} \mathbb{E}^\circ T_k / \mathbb{E}^\circ \tilde{T}_k = 1$, implying that $\mathbb{E}^\circ \tilde{T}_k / k \rightarrow_{k \rightarrow \infty} 1/v_P$.

Next, note that on the event $\{L_{ik} < k \text{ for } i = 0, \dots, m_k\}$, the processes $\{X_n\}$ and $\{Y_n\}$ coincide for $n < T_{m_k k}$. Hence

$$\mathbb{P}^\circ(T_{nv} > n) \leq \mathbb{P}^\circ(\tilde{T}_{m_k k} > n) + m_k \mathbb{P}^\circ(L_0 > k). \tag{2.4.12}$$

But, as in (2.4.6), for k large enough

$$\mathbb{P}^\circ(L_0 > k) \leq E_P(e^{k(\bar{R}_k + \delta)}) \leq e^{\log n(Aa + \delta')} \leq n^{-s + \delta''},$$

where $\delta'' := \delta''(\delta) \rightarrow_{\delta \rightarrow 0} 0$. Since $m_k < n$, the second term in (2.4.12) is of the right order, and the upper bound follows as soon as we prove that, for n large enough

$$\mathbb{P}^\circ(\tilde{T}_{m_k k} > n) \leq n^{1-s+\delta'}. \tag{2.4.13}$$

To see (2.4.13), note that $\tilde{T}_{m_k k} = \sum_{i=1}^{m_k} \tilde{\tau}_k^{(i)}$, with $\mathbb{E}^\circ(\tilde{\tau}_k^{(i)})/k = \mathbb{E}^\circ(\tilde{T}_k)/k \rightarrow 1/v_P$. Hence, for some $\eta > 0$, using that $km_k \leq v < v_P$,

$$\begin{aligned} \mathbb{P}^\circ(\tilde{T}_{m_k k} > n) &\leq \mathbb{P}^\circ\left(\sum_{i=1}^{m_k} (\tilde{\tau}_k^{(i)} - \mathbb{E}^\circ(\tilde{\tau}_k^{(i)})) > 4\eta n\right) \\ &\leq 4\mathbb{P}^\circ\left(\sum_{i=1}^{\lceil m_k/4 \rceil} (\tilde{\tau}_k^{(4i)} - \mathbb{E}^\circ(\tilde{T}_k)) > \eta n\right). \end{aligned}$$

Note that the quenched law of $\tilde{\tau}_k^{(4i)}$ depends on $\{\omega_j, j \in I_i\}$ where $I_i = \{4i - k, 4i - k + 1, \dots, 4i + k\}$. Let $\{\bar{\tau}_k^{(i)}\}$ be i.i.d. random variables such that for any Borel set $G, P(\bar{\tau}_k^{(i)} \in G) = \mathbb{P}^\circ(\tilde{\tau}_k^{(i)} \in G)$. Then, by iterating the definition of $\alpha(\cdot)$, one has that

$$\mathbb{P}^o \left(\sum_{i=1}^{\lceil m_k/4 \rceil} (\tilde{\tau}_k^{(4i)} - \mathbb{E}^o(\tilde{T}_k)) > \eta n \right) \leq P \left(\sum_{i=1}^{\lceil m_k/4 \rceil} (\bar{\tau}_k^{(4i)} - \mathbb{E}^o(\bar{T}_k)) > \eta n \right) + \frac{m_k \alpha(2k)}{4}. \quad (2.4.14)$$

We recall that

$$\frac{m_k \alpha(2k)}{4} \leq o(n^{1-s}). \quad (2.4.15)$$

The following estimate, whose proof is deferred, is crucial to the proof of (2.4.13):

Lemma 2.4.16 *For each $\kappa < s$, there exists a constant $c(\kappa) < \infty$ such that*

$$\mathbb{E}^o(T_k)^\kappa \leq c(\kappa)k^\kappa. \quad (2.4.17)$$

By Markov’s inequality, for any $\kappa < \kappa' < s$,

$$P(\bar{\tau}_k^{(4i)} - E\bar{\tau}_k^{(4i)} > \eta n) \leq \frac{1}{(\eta n)^{\kappa'}} E|\bar{\tau}_k^{(4i)} - E\bar{\tau}_k^{(4i)}|^{\kappa'} \leq n^{-\kappa}$$

where n is large enough and we used Lemma 2.4.16 and the fact that $E((\bar{\tau}_k^{(4i)})^{\kappa'}) = \mathbb{E}^o((\tilde{\tau}_k^{(4i)})^{\kappa'}) \leq \mathbb{E}^o(T_k^{\kappa'})$. Hence, (see [54, (1.3),(1.7a)]),

$$P \left(\sum_{i=1}^{\lceil m_k/4 \rceil} (\bar{\tau}_k^{(4i)} - E\bar{\tau}_k^{(4i)}) > \eta n \right) \leq \lceil \frac{m_k}{4} \rceil P(\bar{\tau}_k^{(4)} - E\bar{\tau}_k^{(4)} > \eta n) + \frac{1}{2} n^{1-\kappa} \leq n^{1-\kappa}.$$

Since $\kappa < s$ is arbitrary, this completes the proof, modulo the

Proof of Lemma 2.4.16

Note first that by Minkowski’s inequality, for any $k \geq 1$,

$$\mathbb{E}^o(T_k^\kappa) = \mathbb{E}^o \left(\sum_{i=1}^k \tau_i \right)^\kappa \leq k^\kappa \mathbb{E}^o \tau_1^\kappa.$$

Hence, it will be enough to prove that

$$\mathbb{E}^o(\tau_1^\kappa) < \infty. \quad (2.4.18)$$

To prove (2.4.18), we build upon the techniques developed in the course of proving Lemma 2.1.21. Indeed, recall the random variables $U_{i,j}, Z_{i,j}$ and N_i defined there, and note that since $\tau_1 = \sum_{i=-\infty}^o N_i$, it is enough to estimate

$$\mathbb{E}^\circ \left(\sum_{i=-\infty}^0 N_i \right)^\kappa = \mathbb{E}^\circ \left(\sum_{i=-\infty}^0 U_i + U_{i+1} + Z_i \right)^\kappa \leq C_\varepsilon \mathbb{E}^\circ \left(\sum_{i=-\infty}^0 U_i \right)^\kappa. \tag{2.4.19}$$

An important step in the evaluation of the RHS in (2.4.19) involves the computation of moments of U_i . To present the idea, consider first the case $\kappa > 2$, and write

$$U_i = \sum_{j=1}^{U_{i+1}} G_j$$

where, under P_ω° , the G_j are i.i.d. geometric random variables, independent of $\{U_{i+1}, \dots, U_0\}$, of parameter $\frac{\omega_i^-}{\omega_i^- + \omega_i^+}$. Hence,

$$\begin{aligned} E_\omega^\circ(U_i^2) &= E_\omega^\circ \left(\sum_{j=1}^{U_{i+1}} (G_j - E_\omega^\circ G_j) + \sum_{j=1}^{U_{i+1}} E_\omega^\circ G_j \right)^2 \\ &\leq c_\delta E_\omega^\circ \left(\sum_{j=1}^{U_{i+1}} (G_j - E_\omega^\circ G_j) \right)^2 + (1 + \delta)(E_\omega^\circ G_j)^2 \cdot E_\omega^\circ(U_{i+1}^2) \\ &\leq c'_\delta E_\omega^\circ(U_{i+1}) \cdot E_\omega^\circ(G_j^2) + (1 + \delta)\rho_i^2 E_\omega^\circ(U_{i+1}^2). \end{aligned} \tag{2.4.20}$$

Here, c_δ, c'_δ are constants which depend on δ only. Since $E_\omega^\circ(G_j^2)$ is uniformly (in ω) bounded, and $E_\omega^\circ(U_{i+1}) = \rho_i$, we get

$$E_\omega^\circ(U_i^2) \leq c''_\delta \rho_i E_\omega^\circ U_{i+1} + (1 + \delta)\rho_i^2 E_\omega^\circ(U_{i+1}^2).$$

Iterating and using (cf. (2.1.24)) that $E_\omega^\circ U_{i+1} = \prod_{j=i+1}^0 \rho_j$, we conclude the existence of a constant c'''_δ such that

$$E_\omega^\circ(U_i^2) \leq c'''_\delta \left(\sum_{j=0}^{|i|} \left(\prod_{k=-j}^0 \rho_k + \prod_{k=-j}^0 (\rho_k^2(1 + \delta)) \right) \right),$$

and hence

$$\mathbb{E}^\circ(U_i^2) \leq c'''_\delta \sum_{j=0}^{|i|} \left(E_P \prod_{k=-j}^0 \rho_k + E_P \prod_{k=-j}^0 (\rho_k^2(1 + \delta)) \right). \tag{2.4.21}$$

Note that, by Varadhan’s lemma (see [19, Theorem 4.3.1]), for any constant β ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{E}^\circ \left(\prod_{k=-n}^0 \rho_k^\beta \right) = \sup_y \left(\beta y - J(y) \right) = \sup_y y \left(\beta - \frac{J(y)}{y} \right) := \beta'(\beta), \tag{2.4.22}$$

and $\beta'(\beta) < 0$ as soon as $\beta < s$. Hence, substituting in (2.4.21), and choosing δ such that $\log(1 + \delta) < \beta'(\beta)/4$, we obtain that for some constant c_δ'''' ,

$$\mathbb{E}^o(U_i^2) \leq c_\delta'''' e^{-i\beta'(2)/2},$$

implying that

$$\sqrt{\mathbb{E}^o\left(\sum_{i=-\infty}^0 N_i\right)^2} \leq C_\varepsilon^{\frac{1}{2}} \sqrt{\sum_{i=-\infty}^0 \mathbb{E}^o(U_i^2)} < \infty.$$

A similar argument holds for any integer $\kappa < s$: mimicking the steps leading to (2.4.21), we get that

$$E_\omega^o(U_i^\kappa) \leq c_\delta'''' \left(\sum_{j=0}^{|i|} \left(\prod_{k=-j}^0 \rho_k^{\kappa/2} + \prod_{k=-j}^0 \rho_k^\kappa \right) \right),$$

and using (2.4.22) and an induction on lower (integer) moments, we get that $\mathbb{E}^o\left(\sum_{i=-\infty}^0 N_i\right)^\kappa < \infty$ for all $\kappa < s$ integer. Finally, to handle $\lfloor s \rfloor < \kappa < s$, we replace (2.4.20) by

$$\begin{aligned} E_\omega^o(U_i^\kappa) &\leq c'_\delta E_\omega^o(U_{i+1}^{\kappa/2\vee 1}) E_\omega^o(G_j^\kappa) + (1 + \delta) \rho_i^\kappa E_\omega^o(U_{i+1}^\kappa) \\ &\leq c_\delta'' (E_\omega^o U_i^{\lceil \kappa/2 \rceil})^{\frac{\kappa/2}{\lceil \kappa/2 \rceil}} + (1 + \delta) \rho_i^\kappa E_\omega^o(U_{i+1}^\kappa), \end{aligned}$$

and one proceeds as before. □

Quenched subexponential estimates

Theorem 2.4.23 *Assume P satisfies Assumption 2.4.2, and $v_P > 0$. Then, for any $v \in (0, v_P)$, any $\eta > 0$, and any $\delta > 0$ small enough,*

$$\liminf_{n \rightarrow \infty} \frac{1}{n^{1-1/s+\eta}} \log P_\omega^o\left(\frac{X_n}{n} \in (v - \delta, v + \delta)\right) = 0, \quad P - a.s. \quad (2.4.24)$$

Further,

$$\limsup_{n \rightarrow \infty} \frac{1}{n^{1-1/s-\eta}} \log P_\omega^o\left(\frac{X_n}{n} < v\right) = -\infty, \quad P - a.s.. \quad (2.4.25)$$

Proof. Starting with the lower bound, we have, using (2.4.4) and (2.4.5), that for some $\delta_1(\omega) > 0$,

$$\begin{aligned} P_\omega^o\left(\frac{X_n}{n} \in (v - \delta, v + \delta)\right) &\geq P_\omega^o\left(\frac{(v - 4\eta)n}{v_P} < T_{(v-2\eta)n} < n\right) \\ &\quad P_\omega^{(v-\eta)n}\left(T_{vn} > \frac{4\eta}{v_P}n, L_{(v-\eta)n} < \eta n\right) - e^{-\delta_1(\omega)n}. \end{aligned}$$

By (2.1.16), $P_\omega^o(\frac{X_n}{n} \in (v - \delta, v + \delta)) \rightarrow_{n \rightarrow \infty} 1$, P -a.s. On the other hand, as in the proof of (2.4.8), fix y such that $J(y)/y \leq s + \delta'/4$, $k = \lfloor (1 - \delta')/ys \rfloor$, and $K = n^{\delta'/4}$. Then, one checks as in the annealed case that

$$P(\forall j \in \{1, \dots, \overline{m}_K\} : \overline{R}_k(\theta^{jK} \omega) \leq y) \leq \frac{1}{n^2},$$

and one concludes by the Borel-Cantelli lemma that there exists an $n_0(\omega)$ such that for all $n_0(\omega)$, there exists a $j \in \{1, \dots, \overline{m}_K\}$ such that $\overline{R}_k(\theta^{jK} \omega) > y$. The lower bound (2.4.24) now follows as in the proof of (2.4.9).

Turning to the proof of the upper bound (2.4.25), as in the annealed setup it is straightforward to reduce the proof to proving

$$\lim_{n \rightarrow \infty} \frac{1}{n^{1-1/s-\delta}} \log P_\omega^o(T_n > n/v) = -\infty. \tag{2.4.26}$$

We provide now a short sketch of the proof of (2.4.26) before getting our hands dirty in the actual computations. Divide the interval $[0, nv]$ into blocks of size roughly $k = k_n := n^{1/s+\delta}$. By using the annealed bounds of Theorem 2.4.3, one knows that $P(T_k > k/v) \sim k^{1-s}$. Hence, taking appropriate subsequences, one applies a Borel-Cantelli argument to control uniformly the probability $P_\omega^{ik}(T_{(i+1)k} > k/v)$, c.f. Lemma 2.4.28.

The next step involves a decoupling argument. Define

$$\overline{T}_{(i+1)k} = \inf \{t > T_{ik} : X_t = (i + 1)k \text{ or } X_t = (i - 1)k\}. \tag{2.4.27}$$

Then one shows that for all relevant blocks, that is $i = \pm 1, \pm 2, \dots, \pm n/k$, $P_\omega^{ik}(\overline{T}_{(i+1)k} \neq T_{(i+1)k})$ is small enough. Therefore, we can consider the random variables $\overline{T}_{(i+1)k} - T_{ik}$ instead of $T_{(i+1)k} - T_{ik}$, which have the advantage that their dependence on the environment is well localized. This allows us to obtain a uniform bound on the tails of $\overline{T}_{(i+1)k} - T_{ik}$, for all relevant i , see (2.4.30).

The final step involves estimating how many of the k -blocks will be traversed from right to left before the RWRE hits the point nv . This is done by constructing a simple random walk (SRW) S_t whose probability of jump to the left dominates $P_\omega^{ik}(T_{(i+1)k} \neq \overline{T}_{(i+1)k})$ for all relevant i . The analysis of this SRW will allow us to claim (c.f. Lemma 2.1.17) that the number of visits to a k -block after entering its right neighbor is negligible. Thus, the original question on the tail of T_n is replaced by a question on the sum of (dominated by i.i.d.) random variables, which is resolved by means of the tail estimates obtained in the second step.

A slight complication is presented by the need to work with subsequences in order to apply the Borel-Cantelli lemma at various places. Going from subsequences to the original n sequence is achieved by means of monotonicity arguments. Indeed, by monotonicity, note that it is enough to prove the result when, for arbitrary δ small enough, n is replaced by the subsequence $n_j = \lfloor j^{2/\delta} \rfloor$, since $n_{j+1}/n_j \rightarrow_{j \rightarrow \infty} 1$.

Turning to the actual proof, fix $C_n = n^\delta$, $k = k_j = \frac{C_{n_j} n_j^{1/s}}{1 - \varepsilon}$ for some $1 > \varepsilon > 0$, $b_n = C_n^{-\delta}$, and $I_j = \left\{ -\left\lfloor \frac{n_j}{k_j} \right\rfloor - 1, \dots, \left\lfloor \frac{n_j}{k_j} \right\rfloor + 1 \right\}$. Finally, fix $v' < v$ and $\bar{T}_{(i+1)k}$ as in (2.4.27). (We will always use $\bar{T}_{(i+1)k}$ in conjunction with the RWRE started at $ik!$). We now claim the:

Lemma 2.4.28 *For P - a.e. ω , there exists a $J_0(\omega)$ such that for all $j > J_0(\omega)$, and all $i \in I_j$,*

$$P_\omega^{ik} \left(\frac{T_{(i+1)k_j}}{k_j} > \frac{1}{v'} \right) \leq b_{n_j} . \tag{2.4.29}$$

Further, for all $j > J_0(\omega)$, and each $i \in I_j$, and for $x \geq 1$,

$$P_\omega^{(ik)} \left(\frac{\bar{T}_{(i+1)k_j}}{k_j} > \frac{x}{v'} \right) \leq (2b_{n_j})^{\lfloor x/2 \rfloor \vee 1} . \tag{2.4.30}$$

Proof of Lemma 2.4.28. By Chebycheff’s bound,

$$\begin{aligned} P \left(P_\omega^{ik} \left(\frac{T_{(i+1)k_j}}{k_j} > \frac{1}{v'} \right) > b_{n_j} \right) &\leq \frac{1}{b_{n_j}} \mathbb{P}^{ik} \left(\frac{T_{(i+1)k_j}}{k_j} > \frac{1}{v'} \right) \\ &\leq \frac{1}{b_{n_j}} k_j^{1-s+o(1)} , \end{aligned}$$

where the last inequality follows from Theorem 2.4.3. Hence,

$$\begin{aligned} P \left(P_\omega^{ik} \left(\frac{T_{(i+1)k_j}}{k_j} > \frac{1}{v'} \right) > b_{n_j} \text{ for some } i \in I_j \right) &\leq 3 \left\lfloor \frac{n_j}{k_j} \right\rfloor \cdot \frac{1}{b_{n_j}} \cdot k_j^{1-s+o(1)} \\ &\leq \frac{3}{n_j^{\delta(s-o(1)-\delta)}} \leq \frac{4}{j^{2(s-o(1)-\delta)}} \end{aligned}$$

and (2.4.29) follows from the Borel-Cantelli lemma. (2.4.30) follows by iterating this inequality and using the Markov property. \square

Recall that $a = \sup_y (y - J(y)) < 0$ and let $0 < \theta < -\frac{a}{1-\varepsilon/4}$, $d_n^\theta = e^{-\theta n^{1/s} C_n}$. We now have:

Lemma 2.4.31 *For P - a.e. ω , there is a $J_1(\omega)$ s.t. for all $j \geq J_1(\omega)$, all $i \in I_j$,*

$$P_\omega^{ik} \left(\bar{T}_{(i+1)k_j} \neq T_{(i+1)k_j} \right) \leq d_{n_j}^\theta .$$

Proof of Lemma 2.4.31. Again, we use the Chebycheff bound:

$$\begin{aligned}
 &P\left(P_\omega^{ik}\left(\overline{T}_{(i+1)k_j} \neq T_{(i+1)k_j}\right) > d_{n_j}^\theta, \text{ some } i \in I_j\right) \\
 &\leq \frac{1}{d_{n_j}^\theta} \cdot \frac{3n_j}{k_j} \mathbb{P}^\circ\left(\overline{T}_{k_j} \neq T_{k_j}\right) \\
 &\leq \frac{1}{d_{n_j}^\theta} \cdot \frac{3n_j}{k_j} \cdot \exp(-k_j a(1 - \varepsilon/2)) \\
 &\leq 3 n_j^{1-\frac{1}{s}-\delta} \exp\left(n_j^{\frac{1}{s}+\delta} \left(\frac{a}{1-\varepsilon/4} + \theta\right)\right),
 \end{aligned}$$

where the second inequality follows again from (2.1.4) and the LDP for \overline{R}_k . The conclusion follows from the Borel-Cantelli lemma. \square

We need one more preliminary computation related to the bounds in (2.4.30). Let $\{Z_{k_j}^{(i)}\}$, $i = 1, 2, \dots$ denote a sequence of i.i.d. positive random variables, with

$$P\left(\frac{Z_{k_j}^{(i)}}{k_j} < \mu'\right) = 0, \quad P\left(\frac{Z_{k_j}^{(i)}}{k_j} > \mu'x\right) = (2b_{n_j})^{\lfloor x/2 \rfloor \vee 1}, \quad x \geq 1.$$

Note now that for any $\lambda > 0$, and any $\varepsilon > 0$,

$$\begin{aligned}
 E\left(\exp\left(\lambda \frac{Z_{k_j}^{(i)}}{k_j}\right)\right) &= \int_0^\infty P\left(\frac{Z_{k_j}^{(i)}}{k_j} > \frac{\log u}{\lambda}\right) du \\
 &\leq e^{\lambda\mu'(1+\varepsilon)} + \int_{e^{\lambda\mu'(1+\varepsilon)}}^\infty (2b_{n_j})^{\left\lceil \frac{\log u}{2\lambda\mu'(1+\varepsilon)} \right\rceil \vee 1} du \\
 &= e^{\lambda\mu'(1+\varepsilon)} + g_j.
 \end{aligned} \tag{2.4.32}$$

where $g_j \rightarrow_{j \rightarrow \infty} 0$.

In order to control the number of repetitions of visits to k_j -blocks, we introduce an auxiliary random walk. Let S_t , $t = 0, 1, \dots$, denote a simple random walk with $S_0 = 0$ and

$$P\left(S_{t+1} = S_t + 1 \mid S_t\right) = 1 - P\left(S_{t+1} = S_t - 1 \mid S_t\right) = 1 - d_n^\theta.$$

Set $M_{n_j} = \frac{1}{C_{n_j}} n_j^{1-\frac{1}{s}}$.

Lemma 2.4.33 *For θ as in Lemma 2.4.31, and n large enough,*

$$P\left(\inf\{t : S_t = \lfloor \frac{n_j}{k_j} \rfloor\} > M_{n_j}\right) \leq \exp\left(-\frac{\theta\varepsilon}{2} n_j\right).$$

Proof of Lemma 2.4.33.

$$\begin{aligned}
 P\left(\inf\left\{t : S_t = \left\lceil \frac{n_j}{k_j} \right\rceil\right\} > M_{n_j}\right) &\leq P\left(\frac{S_{\lceil M_{n_j} \rceil}}{M_{n_j}} < \frac{n_j}{k_j M_{n_j}}\right) \\
 &= P\left(\frac{S_{\lceil M_{n_j} \rceil}}{M_{n_j}} < 1 - \varepsilon\right) \leq 2e^{-M_{n_j} h_{n_j}(1-\varepsilon)},
 \end{aligned}$$

where the last inequality is a consequence of Chebycheff’s inequality and the fact that $d_n^\theta < \varepsilon$. Here,

$$h_n(1-x) = (1-x) \log\left(\frac{1-x}{1-d_n^\theta}\right) + x \log\frac{x}{d_n^\theta}.$$

Using $h_n(1-x) \geq -\frac{2}{\varepsilon} - x \log d_n^\theta$, we get

$$P\left(\frac{S_{\lceil M_{n_j} \rceil}}{M_{n_j}} < 1 - \varepsilon\right) \leq 2e^{2M_{n_j}/\varepsilon} e^{+\varepsilon M_{n_j} \log d_n^\theta} \leq e^{-\frac{\varepsilon}{2} \theta n_j}. \quad \square$$

We are now ready to prove (2.4.26). Note that, for all $j > J_0(\omega)$, and all $i \in I_j$, we may, due to (2.4.30), construct $\{Z_{k_j}^{(i)}\}$ and $\{\bar{T}_{(i+1)k_j}\}$ on the same probability space such that for all $i \in I_j$, $P_\omega^{ik}(Z_{k_j}^{(i)} \geq \bar{T}_{(i+1)k_j}) = 1$. Fix $1/v_P > 1/v' > 1/v$ and $\varepsilon > 0$ small enough. Recalling that under the law P_ω^o , the random variables $\bar{T}_{k_j}^{(i)} := \bar{T}_{(i+1)k_j} - T_{ik_j}$ are independent, we obtain, with $\{S_t\}$ defined before Lemma 2.4.33, and j large enough,

$$\begin{aligned}
 P_\omega^o(T_{n_j} > n_j/v) &\leq P\left(\inf\left\{t : S_t = \left\lceil \frac{n_j}{k_j} \right\rceil\right\} > M_{n_j}\right) + P\left(\sum_{i=1}^{M_{n_j}} Z_{k_j}^{(i)} > n_j/v\right) \\
 &\leq e^{-\theta \varepsilon n_j/2} + P\left(\frac{1}{M_{n_j}} \sum_{i=1}^{M_{n_j}} \frac{Z_{k_j}^{(i)}}{k_j} > 1/v(1-\varepsilon)\right) \\
 &\leq e^{-\theta \varepsilon n_j/2} + \left[E\left(\exp\left(\lambda \frac{Z_{k_j}^{(i)}}{k_j^{(i)}}\right)\right) \cdot e^{-\lambda(1-\varepsilon)/v}\right]^{M_{n_j}} \\
 &\leq e^{-\theta \varepsilon n_j/2} + \left(e^{\lambda(1/v'+2\varepsilon/v-1/v)} + g_j e^{-\lambda(1-\varepsilon)/v}\right)^{M_{n_j}} \\
 &\leq e^{-\theta \varepsilon n_j/2} + \left(e^{-\lambda \varepsilon/v}\right)^{M_{n_j}},
 \end{aligned}$$

where Lemma 2.4.33 was used in the second inequality and (2.4.32) in the fourth. Since $\lambda > 0$ is arbitrary, (2.4.26) follows. \square

Remarks: 1. A study of the proof of the annealed estimates shows that the strong mixing condition **(D3)** can be replaced by the slightly milder one that $\alpha(n) = \exp(-Cn)$ for some C large enough such that (2.4.15) holds, if one also assumes the existence of a LDP for \bar{R}_k . In this form, the assumption is satisfied for many Markov chains satisfying a Doeblin condition.

2. It is worthwhile noting that the transfer of the annealed estimates to the quenched setting required very few assumptions on the environment, besides the existence of a LDP for \overline{R}_k . This technique, as we will see, is not limited to the one-dimensional setup, and works well in situations where a drift is present.

3. One may study by similar techniques also the case where $E_P(\overline{S}) < \infty$ but $\rho_{\max} = 1$ with $\alpha := P(\rho_{\max} = 1) > 0$. The rate of decay is then quite different: at least when the environment is i.i.d., the annealed rate of decay in Theorem 2.4.3 is exponential with exponent $n^{1/3}$, see [18], whereas the quenched one has exponent $n/(\log n)^2$, see [30], and it seems both proofs extend to the mixing setup. By adapting the method of enlargement of obstacles to this setup, one actually can show more in the i.i.d. environment case: it holds then that,

$$\lim_{\delta \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n^{1/3}} \log \mathbb{P}^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) = -\frac{3}{2} \left| \frac{\pi \log \alpha}{2} \right|^{2/3}, \quad (2.4.34)$$

and

$$\lim_{\delta \rightarrow 0} \lim_{n \rightarrow \infty} \frac{(\log n)^2}{n} \log P^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) = -\frac{|\pi \log \alpha|^2}{8} \left(1 - \frac{v}{v_P} \right), \quad (2.4.35)$$

see [60] and [61]. (Note that the lower bounds in (2.4.34) and (2.4.35) are not hard to obtain, by constructing “neutral” traps. The difficulty lies in matching the constants in the upper bound to the ones in the lower bound.) The technique of enlargement of obstacles in this context is based on considerably refining the classification of blocks used above when going from annealed to quenched estimates, by introducing the notion of “good” and “bad” blocks (and double blocks...)

4. One can check, at least in the i.i.d. environment case, that when $\rho_{\max} = 1$ with $\alpha = 0$ then intermediate decay rates, between Theorems 2.4.3, 2.4.23 and (2.4.34), (2.4.35) can be achieved. We do not elaborate further here.

5. Again in the case of i.i.d. environment and the setup of Theorem 2.4.23, one can show, c.f. [30], that

$$\limsup_{n \rightarrow \infty} \frac{1}{n^{1-1/s}} \log P_\omega^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) = 0, \quad P - a.s. \quad (2.4.36)$$

This is due to fluctuations in the length of the “significant” trap where the walk may stay for large time. Based on the study of these fluctuations, it is reasonable to conjecture that

$$\liminf_{n \rightarrow \infty} \frac{1}{n^{1-1/s}} \log P_\omega^o \left(\frac{X_n}{n} \in (v - \delta, v + \delta) \right) = -\infty, \quad P - a.s.,$$

explaining the need for δ in the statement of Theorem 2.4.23. This conjecture has been verified only in the case where $P(\rho_{\min} = 0) > 0$, i.e. in the presence of “reflecting nodes”, c.f. [29, 28].

Bibliographical notes: The derivation in this section is based on [18] and [30]. Other relevant references, giving additional information not described here, are described in the remarks at the end of the section, so we only mention them here without repeating the description given there: [29, 60, 61].

2.5 Sinai’s model: non standard limit laws and aging properties

Throughout this section, define $\overline{R}_k = k^{-1} \sum_{i=1}^{k-1} \log \rho_i(\text{sign } i)$. We assume the following

Assumption 2.5.1

- (E1) Assumption 2.1.1 holds.
- (E2) $E_P \log \rho_0 = 0$, and there exists an $\varepsilon > 0$ such that $E_P |\log \rho_0|^{2+\varepsilon} < \infty$.
- (E3) P is strongly mixing, and the functional invariance principle holds for $\sqrt{k} \overline{R}_k / \sigma_P$; that is, $\{\sqrt{k} \overline{R}_{[kt]} / \sigma_P\}_{t \in \mathbb{R}}$ converges weakly to a Brownian motion for some $\sigma_P > 0$ (sufficient conditions for such convergence are as in Lemma 2.2.4).

(In the i.i.d. case, note that $\sigma_P^2 = E_P(\log \rho_0)^2$). Define

$$W^n(t) = \frac{1}{\log n} \sum_{i=0}^{\lfloor (\log n)^2 t \rfloor} \log \rho_i \cdot (\text{sign } t)$$

with $t \in \mathbb{R}$. By Assumption 2.5.1, $\{W^n(t)\}_{t \in \mathbb{R}}$ converges weakly to $\{\sigma_P B_t\}$, where $\{B_t\}$ is a two sided Brownian motion.

Next, we call a triple (a, b, c) with $a < b < c$ a valley of the path $\{W^n(\cdot)\}$ if

$$\begin{aligned} W^n(b) &= \min_{a \leq t \leq c} W^n(t), \\ W^n(a) &= \max_{a \leq t \leq b} W^n(t), \\ W^n(c) &= \max_{b \leq t \leq c} W^n(t). \end{aligned}$$

The *depth* of the valley is defined as

$$d_{(a,b,c)} = \min(W^n(a) - W^n(b), W^n(c) - W^n(b)).$$

If (a, b, c) is a valley, and $a < d < e < b$ are such that

$$W^n(e) - W^n(d) = \max_{a \leq x < y \leq b} W^n(y) - W^n(x)$$

then (a, d, e) and (e, b, c) are again valleys, which are obtained from (a, b, c) by a *left refinement*. One defines similarly a *right refinement*. Define

$$\begin{aligned}
 c_0^n &= \min\{t \geq 0 : W^n(t) \geq 1\} \\
 a_0^n &= \max\{t \leq 0 : W^n(t) \geq 1\} \\
 W^n(b_0^n) &= \min_{a_0^n \leq t \leq c_0^n} W^n(t).
 \end{aligned}$$

(b_0^n is not uniquely defined, however, due to Assumption 2.5.1, with P -probability approaching 1 as $n \rightarrow \infty$, all candidates for b_0^n are within distance converging to 0 as $n \rightarrow \infty$; we define b_0^n then as the smallest one in absolute value.)

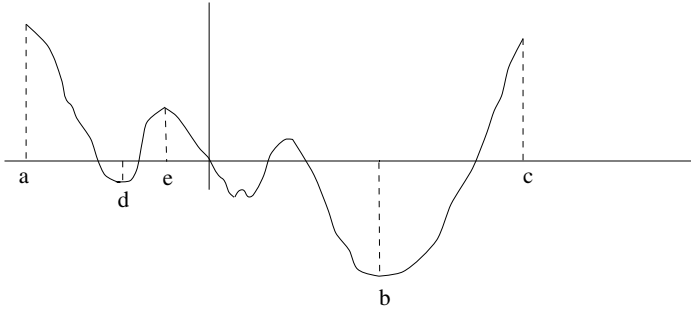


Fig. 2.5.1. Left refinement of (a, b, c)

One may now apply a (finite) sequence of refinements to find the *smallest* valley $(\bar{a}^n, \bar{b}^n, \bar{c}^n)$ with $\bar{a}^n < 0 < \bar{c}^n$, while $d_{(\bar{a}^n, \bar{b}^n, \bar{c}^n)} \geq 1$. We define similarly the smallest valley $(\bar{a}_\delta^n, \bar{b}_\delta^n, \bar{c}_\delta^n)$ such that $d_{(\bar{a}_\delta^n, \bar{b}_\delta^n, \bar{c}_\delta^n)} \geq 1 + \delta$. Let

$$A_n^{J,\delta} = \left\{ \begin{aligned}
 \omega \in \Omega : \bar{b}^n &= \bar{b}_\delta^n, \text{ any refinement } (a, b, c) \text{ of } (\bar{a}_\delta^n, \bar{b}_\delta^n, \bar{c}_\delta^n) \text{ with} \\
 b \neq \bar{b}^n &\text{ has depth } < 1 - \delta, |\bar{a}_\delta^n| + |\bar{c}_\delta^n| \leq J, \\
 \min_{t \in [\bar{a}^n, \bar{c}^n] \setminus [\bar{b}^n - \delta, \bar{b}^n + \delta]} &W^n(t) - W^n(\bar{b}^n) > \delta^3
 \end{aligned} \right\}$$

then it is easy to check by the properties of Brownian motion that

$$\lim_{\delta \rightarrow 0} \lim_{J \rightarrow \infty} \lim_{n \rightarrow \infty} P(A_n^{J,\delta}) = 1. \tag{2.5.2}$$

The following is the main result of this section:

Theorem 2.5.3 *Assume $P(\min(\omega_0^-, \omega_0^+) < \varepsilon) = 0$ and Assumption 2.5.1. For any $\eta > 0$,*

$$\mathbb{P}^\circ \left(\left| \frac{X_n}{(\log n)^2} - \bar{b}^n \right| > \eta \right) \xrightarrow{n \rightarrow \infty} 0.$$

Proof. Fix $\delta < \eta/2, J$ and n large enough with $\omega \in A_n^{J,\delta}$. For simplicity of notations, assume in the sequel that ω is such that $\bar{b}^n > 0$. Write

$a^n = \bar{a}^n(\log n)^2, b^n = \bar{b}^n(\log n)^2, c^n = \bar{c}^n(\log n)^2$, with similar notations for $a_\delta^n, b_\delta^n, c_\delta^n$. Define

$$\bar{T}_{b,n} = \min\{t \geq 0 : X_t = b^n \text{ or } X_t = a_\delta^n\}.$$

By (2.1.4),

$$P_\omega^o\left(X_{\bar{T}_{b,n}} = a_\delta^n\right) \leq \frac{1}{1 + \frac{\exp\{(\log n)(W^n(\bar{a}_\delta^n) - W^n(\bar{b}^n))\}}{Jn(\log n)^2}} \leq \frac{J(\log n)^2}{n^\delta}. \quad (2.5.4)$$

On the other hand, let $\tilde{T}_{b,n}$ have the law of $\bar{T}_{b,n}$ except that the walk $\{X_t\}$ is reflected at a_δ^n , and define similarly $\tilde{\tau}_1$. Using the same recursions as in (2.1.14), we have that

$$E_\omega^o(\tilde{\tau}_1) = \frac{1}{\omega_0^+} + \frac{\rho_0}{\omega_{(-1)}^+} + \cdots + \frac{\prod_{i=0}^{a_\delta^n+2} \rho_{-i}}{\omega_{a_\delta^n-1}^+} + \prod_{i=0}^{a_\delta^n+1} \rho_{-i}.$$

Hence, with $\tilde{\omega}_i = \omega_i$ for $i \neq a_\delta^n$ and $\tilde{\omega}_{a_\delta^n}^+ = 1$, for all n large enough,

$$\begin{aligned} E_\omega^o(\bar{T}_{b,n}) &\leq E_\omega^o(\tilde{T}_{b,n}) = \sum_{i=1}^{b^n} \sum_{j=0}^{i-1-a_\delta^n} \frac{\prod_{k=1}^j \rho_{i-k}}{\omega_{(i-j-1)}^+} \\ &\leq \frac{1}{\varepsilon} \sum_{i=1}^{b^n} \sum_{j=0}^{i-1-a_\delta^n} e^{(\log n)(W^n(i) - W^n(i-j))} \leq \frac{2J^2}{\varepsilon} e^{\log n(1-\delta)} \leq n^{1-\frac{\delta}{2}}. \end{aligned}$$

We thus conclude that

$$P_\omega^o\left(\bar{T}_{b,n} < n, \quad X_{\bar{T}_{b,n}} = b^n\right) \xrightarrow{n \rightarrow \infty} 1$$

implying that

$$P_\omega^o\left(T_{b^n} < n\right) \xrightarrow{n \rightarrow \infty} 1. \quad (2.5.5)$$

Next note that another application of (2.1.4) yields

$$\begin{aligned} P_\omega^{b^n-1}(X \text{ hits } b^n \text{ before } a_\delta^n) &\geq 1 - n^{-(1+\frac{\delta}{2})} \\ P_\omega^{b^n+1}(X \text{ hits } b^n \text{ before } c_\delta^n) &\geq 1 - n^{-(1+\frac{\delta}{2})}. \end{aligned} \quad (2.5.6)$$

On the same probability space, construct a RWRE $\{\tilde{X}_t\}$ with the same transition mechanism as $\{X_t\}$ except that it is reflected at a_δ^n , i.e. replace ω by $\tilde{\omega}$. Then, using (2.5.6),

$$\begin{aligned} P_\omega^o\left(\left|\frac{X_n}{(\log n)^2} - \bar{b}^n\right| > \delta\right) &\leq P_\omega^o(T_{b^n} > n) + \max_{t \leq n} P_\omega^{b^n}\left(\left|\frac{X_t}{(\log n)^2} - \bar{b}^n\right| > \delta\right) \\ &\leq P_\omega^o(T_{b^n} > n) + \left[1 - (1 - n^{-(1+\frac{\delta}{2})})^n\right] \\ &\quad + \max_{t \leq n} P_\omega^{b^n}\left(\left|\frac{\tilde{X}_t}{(\log n)^2} - \bar{b}^n\right| > \delta\right). \end{aligned}$$

Hence, in view of (2.5.2) and (2.5.5), the theorem holds as soon as we show that

$$\sup_{\omega \in A_n^{J,\delta}} \max_{t \leq n} P_\omega^{\bar{b}^n} \left(\left| \frac{\tilde{X}_t}{(\log n)^2} - \bar{b}^n \right| > \delta \right) \xrightarrow{n \rightarrow \infty} 0. \tag{2.5.7}$$

To see (2.5.7), define

$$f(z) = \frac{\prod_{a_\delta^n + 1 \leq i < z} \omega_i^+}{\prod_{a_\delta^n + 1 \leq i < z} \omega_{i+1}^-}, \quad \bar{f}(z) = \frac{f(z)}{f(b^n)}$$

(as usual, the product over an empty set of indices is taken as 1. $\bar{f}(\cdot)$ corresponds to the invariant measure for the resistor network corresponding to \tilde{X}). Next, define the operator

$$(Ag)(z) = \bar{\omega}_{z-1}^+ g(z-1) + \bar{\omega}_{z+1}^- g(z+1) + \bar{\omega}_z^0 g(z) \tag{2.5.8}$$

where $\bar{\omega}_z = \omega_z$ for $z > a_\delta^n$, $\bar{\omega}_{a_\delta^n}^+ = 1, \bar{\omega}_{a_\delta^n-1}^+ = 0$. Note that $A\bar{f} = \bar{f}$, and further that

$$P_\omega^{b^n}(\tilde{X}_t = z) = A^t \mathbf{1}_{b^n}(z).$$

Since $\bar{f}(z) \geq \mathbf{1}_{b^n}(z)$ and A is a positive operator, we conclude that

$$P_\omega^{b^n}(\tilde{X}_t = z) \leq \bar{f}(z).$$

But, for z with $|z/(\log n)^2 - \bar{b}^n| > \delta$, it holds that $\bar{f}(z) \leq e^{-\delta^3 \log n}$, and hence

$$P_\omega^{b^n}(\tilde{X}_t = z) \leq n^{-\delta^3}.$$

Thus, for $\omega \in A_n^{J,\delta}$, using the fact that the second inequality in (2.5.6) still applies for \tilde{X} ,

$$\max_{t \leq n} P_\omega^{\bar{b}^n} \left(\left| \frac{\tilde{X}_t}{(\log n)^2} - \bar{b}^n \right| > \delta \right) \leq (\bar{b}^n + \delta)(\log n)^2 n^{-\delta^3} + 1 - \left(1 - n^{-(1+\delta/2)}\right)^n,$$

yielding (2.5.7) and completing the proof of the theorem. □

We next turn to a somewhat more detailed study of the random variable \bar{b}^n . By replacing 1 with t in the definition of \bar{b}^n , one obtains a process $\{\bar{b}^n(t)\}_{t \geq 0}$. Further, due to Assumption 2.5.1, the process $\{\bar{b}^n(t/\sigma_P)\}_{t \geq 0}$ converges weakly to a process $\{\bar{b}(t)\}_{t \geq 0}$, defined in terms of the Brownian motion $\{B_t\}_{t \geq 0}$; Indeed, $\bar{b}(t)$ is the location of the bottom of the smallest valley of $\{B_t\}_{t \geq 0}$, which surrounds 0 and has depth t . Throughout this section we denote by \mathcal{Q} the law of the Brownian motion B . Our next goal is to characterize the process $\{\bar{b}(t)\}_{t \geq 0}$. Toward this end, define

$$\begin{aligned} m_+(t) &= \min\{B_s : 0 \leq s \leq t\}, \quad m_-(t) = \min\{B_{-s} : 0 \leq s \leq t\} \\ T_+(a) &= \inf\{s \geq 0 : B_s - m_+(s) = a\}, \\ T_-(a) &= \inf\{s \geq 0 : B_{-s} - m_-(s) = a\} \\ s_\pm(a) &= \inf\{s \geq 0 : m_\pm(T_\pm(a)) = B_{\pm s}\}, \\ M_\pm(a) &= \sup\{B_{\pm \eta} : 0 \leq \eta \leq s_\pm(a)\}. \end{aligned}$$

Next, define $W_{\pm}(a) = B_{s_{\pm}(a)}$. It is not hard to check that the pairs $(M_+(\cdot), W_+(\cdot))$ and $(M_-(\cdot), \bar{W}_-(\cdot))$ form independent Markov processes. Define finally

$$H_{\pm}(a) = (W_{\pm}(a) + a) \vee M_{\pm}(a).$$

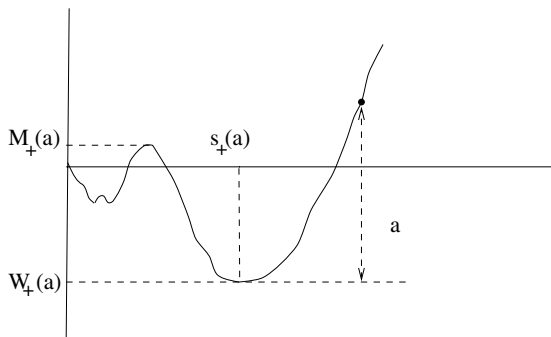


Fig. 2.5.2. The random variables $(M_+(a), W_+(a), s_+(a))$

We now have the

Theorem 2.5.9 For each $a > 0$, $\mathcal{Q}(\bar{b}(a) \in \{s_+(a), -s_-(a)\}) = 1$. Further, $\bar{b}(a) = s_+(a)$ iff $H_+(a) < H_-(a)$.

Proof. Note that $\mathcal{Q}(H_+(a) = H_-(a)) = 0$. That $\bar{b}(a) \in \{s_+(a), -s_-(a)\}$ is a direct consequence of the definitions, i.e. assuming $\bar{b}(a) > 0$ and $\bar{b}(a) \neq s_+(a)$ it is easy to show that one may refine from the right the valley defining $\bar{b}(a)$, contradicting minimality. We begin by showing, after Kesten [41], that $\bar{b}(a) = s_+(a)$ iff either

$$W_-(a) > W_+(a), \quad M_+(a) < (W_-(a) + a) \vee M_-(a) \tag{2.5.10}$$

or

$$W_-(a) < W_+(a), \quad M_-(a) > (W_+(a) + a) \vee M_+(a). \tag{2.5.11}$$

Indeed, assume $\bar{b}(a) = s_+(a)$, and $W_-(a) > W_+(a)$. Let $(\alpha, \bar{b}(a), \gamma)$ denote the minimal valley defining $\bar{b}(a)$. If $-s_-(a) \leq \alpha$, then

$$\begin{aligned} M_-(a) &= \max\{B_{-s} : s \in (0, s_-(a))\} \geq B_{-\alpha} \\ &= \max\{B_s : -\alpha \leq s \leq \bar{b}(a)\} \geq M_+(a) \end{aligned} \tag{2.5.12}$$

implying (2.5.10). On the other hand, if $-s_-(a) > \alpha$, refine $(\alpha, \bar{b}(a), \gamma)$ on the left (find α', β' with $\alpha < \alpha' < \beta' < \bar{b}(a)$), such that

$$B_{\beta'} - B_{\alpha'} = \max_{\alpha < x < y < \bar{b}(a)} (B_y - B_x) \geq M_+(a) - W_-(a)$$

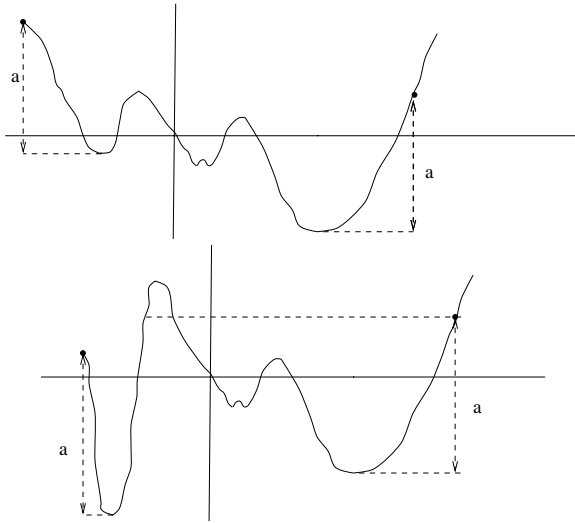


Fig. 2.5.3. $\bar{b}(a) = s_+(a)$

and thus minimality of $(\alpha, \bar{b}(a), \gamma)$ implies that $M_+(a) - W_-(a) < a$, implying (2.5.10).

We thus showed that if $\bar{b}(a) = s_+(a)$ and $W_-(a) > W_+(a)$ then (2.5.10) holds. On the other hand, if (2.5.10) holds, we show that $\bar{b}(a) = s_+(a)$ by considering the cases $\alpha \leq -s_-(a)$ and $-s_-(a) < \alpha$ separately. In the former case, necessarily $\gamma > s_+$, for otherwise $M_-(\alpha) \leq B_\gamma \leq M_+(a) \leq W_-(a) + a$ which together with $\bar{b}(a) = -s_-(a)$ would imply that the depth of $(\alpha, \bar{b}(a), \gamma)$ is smaller than a . Thus, under (2.5.10) if $\alpha \leq -s_-(a)$ then $\gamma > s_+$, and in this case $\bar{b}(a) = s_+(a)$ since $B_{s_+(a)} < B_{-s_-(a)}$. Finally, if $\alpha > -s_-(a)$ then $\bar{b}(a) \neq -s_-(a)$ and hence $\bar{b}(a) = s_+(a)$.

Hence, we showed that if $W_-(a) > W_+(a)$ then (2.5.10) is equivalent to $\bar{b}(a) = s_+(a)$. Interchanging the positive and negative axis, we conclude that if $W_-(a) < W_+(a)$, then $\bar{b}(a) = -s_-(a)$ iff $M_+(a) < (W_+(a) + 1) \vee M_-(a)$. This completes the proof that $\bar{b}(a) = s_+(a)$ is equivalent to (2.5.10) or (2.5.11).

To complete the proof of the theorem, assume first $W_-(a) > W_+(a)$. Then, $\bar{b}(a) = s_+(a)$ iff (2.5.10) holds, i.e. $M_+(a) < (W_-(a) + a) \vee M_-(a) = H_-(a)$. But $H_-(a) \geq W_-(a) + a \geq W_+(a) + a$, and hence $M_+(a) < H_-(a)$ is equivalent to $M_+(a) \vee (W_+(a) + a) < H_-(a)$, i.e. $H_+(a) < H_-(a)$. The case $W_+(a) < W_-(a)$ is handled similarly by using (2.5.11). \square

One may use the representation in Theorem 2.5.9 in order to evaluate explicitly the law of $\bar{b}(a)$ (note that $\bar{b}(a) \stackrel{L}{\sim} a^2 \bar{b}(1)$ by Brownian scaling). This is done in [41], and we do not repeat the construction here. Our goal is to use Theorem 2.5.9 to show that Sinai’s model exhibits *aging* properties. More precisely, we claim that

Theorem 2.5.13 *Assume $P(\min(\omega_0, \omega_0^+) < \varepsilon) = 0$ and Assumption 2.5.1. Then, for $h > 1$,*

$$\lim_{\eta \rightarrow 0} \lim_{n \rightarrow \infty} \mathbb{P}^o \left(\frac{|X_{n^h} - X_n|}{(\log n)^2} < \eta \right) = \frac{1}{h^2} \left[\frac{5}{3} - \frac{2}{3} e^{-(h-1)} \right]. \tag{2.5.14}$$

Proof. Applying Theorem 2.5.3, the limit in the left hand side of (2.5.14) equals

$$\mathcal{Q}(\bar{b}(h) = \bar{b}(1)) = 2\mathcal{Q}(\bar{b}(h) = \bar{b}(1) = s_+(1) = s_+(h)).$$

Note that

$$\mathcal{Q}(s_+(h) = s_+(1)) = \mathcal{Q} \left(\begin{array}{c} \text{Brownian motion, started at height 1,} \\ \text{hits } h \text{ before hitting 0} \end{array} \right) = \frac{1}{h}.$$

Hence, using that on $s_+(1) = s_+(h)$ one has $W_+(1) = W_+(h), M_+(1) = M_+(h)$, and using that the event $\{s_+(h) = s_+(1)\}$ depends only on increments of the path of the Brownian motion after time $T_+(1)$, one gets

$$\mathcal{Q}(\bar{b}(h) = \bar{b}(1)) = \frac{2}{h} \mathcal{Q}(H_+(1) < H_-(1), (W_+(1) + h) \vee M_+(1) < H_-(h)). \tag{2.5.15}$$

Next, let

$$\begin{aligned} \tau_0 &= \min\{t > s_-(1) : B_{-t} = W_-(1) + 1\} \\ \tau_h &= \min\{t > \tau_0 : B_{-t} = W_-(1) + h \text{ or } B_t = W_-(1)\}. \end{aligned}$$

Note that $\tau_h - \tau_0$ has the same law as that of the hitting time of $\{0, h\}$ by a Brownian motion Z_t with $Z_0 = 1$. (Here, $Z_t = B_{-(\tau_0+t)} - W_-(1)$!). Further, letting $I_h = \mathbf{1}_{\{B_{\tau_h} = W_-(1)\}} (= \mathbf{1}_{\{Z_{\tau_h - \tau_0} = 0\}})$, it holds that

$$\begin{aligned} W_-(h) &= W_-(1) + I_h \tilde{W}_-(h) \\ M_-(h) &= \begin{cases} M_-(1), & I_h = 0 \\ M_-(1) \vee (\overline{M}_-(h) + W_-(1) + 1) \vee (\tilde{M}_-(h) + W_-(1)), & I_h = 1 \end{cases} \end{aligned}$$

where $(\tilde{W}_-(h), \tilde{M}_-(h))$ are independent of $(W_-(1), M_-(1))$ and possess the same law as $(W_-(h), M_-(h))$, while $\overline{M}_-(h)$ is independent of both $(W_-(1), M_-(1))$ and $(\tilde{W}_-(h), \tilde{M}_-(h))$ and has the law of the maximum of a Brownian motion, started at 0, killed at hitting -1 and conditioned not to hit $h - 1$. (See figure 2.5.4 for a graphical description of these random variables.)

Set now

$$\hat{M}_-(h) = \begin{cases} h, & I_h = 0 \\ 1 + \overline{M}_-(h), & I_h = 1, \end{cases}$$

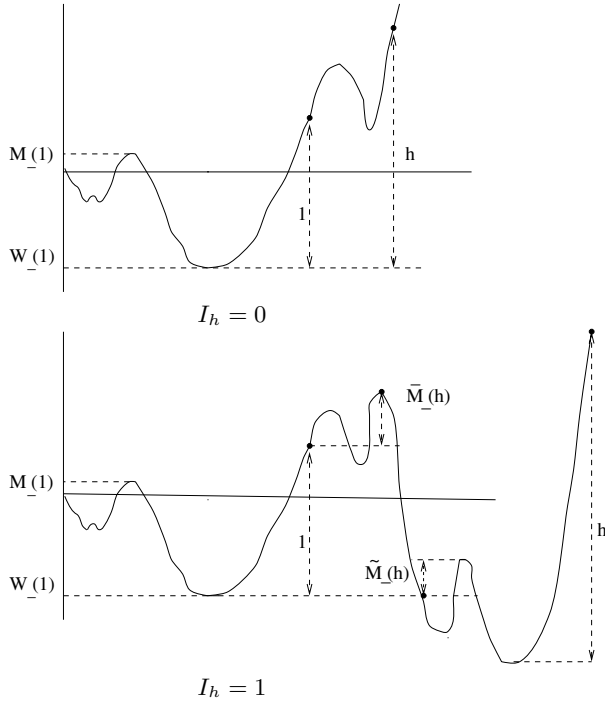


Fig. 2.5.4. Definition of auxiliary variables

$\tilde{H}_-(h) = (\tilde{W}_-(h) + h) \vee \tilde{M}_-(h)$ and $\Gamma(h) = \max(\tilde{H}_-(h), \hat{M}_-(h))$. Note that $\tilde{H}_-(h)$ has the same law as $H_-(h)$ but is independent of $\bar{M}_-(h)$. Further, it is easy to check that $(W_-(h) + h) \vee M_-(h) = (W_-(1) + \Gamma(h)) \vee M_-(1)$ (note that either $M_-(h) = M_-(1)$ or $M_-(h) > M_-(1)$ but in the latter case, $M_-(h) \leq W_-(1) + \Gamma(h)$.) We have the following lemma, whose proof is deferred:

Lemma 2.5.16 *The law of $\Gamma(h)$ is $\frac{1}{h}\delta_h + \frac{h-1}{h}U[1, h]$, where $U[1, h]$ denotes the uniform law on $[1, h]$.*

Substituting in (2.5.15), we get that

$$\mathcal{Q}(\bar{b}(h) = \bar{b}(1)) = \mathcal{Q}(E_{\mathcal{Q}}(\bar{b}(h) = \bar{b}(1)|\Gamma(h))) = \frac{2}{h^2} \left[\int_1^h \mathcal{Q}(t)dt + \mathcal{Q}(h) \right] \tag{2.5.17}$$

where

$$\mathcal{Q}(t) = \mathcal{Q}(H_+(1) < H_-(1), H_+(h) < H_-(t) | s_+(h) = s_+(1), s_-(1) = s_-(t)) .$$

In order to evaluate the integral in (2.5.17), we need to evaluate the joint law of $(H_+(1), H_+(t))$ (the joint law of $(H_-(1), H_-(t))$ being identical). Since

$0 \leq H_+(1) \leq 1$ and $H_+(1) \leq H_+(t) \leq H_+(1) + t - 1$, the support of the law of $(H_+(1), H_+(t))$ is the domain A defined by $0 \leq x \leq 1, x \leq y \leq x + t - 1$. Note that for $(z, w) \in A$,

$$\begin{aligned} & \mathbb{Q}(H_+(1) \leq z, H_+(t) \leq w \mid s_+(1) = s_+(h)) \\ &= \mathbb{Q}(M_+(1) \leq z \wedge w, W_+(1) \leq -[(1 - z) \vee (t - w)]) \\ &= \mathbb{Q}\left(M_+(1) \leq z, W_+(1) \leq -(t - w)\right). \end{aligned}$$

We now have the following well known lemma. For completeness, the proof is given at the end of this section:

Lemma 2.5.18 *For $z + y \geq 1, 0 \leq z \leq 1, y \geq 0$,*

$$\mathbb{Q}(M_+(1) \leq z, W_+(1) \leq -y) = ze^{-(z+y-1)}.$$

Lemma 2.5.18 implies that, for $(z, w) \in A, t > 1$,

$$\mathbb{Q}(H_+(1) \leq z, H_+(t) \leq w \mid s_+(1) = s_+(h)) = ze^{-(z+t-w-1)}. \tag{2.5.19}$$

Denote by B_1 the segment $\{0 \leq x = y \leq 1\}$ and by B_2 the segment $\{t - 1 \leq y = x + t - 1 \leq t\}$. We conclude, after some tedious computations, that the conditional law of $(H_+(1), H_+(t))$:

- possesses the density $f(z, \omega) = (1 - z)e^{-z}e^{-w-(t-1)}, (z, w) \in A \setminus (B_1 \cup B_2)$
- possesses the density $\tilde{f}(z, y) = (1 - z)e^{-(t-1)}, z = w \in B_1$
- possesses the density $\tilde{f}(z, z + t - 1) = z, w = z + t - 1 \in B_2$.

Substituting in the expression for $Q(t)$, we find that

$$Q(t) = \frac{5}{12}e^{-(h-t)} + \frac{1}{12}e^{-(h+t-2)}.$$

Substituting in (2.5.19), the theorem follows. □

Proof of Lemma 2.5.16: Note that $\mathbb{Q}(I_h = 0) = 1/h$, and in this case $\Gamma_h = h$. Thus, we only need to consider the case where $I_h = 1$ and show that under this conditioning, $\max(H_-(h), 1 + \overline{M}_-(h))$ possesses the law $U[1, h]$. Note that by standard properties of Brownian motion,

$$\mathbb{Q}(\hat{M}_-(h) \leq \xi \mid I_h = 1) = \frac{\xi-1}{\frac{\xi}{h-1}}.$$

We show below that the law of $\tilde{H}_-(h)$, which is identical to the law of $H_-(h)$, is uniform on $[0, h]$. Thus, using independence, for $\xi \in [1, h]$,

$$\mathbb{Q}(\Gamma_h < \xi \mid I_h = 1) = \frac{h(\xi - 1)\xi}{\xi(h - 1)h} = \frac{\xi - 1}{h - 1},$$

i.e. the law of Γ_h conditioned on $I_h = 1$ is indeed $U[1, h]$.

It thus only remains to evaluate the law of $H_-(h)$. By Brownian scaling, the law of $H_-(h)$ is identical to the law of $hH_+(1)$, so we only need show that the law of $H_+(1)$ is uniform on $[0, 1]$. This in fact is a direct consequence of Lemma 2.5.18. \square

Proof of Lemma 2.5.18: Let \mathbb{Q}^x denote the law of a Brownian motion $\{Z_t\}$ starting at time 0 at x . The Markov property now yields

$$\begin{aligned} \mathbb{Q}(M_+(1) \leq z, W_+(1) \leq -y) &= \mathbb{Q}^\circ(\{Z_t\} \text{ hits } z - 1 \text{ before hitting } z) \\ &= \mathbb{Q}^{z-1}(M_+(1) \leq z, W_+(1) \leq -y) \\ &= z\mathbb{Q}^\circ(M_+(1) \leq 1, W_+(1) \leq -y - z + 1) \\ &= z\mathbb{Q}^\circ(W_+(1) \leq -(y + z - 1)). \end{aligned} \tag{2.5.20}$$

For $x \geq 0$, let $f(x) := \mathbb{Q}(W_+(1) \leq -x)$. The Markov property now implies

$$f(x + \epsilon) = f(x)\mathbb{Q}^{-x}(W_+(1) \leq -(x + \epsilon)) = f(x)f(\epsilon).$$

Since $f(0) = 1$ and $f(\epsilon) = 1 - \epsilon + o(\epsilon)$, it follows that $f(x) = e^{-x}$. Substituting in (2.5.20), the lemma follows. \square

Bibliographical notes: Theorem 2.5.3 is due to [66]. The proof here follows the approach of Golosov [31], who dealt with a RWRE reflected at 0, i.e. with state space \mathbb{Z}_+ . In the same paper, Golosov evaluates the analogue of Theorem 2.5.9 in this reflected setup, and in [32] he provides sharp (pathwise) localization results. These are extended to the case of a walk on \mathbb{Z} in [33]. The statement of Theorem 2.5.9 and the proof here follow the article [41], where an explicit characterization of the law of $\bar{b}(1)$ is provided. The same characterization appears also in [33]. The aging properties of RWRE (Theorem 2.5.13) were first derived heuristically in [24], to which we refer for additional aging properties and discussion. The derivation here is based on [17]. The right hand side of formula (2.5.14) appears also in [33], in a slightly different context. We mention that results of iterated logarithm types, and results concerning most visited sites for Sinai’s RWRE, can be found in [35], [36]. See [65] for a recent review. Finally, extensions of the results in this section and a theorem concerning the dichotomy between Sinai’s regime and the classical CLT for ergodic environments can be found in [7].

Limit laws for transient RWRE in an i.i.d. environment appear in [42]. One distinguishes between CLT limit laws and stable laws: recall the parameter s introduced in Section 2.4. The main result of [42] is that if $s > 2$, a CLT holds true (see Section 2.2 for other approaches), whereas for $s \in (0, 2)$ a Stable(s) limit law holds true. Note that this is valid even when $s < 1$, i.e. when $v_P = 0$! It is an interesting open problem to extend the results concerning stable limit laws to non i.i.d. environments. Some results in this direction are forthcoming in the Technion thesis of A. Roitershtein.

3 RWRE – $d > 1$

3.1 Ergodic Theorems

In this section we present some of the general results known concerning 0 – 1 laws and laws of large numbers for nearest neighbour RWRE in \mathbb{Z}^d . Even if considerable progress was achieved in recent years, the situation here is, unfortunately, much less satisfying than for $d = 1$.

A standing assumption throughout this section is the following:

Assumption 3.1.1

(A1) P is stationary and ergodic, and satisfies a ϕ -mixing condition: there exists a function $\phi(l) \xrightarrow{l \rightarrow \infty} 0$ such that any two l -separated events A, B with $P(A) > 0$,

$$\left| \frac{P(A \cap B)}{P(A)} - P(B) \right| \leq \phi(l).$$

(A2) P is uniformly elliptic: there exists an $\varepsilon > 0$ such that

$$P(\omega(0, e) \geq \varepsilon) = 1, \quad \forall e \in \{\pm e_i\}_{i=1}^d.$$

(Events A, B are l -separated if the shortest lattice path connecting A and B is of length l or more.)

Remark: I have recently learnt that Assumption (A1) implies, in fact, that P is finitely dependent, c.f. [5]. On the other hand, the basic structure of what appears in the rest of this section remains unchanged if P is *mixing on cones*, see [13], and thus I have kept the proof in its original form.

Fix $\ell \in \mathbb{R}^d \setminus \{0\}$, and consider the events

$$A_{\pm\ell} = \left\{ \lim_{n \rightarrow \infty} X_n \cdot \ell = \pm\infty \right\}.$$

We have the

Theorem 3.1.2 *Assume Assumption 3.1.1. Then*

$$\mathbb{P}^\circ(A_\ell \cup A_{-\ell}) \in \{0, 1\}.$$

Proof. We begin by constructing an extension of our probability space: recall that the RWRE was defined by means of the law $\mathbb{P}^\circ = P \otimes P_\omega^\circ$ on $(\Omega \times (\mathbb{Z}^d)^\mathbb{N}, \mathcal{F} \times \mathcal{G})$. Set $W = \{0\} \cup \{\pm e_i\}_{i=1}^d$ and \mathcal{W} the cylinder σ -algebra on $W^\mathbb{N}$. We now define the measure

$$\overline{\mathbb{P}}^\circ = P \otimes Q_\varepsilon \otimes \overline{P}_{\omega, \varepsilon}^\circ$$

on

$$\left(\Omega \times W^\mathbb{N} \times (\mathbb{Z}^d)^\mathbb{N}, \mathcal{F} \times \mathcal{W} \times \mathcal{G} \right)$$

in the following way: Q_ε is a product measure, such that with $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots)$ denoting an element of $W^{\mathbb{N}}$, $Q_\varepsilon(\varepsilon_1 = \pm e_i) = \varepsilon$, $i = 1, \dots, d$, $Q_\varepsilon(\varepsilon_1 = 0) = 1 - 2\varepsilon d$. For each fixed ω, ε , $\overline{P}_{\omega, \varepsilon}^o$ is the law of the Markov chain $\{X_n\}$ with state space \mathbb{Z}^d , such that $X_0 = 0$ and, for each $e \in W$, $e \neq 0$,

$$\overline{P}_{\omega, \varepsilon}^o(X_{n+1} = z + e | X_n = z) = \mathbf{1}_{\{\varepsilon_{n+1}=e\}} + \frac{\mathbf{1}_{\{\varepsilon_{n+1}=0\}}}{1 - 2d\varepsilon} [\omega(z, z + e) - \varepsilon].$$

It is not hard to check that the law of $\{X_n\}$ under $\overline{\mathbb{P}}^o$ coincides with its law under \mathbb{P}^o , while its law under $Q_\varepsilon \otimes \overline{P}_{\omega, \varepsilon}^o$ coincides with its law under P_ω^o .

We will prove the theorem for $\ell = (1, 0 \dots 0)$, the general case being similar but requiring more cumbersome notations. Note that for any $u < v$, the walk cannot visit infinitely often the strip $u \leq z \cdot \ell \leq v$ without crossing the line $z \cdot \ell = v$. More precisely, with

$$T_v = \inf\{n \geq 0 : X_n \cdot \ell \geq v\}, \tag{3.1.3}$$

we have

$$\mathbb{P}^o(\#\{n > 0 : X_n \cdot \ell \geq u\} = \infty, T_v = \infty) = 0. \tag{3.1.4}$$

Indeed, note that for any z with $u \leq z \cdot \ell \leq v$, and any ω ,

$$P_\omega^z(X_{v-u} \cdot \ell \geq v) = Q_\varepsilon \otimes P_{\omega, \varepsilon}^z(X_{v-u} \cdot \ell \geq v) \geq \varepsilon^{v-u},$$

yielding (3.1.4) by the strong Markov property.

Assume next that $\mathbb{P}^o(A_\ell) > 0$. Set $D = \inf\{n \geq 0 : X_n \cdot \ell < X_0 \cdot \ell\}$. Clearly, $\mathbb{P}^o(D = \infty) > 0$, because if $\mathbb{P}^o(D = \infty) = 0$ then $\mathbb{P}^z(D < \infty) = 1 \forall z \in \mathbb{Z}^d$, and thus P -a.s., for all $z \in \mathbb{Z}^d$, $P_\omega^z(D < \infty) = 1$. This implies by the Markov property that

$$\liminf_{n \rightarrow \infty} X_n \cdot \ell \leq 0, \quad \mathbb{P}^o\text{-a.s.},$$

contradicting $\mathbb{P}^o(A_\ell) > 0$.

Define \mathcal{O}_ℓ to be the event that $X_n \cdot \ell$ changes its sign infinitely often. We next show that whenever $\mathbb{P}^o(A_\ell) > 0$, then $\mathbb{P}^o(\mathcal{O}_\ell) = 0$. Set $M = \sup_n X_n \cdot \ell$, fix $v > 0$ and note by (3.1.4) that

$$\mathbb{P}^o(\mathcal{O}_\ell \cap \{M < v\}) = 0. \tag{3.1.5}$$

We next prove that if $\mathbb{P}^o(A_\ell) > 0$ then $\mathbb{P}^o(\mathcal{O}_\ell \cap \{M = \infty\}) = 0$, by first noting that

$$\mathbb{P}^o(\mathcal{O}_\ell \cap \{M = \infty\}) = \overline{\mathbb{P}}^o(\mathcal{O}_\ell \cap \{M = \infty\}).$$

Then, set $\mathcal{G}_n = \sigma((\varepsilon_i, X_i), i \leq n)$, fix $L > 0$ and, setting $S_0 = 0$, define recursively \mathcal{G}_n stopping times as follows:

$$\begin{aligned} R_k &= \inf\{n \geq S_k : X_n \cdot \ell < 0\}, \\ S_{k+1} &= \inf\{n \geq R_k : X_{n-L} \cdot \ell \\ &\geq \max\{X_m \cdot \ell : m \leq n - L\}, \varepsilon_{n-1} = \varepsilon_{n-2} = \dots = \varepsilon_{n-L} = e_1\}. \end{aligned}$$

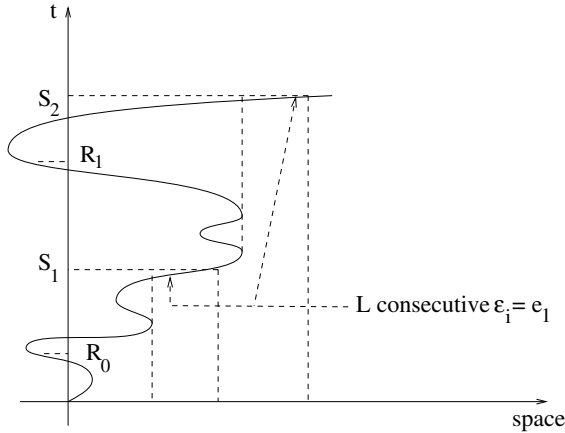


Fig. 3.1.1. Definition of the hitting times (S_k, R_k)

On $\mathcal{O}_\ell \cap \{M = \infty\}$, all these stopping times are finite. Now, at each time $S_k - L$ the walk enters a half space it never visited before, and then *due to the action of the ϵ sequence alone*, it proceeds L steps in the direction e_1 . Formally, “events in the σ -algebra \mathcal{G}_{S_k} are L -separated from $\sigma(\omega_z : z \cdot \ell \geq X_{S_k} \cdot \ell)$ ”. Note that, using $\mathbb{P}^o(A_\ell) > 0$ in the second inequality,

$$\overline{\mathbb{P}}^o(R_0 < \infty) = \overline{\mathbb{P}}^o(D < \infty) < 1,$$

whereas, using θ to denote both time and space shifts as needed from the context,

$$\begin{aligned} \overline{\mathbb{P}}^o(R_1 < \infty) &\leq \overline{\mathbb{P}}^o(R_0 < \infty, R_0 \circ \theta_{X_{S_1}} < \infty) \\ &= \sum_{z \in \mathbb{Z}^d} \overline{\mathbb{P}}^o(R_0 < \infty, R_0 \circ \theta_z < \infty, X_{S_1} = z) \\ &= \sum_{z \in \mathbb{Z}^d} \sum_{n \in \mathbb{N}} E_{P \otimes Q_\epsilon} \left(\overline{\mathbb{P}}_{\omega, \epsilon}^o(R_0 < \infty, X_{S_1} = z, S_1 = n) \cdot \overline{\mathbb{P}}_{\theta^z \omega, \theta^n \epsilon}^o(R_0 < \infty) \right). \end{aligned}$$

Note that $\overline{\mathbb{P}}_{\theta^z \omega, \theta^n \epsilon}^o(R_0 < \infty)$ is measurable on $\sigma(\omega_x : x \cdot \ell \geq z \cdot \ell) \times \sigma(\epsilon_i, i > n)$, whereas $\overline{\mathbb{P}}_{\omega, \epsilon}^o(R_0 < \infty, X_{S_1} = z, S_1 = n)$ is measurable on $\sigma(\omega_x : x \cdot \ell \leq z \cdot \ell - L) \times \sigma(\epsilon_i, i \leq n)$. Hence, by the ϕ -mixing property of P and the product structure of Q_ϵ ,

$$\begin{aligned} \bar{\mathbb{P}}^\circ(R_1 < \infty) &\leq \sum_{z \in \mathbb{Z}^d} \sum_{n \in \mathbb{N}} \left[E_{P \otimes Q_\varepsilon} \left(\bar{P}_{\omega, \varepsilon}^\circ(R_0 < \infty, X_{S_1} = z, S_1 = n) \right) \right. \\ &\quad \left. \cdot E_{P \otimes Q_\varepsilon} \left(\bar{P}_{\omega, \varepsilon}^\circ(R_0 < \infty) \right) \right] \\ &+ \phi(L) \sum_{z \in \mathbb{Z}^d} \sum_{n \in \mathbb{N}} E_{P \otimes Q_\varepsilon} \left(\bar{P}_{\omega, \varepsilon}^\circ(R_0 < \infty, X_{S_1} = z, S_1 = n) \right) \\ &\leq (\bar{\mathbb{P}}^\circ(R_0 < \infty))^2 + \phi(L) \bar{\mathbb{P}}^\circ(R_0 < \infty) \leq (\bar{\mathbb{P}}^\circ(D < \infty) + \phi(L))^2. \end{aligned}$$

Repeating this procedure, we conclude that $\bar{\mathbb{P}}^\circ(\mathcal{O}_\ell \cap \{M = \infty\}) \leq \bar{\mathbb{P}}^\circ(R_k < \infty) \leq (\mathbb{P}^\circ(D < \infty) + \phi(L))^{k+1}$. Since k is arbitrary and $\phi(L) \xrightarrow{L \rightarrow \infty} 0$, we conclude that $\bar{\mathbb{P}}^\circ(\mathcal{O}_\ell \cap \{M = \infty\}) = 0$, yielding with the above that $\mathbb{P}^\circ(\mathcal{O}_\ell) = 0$ as soon as $\mathbb{P}^\circ(A_\ell) > 0$. In a similar manner one proves that $\mathbb{P}^\circ(A_{-\ell}) > 0$ also implies $\mathbb{P}^\circ(\mathcal{O}_\ell) = 0$.

Assume now $1 > \mathbb{P}^\circ(A_\ell \cup A_{-\ell})$. Then one can find a v such that $\mathbb{P}^\circ(X_n \cdot \ell \in [-v, v] \text{ infinitely often}) > 0$. Therefore, $\mathbb{P}^\circ(\mathcal{O}_\ell) > 0$, implying by the above $\mathbb{P}^\circ(A_\ell) = \mathbb{P}^\circ(A_{-\ell}) = 0$. □

Remark: It should be obvious that one does not need the full strength of **(A1)** in Assumption 3.1.1, and weaker forms of mixing suffice. For an example of how this can be relaxed, see [13].

Bibliographical notes: The 0-1 law described in this section is due to Kalikow [38], who handled the i.i.d. setup. Our proof borrows from [82], which, still in the i.i.d. case, relaxes the uniform ellipticity assumption A2. In that paper, they show that a stronger 0-1 law holds if P is a product measure and $d = 2$, namely they show that $\mathbb{P}^\circ(A_\ell) \in \{0, 1\}$, while that last conclusion is false for certain mixing environments with elliptic, but not uniformly elliptic, environments.

3.2 A Law of Large Numbers in \mathbb{Z}^d

Our next goal is to prove a law of large numbers. Unfortunately, at this point we are not able to deal with general non i.i.d. environments (see however Remark 2 following the proof of Theorem 3.2.2), and further the case of i.i.d. environments does offer some simplifications. Thus, throughout this section we make the following assumptions:

Assumption 3.2.1 P is a uniformly elliptic, i.i.d. law on Ω .

The main result of this section is the following:

Theorem 3.2.2 Assume Assumption 3.2.1 and that $\mathbb{P}^\circ(A_\ell \cup A_{-\ell}) = 1$. Then, there exist deterministic $v_\ell, v_{-\ell}$ (possibly zero) such that

$$\lim_{n \rightarrow \infty} \frac{X_n \cdot \ell}{n} = v_\ell \mathbf{1}_{A_\ell} + v_{-\ell} \mathbf{1}_{A_{-\ell}}, \quad \mathbb{P}^\circ\text{-a.s.}$$

(See (3.2.8) for an expression for v_ℓ . When $v_\ell \neq 0$ for some ℓ , we say that the walk is *ballistic*).

Proof. As in Section 3.1 we will take here $\ell = (1, 0, \dots, 0)$. Further, we assume throughout that $\mathbb{P}^o(A_\ell) > 0$. The proof is based on introducing a renewal structure, as follows: Define $\bar{S}_0 = 0, M_0 = \ell \cdot X_0$,

$$\begin{aligned} \bar{S}_1 &= T_{M_0+1} \leq \infty, \quad \bar{R}_1 = D \circ \theta_{\bar{S}_1} + \bar{S}_1 \leq \infty, \\ M_1 &= \sup\{\ell \cdot X_m, \quad 0 \leq m \leq \bar{R}_1\} \leq \infty \end{aligned}$$

and by induction, for $k \geq 1$,

$$\begin{aligned} \bar{S}_{k+1} &= T_{M_k+1} \leq \infty, \quad \bar{R}_{k+1} = D \circ \theta_{\bar{S}_{k+1}} + \bar{S}_{k+1} \leq \infty, \\ M_{k+1} &= \sup\{\ell \cdot X_m, \quad 0 \leq m \leq \bar{R}_{k+1}\} \leq \infty. \end{aligned}$$

The times $\bar{S}_1, \bar{S}_2, \dots$, are called “fresh times”, and the locations $X_{\bar{S}_1}, X_{\bar{S}_2}, \dots$, are “fresh points”: at the time \bar{S}_k , the path $X \cdot$ visits for the first time after \bar{S}_{k-1} and after hitting again the hyperplane $X_{\bar{S}_{k-1}} \cdot \ell - 1$, a fresh part of the environment. Note that (\bar{S}_i, \bar{R}_i) are related to, but differ slightly from, (S_i, R_i) introduced in Section 3.1. Clearly,

$$0 = \bar{S}_0 \leq \bar{S}_1 \leq \bar{R}_1 \leq \bar{S}_2 \leq \dots \leq \infty$$

and the inequalities are strict if the left member is finite. Define:

$$\begin{aligned} K &= \inf\{k \geq 1 : \bar{S}_k < \infty, \bar{R}_k = \infty\} \leq \infty, \\ \tau_1 &= \bar{S}_K \leq \infty. \end{aligned}$$

τ_1 is called a “regeneration time”, because after τ_1 , $X \cdot \ell$ never falls behind $X_{\tau_1} \cdot \ell$.

By the same argument as in the proof of Theorem 3.1.2, $\mathbb{P}^o(\bar{R}_k < \infty) \leq \mathbb{P}^o(D < \infty)^k \xrightarrow[k \rightarrow \infty]{} 0$ because $\mathbb{P}^o(A_\ell) > 0$ implies $\mathbb{P}^o(D < \infty) < 1$. On the other hand, on A_ℓ , $\bar{R}_k < \infty \Rightarrow S_{k+1} < \infty$, \mathbb{P}^o -a.s., and hence

$$\mathbb{P}^o(A_\ell \cap \{K = \infty\}) = \mathbb{P}^o(A_\ell \cap \{\tau_1 = \infty\}) = 0.$$

Define now the measure

$$\mathbb{Q}^o(\cdot) = \mathbb{P}^o(\cdot \mid \{\tau_1 < \infty\}) = \mathbb{P}^o(\cdot \mid A_\ell)$$

and set

$$\mathcal{G}_1 = \sigma(\tau_1, X_0, \dots, X_{\tau_1}, \{\omega(y, \cdot)\}_{\ell \cdot y < \ell \cdot X_{\tau_1}}).$$

Note that since $\{D = \infty\} \subset \{\tau_1 < \infty\}$, we have that $\{D = \infty\} \in \mathcal{G}_1$. We have the following crucial lemma, whose proof is a simple exercise in the application of the Markov property, is omitted. It is here that the i.i.d. assumption on the environment plays a crucial role:

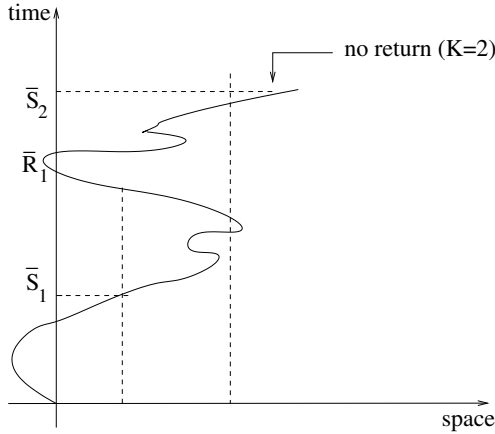


Fig. 3.2.1. Regeneration structure

Lemma 3.2.3 For any measurable sets A, B ,

$$\begin{aligned} & \mathbb{Q}^o \left(\{X_{\tau_1+n} - X_{\tau_1}\}_{n \geq 0} \in A, \{\omega(X_{\tau_1} + y, \cdot)\}_{y \cdot \ell \geq 0} \in B \right) \\ &= \mathbb{P}^o \left(\{X_n\}_{n \geq 0} \in A, \{\omega(y, \cdot)\}_{y \cdot \ell \geq 0} \in B \mid \{D = \infty\} \right). \end{aligned}$$

In fact,

$$\begin{aligned} & \mathbb{Q}^o \left(\{X_{\tau_1+n} - X_{\tau_1}\}_{n \geq 0} \in A, \{\omega(X_{\tau_1} + y, \cdot)\}_{y \cdot \ell \geq 0} \in B \mid \mathcal{G}_1 \right) \\ &= \mathbb{P}^o \left(\{X_n\}_{n \geq 0} \in A, \{\omega(y, \cdot)\}_{y \cdot \ell \geq 0} \in B \mid \{D = \infty\} \right). \end{aligned} \tag{3.2.4}$$

Proof of Lemma 3.2.3 Clearly, it suffices to prove (3.2.4). Let h denote a \mathcal{G}_1 measurable random variable. Set $1_A := 1_{\{X_n - X_0\}_{n \geq 0} \in A}$, $1_B := 1_{\{\omega(y, \cdot)\}_{y \cdot \ell \geq 0}}$. Further, note that for each $k \in \mathbb{N}$, $x \in \mathbb{Z}^d$, there exists a random variable $h_{x,k}$, measurable with respect to $\sigma(\{\omega(y, \cdot)\}_{\ell \cdot y < x \cdot \ell}, \{X_i\}_{i \leq \bar{S}_k}, \bar{S}_k)$, such that on the event $\{\tau_1 = \bar{S}_k, X_{\bar{S}_k} = x\}$, $h = h_{x,k}$ (this follows from the \mathcal{G}_1 measurability of h). Then, using θ to denote spatial shift and θ to denote temporal shift,

$$\begin{aligned}
 E_{\mathbb{P}^o} & \left(1_A \circ \theta^{\tau_1} \cdot 1_B \circ \bar{\theta}^{X_{\tau_1}} \cdot h \cdot \mathbf{1}_{\tau_1 < \infty} \right) \\
 & = \sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} E_P \left(E_{\omega}^o \left(\mathbf{1}_{\bar{S}_k < \infty} \mathbf{1}_{\bar{R}_k = \infty} \mathbf{1}_{X_{\bar{S}_k} = x} 1_A \circ \theta^{\bar{S}_k} \cdot 1_B \circ \bar{\theta}^x \cdot h_{x,k} \right) \right) \\
 & = \sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} E_P \left(1_B \circ \bar{\theta}^x E_{\omega}^o \left(\mathbf{1}_{\bar{S}_k < \infty} \mathbf{1}_{D \circ \bar{\theta}_{\bar{S}_k} = \infty} \mathbf{1}_{X_{\bar{S}_k} = x} 1_A \circ \theta^{\bar{S}_k} \cdot h_{x,k} \right) \right) \\
 & = \sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} E_P \left(1_B \circ \bar{\theta}^x E_{\omega}^x \left(\mathbf{1}_{D = \infty} 1_A \right) E_{\omega}^o \left(h_{x,k} \mathbf{1}_{\bar{S}_k < \infty} \mathbf{1}_{X_{\bar{S}_k} = x} \right) \right) \\
 & = \sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} E_P \left(1_B \mathbf{1}_{D = \infty} 1_A \right) E_P \left(h_{x,k} \mathbf{1}_{\bar{S}_k < \infty} \mathbf{1}_{X_{\bar{S}_k} = x} \right),
 \end{aligned}$$

where we used the Markov property in the next to last equality and the i.i.d. structure of the environment in the last one. Substituting in the above trivial A, B , one concludes that

$$E_{\mathbb{P}^o} \left(h \cdot \mathbf{1}_{\tau_1 < \infty} \right) = P(\{D = \infty\}) \sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} E_P \left(h_{x,k} \mathbf{1}_{\bar{S}_k < \infty} \mathbf{1}_{X_{\bar{S}_k} = x} \right).$$

Hence,

$$E_{\mathbb{Q}^o} \left(1_A \circ \theta^{\tau_1} \cdot 1_B \circ \bar{\theta}^{X_{\tau_1}} \cdot h \right) = E_{\mathbb{Q}^o}(h) E_{\mathbb{P}^o} \left(1_A 1_B | \{D = \infty\} \right),$$

concluding the proof of the lemma. □

Consider now τ_1 as a function of the path $(X_n)_{n \geq 0}$ and set

$$\tau_{k+1} = \tau_k(X_{\cdot}) + \tau_1(X_{\tau_k+ \cdot} - X_{\tau_k}),$$

with $\tau_{k+1} = \infty$ on $\{\tau_k = \infty\}$ (the sequence $\{\tau_k\}$ enumerates times such that for all $k < m < n$, $X_k \cdot \ell < X_m \cdot \ell \leq X_n \cdot \ell$). By the definition and the fact that $\mathbb{P}^o(A_{\ell} \cap \{\tau_1 = \infty\}) = 0$, we have that $\mathbb{P}^o(A_{\ell} \cap \{\tau_k = \infty\}) = 0$. Setting

$$\mathcal{G}_k = \sigma \left(\tau_1, \dots, \tau_k, \quad X_0, \dots, X_{\tau_k}, \quad \{\omega(y, \cdot)\}_{\ell \cdot y < \ell \cdot X_{\tau_k}} \right),$$

an obvious rerun of the proof of Lemma 3.2.3 yields that

$$\begin{aligned}
 & \mathbb{Q}^o \left(\{X_{\tau_k+n} - X_{\tau_k}\}_{n \geq 0} \in A, \{\omega(X_{\tau_k+y}, \cdot)\}_{\ell \cdot y \geq 0} \in B | \mathcal{G}_k \right) \\
 & = \mathbb{P}^o \left(\{X_n\}_{n \geq 0} \in A, \{\omega(y, \cdot)\}_{\ell \cdot y \geq 0} \in B | \{D = \infty\} \right).
 \end{aligned}$$

We thus conclude that under \mathbb{Q}^o ,

$$(X_{\tau_2} - X_{\tau_1}, \tau_2 - \tau_1), \dots, (X_{\tau_{k+1}} - X_{\tau_k}, \tau_{k+1} - \tau_k)$$

are i.i.d. pairs of random variables, independent of (X_{τ_1}, τ_1) , such that

$$\mathbb{Q}^o(X_{\tau_2} - X_{\tau_1} \in C_1, \tau_2 - \tau_1 \in C_2) = \mathbb{P}^o(X_{\tau_1} \in C_1, \tau_1 \in C_2 | \{D = \infty\}).$$

Next, we have the following lemma, whose proof is deferred:

Lemma 3.2.5 (Zerner)

$$\mathbb{E}^o(X_{\tau_1} \cdot \ell | \{D = \infty\}) = \frac{1}{\mathbb{P}^o(D = \infty)}.$$

We are now ready to complete the proof of Theorem 3.2.2. Assume first that $\mathbb{E}^o(\tau_1 | \{D = \infty\}) < \infty$. Then, by the law of large numbers, and the renewal structure,

$$\frac{\tau_k}{k} \xrightarrow{k \rightarrow \infty} \mathbb{E}^o(\tau_1 | \{D = \infty\}), \quad \mathbb{Q}^o\text{-a.s.} \tag{3.2.6}$$

$$\frac{X_{\tau_k} \cdot \ell}{k} \xrightarrow{k \rightarrow \infty} \mathbb{E}^o(X_{\tau_1} \cdot \ell | \{D = \infty\}), \quad \mathbb{Q}^o\text{-a.s.} \tag{3.2.7}$$

(note that the finiteness of the expression in the right hand side of (3.2.7) is trivial if the right hand side of (3.2.6) is finite, and Lemma 3.2.5 is not needed in this case).

Hence,

$$\frac{X_{\tau_k} \cdot \ell}{k} \xrightarrow{k \rightarrow \infty} \frac{\mathbb{E}^o(X_{\tau_1} \cdot \ell | \{D = \infty\})}{\mathbb{E}^o(\tau_1 | \{D = \infty\})} =: v_\ell, \quad \mathbb{Q}\text{-a.s.} \tag{3.2.8}$$

Mimicking now the argument at the end of the proof of Lemma 2.1.5, we conclude that $\frac{X_n \cdot \ell}{n} \xrightarrow{n \rightarrow \infty} v_\ell, \mathbb{Q}^o\text{-a.s.}$, in the case $\mathbb{E}^o(\tau_1 | \{D = \infty\}) < \infty$.

On the other hand, Lemma 3.2.5 implies that (3.2.7) holds true even when $\mathbb{E}^o(\tau_1 | \{D = \infty\}) = \infty$. But then, $\tau_k/k \xrightarrow{k \rightarrow \infty} \infty, \mathbb{Q}^o\text{-a.s.}$ With $v_\ell = 0$ in this case, we conclude that

$$\frac{X_{\tau_k} \cdot \ell}{\tau_k} \xrightarrow{k \rightarrow \infty} v_\ell = 0, \quad \mathbb{Q}^o\text{-a.s.}$$

Finally, setting k_n such that $\tau_{k_n} \leq n < \tau_{k_n+1}$, we have that $k_n \xrightarrow{n \rightarrow \infty} \infty$ and $k_n/n \xrightarrow[n \rightarrow \infty]{} 0, \mathbb{Q}^o\text{-a.s.}$ because $n/k_n \geq \tau_{k_n}/k_n$. Thus,

$$\frac{X_n \cdot \ell}{n} \leq \frac{X_{\tau_{k_n+1}} \cdot \ell}{k_n + 1} \cdot \frac{k_n + 1}{n} \xrightarrow[n \rightarrow \infty]{} 0, \quad \mathbb{Q}^o\text{-a.s.}$$

Since $\liminf_{n \rightarrow \infty} \frac{X_n \cdot \ell}{n} \geq 0, \mathbb{Q}^o\text{-a.s.}$, we conclude

$$\frac{X_n \cdot \ell}{n} \xrightarrow[n \rightarrow \infty]{} 0, \quad \mathbb{Q}^o\text{-a.s.} \quad \square$$

Remarks:

1. Note that on $A_\ell, v_\ell > 0$ if $\mathbb{E}^o(\tau_1 | \{D = \infty\}) < \infty$ and $v_\ell = 0$ otherwise.
2. It is clear from the proof that in fact, if $\mathbb{E}^o(\tau_1 | \{D = \infty\}) < \infty$, then the result of Theorem 3.2.2 can be strengthened to

$$\frac{X_n}{n} \xrightarrow{n \rightarrow \infty} \frac{\mathbb{E}^o(X_{\tau_1} | \{D = \infty\})}{\mathbb{E}^o(\tau_1 | \{D = \infty\})}, \quad \mathbb{Q}^o\text{-a.s.}$$

3. In the stationary, ϕ -mixing case, one can prove that the times $\{\tau_i\}$ are well defined, and form a mixing sequence. What I have not been able to show is that they are identically distributed under \mathbb{Q}^o (they seem not!). Modifying slightly the definition of $(\overline{R}_k, \overline{S}_k)$ by adding an L -safeguard as in Section 3.1, the results of this section extend immediately to the case where the environment is K -dependent (i.e., $\{\omega(x, \cdot)\}_{x \cdot \ell \leq 0}$ and $\{\omega(x, \cdot)\}_{x \cdot \ell > K}$ are independent). This applies, e.g., in the setup considered in [63]. The extension to a mixing setup is more complicated, and some results applicable there can be found in [13].
4. Still discussing mixing environments, some progress has been made using the approach of the environment viewed from the particle. We mention here [44] and in particular the recent preprint [62]. The latter preprint uses a-priori estimates concerning regeneration times in the ballistic case to construct an invariant measure for the environment viewed from the particle which is absolutely continuous with respect to P on certain half-spaces, and deduces a LLN using that measure.

Proof of Lemma 3.2.5

Recall that we consider $\ell = (1, 0, \dots, 0)$. Then,

$$\begin{aligned}
 & \mathbb{Q}^o\left(\{\exists k : X_{\tau_k} \cdot \ell = i\}\right) \\
 &= \frac{\sum_{y \in \mathbb{Z}^{d-1}} E\left(P_\omega^o(\{\exists k : X_{\tau_k} = (i, y)\}, A_\ell)\right)}{\mathbb{P}^o(A_\ell)} \\
 &= \frac{\sum_{y \in \mathbb{Z}^{d-1}} E\left(P_\omega^o(T_i < \infty, X_{T_i} = (i, y), D \circ \theta_{T_i} = \infty)\right)}{\mathbb{P}^o(A_\ell)} \\
 &= \frac{\sum_{y \in \mathbb{Z}^{d-1}} E\left(P_\omega^o(T_i < \infty, X_{T_i} = (i, y))P_\omega^{(i,y)}(D = \infty)\right)}{\mathbb{P}^o(A_\ell)} \\
 &= \frac{\mathbb{P}^o(T_i < \infty)\mathbb{P}^o(D = \infty)}{\mathbb{P}^o(A_\ell)} \xrightarrow{i \rightarrow \infty} \mathbb{P}^o(D = \infty) \tag{3.2.9}
 \end{aligned}$$

(since $\mathbb{P}^o(A_\ell \cup A_{-\ell}) = 1$ and $\lim_{i \rightarrow \infty} \mathbb{P}^o(\{T_i < \infty\} \cap A_{-\ell}) = 0$). On the other hand,

$$\lim_{i \rightarrow \infty} \mathbb{Q}^o\left(\{\exists k : X_{\tau_k} \cdot \ell = i\}\right) = \lim_{i \rightarrow \infty} \mathbb{Q}^o\left(\{\exists k \geq 2 : X_{\tau_k} \cdot \ell = i\}\right)$$

(because $\mathbb{Q}^o(\tau_k > i) \xrightarrow{i \rightarrow \infty} 0$)

$$\begin{aligned}
 &= \lim_{i \rightarrow \infty} \sum_{n \geq 1} \mathbb{Q}^o \left(\{ \exists k \geq 2 : X_{\tau_k} \cdot \ell = i, X_{\tau_1} \cdot \ell = n \} \right) \\
 &= \lim_{i \rightarrow \infty} \sum_{n \geq 1} \mathbb{Q}^o \left(\{ \exists k \geq 2 : (X_{\tau_k} - X_{\tau_1}) \cdot \ell = i - n, X_{\tau_1} \cdot \ell = n \} \right) \\
 &= \lim_{i \rightarrow \infty} \sum_{n \geq 1} \mathbb{Q}^o (X_{\tau_1} \cdot \ell = n) \cdot \mathbb{Q}^o \left(\{ \exists k \geq 2 : (X_{\tau_k} - X_{\tau_1}) \cdot \ell = i - n \} \right).
 \end{aligned}$$

But, recall that by the renewal theorem,

$$\mathbb{Q}^o \left(\exists k \geq 2 : (X_{\tau_k} - X_{\tau_1}) \cdot \ell = i - n \right) \xrightarrow{i \rightarrow \infty} \frac{1}{E_{\mathbb{Q}^o}((X_{\tau_2} - X_{\tau_1}) \cdot \ell)}$$

and hence, by dominated convergence,

$$\lim_{i \rightarrow \infty} \mathbb{Q}^o \left(\{ \exists k : X_{\tau_k} \cdot \ell = i \} \right) = \frac{\sum_{n \geq 1} \mathbb{Q}^o (X_{\tau_1} \cdot \ell = n)}{E_{\mathbb{Q}^o}((X_{\tau_2} - X_{\tau_1}) \cdot \ell)} = \frac{1}{E_{\mathbb{Q}^o}((X_{\tau_2} - X_{\tau_1}) \cdot \ell)}. \tag{3.2.10}$$

Comparing (3.2.9) and (3.2.10), we conclude that

$$E_{\mathbb{Q}^o} \left((X_{\tau_2} - X_{\tau_1}) \cdot \ell \right) = \frac{1}{\mathbb{P}^o(D = \infty)} < \infty. \quad \square$$

Theorem 3.2.2 assumes that $\mathbb{P}^o(A_\ell \cup A_{-\ell}) = 1$, and in that situation provided a LLN if $\mathbb{P}^o(A_\ell) \in \{0, 1\}$. A recent improvement to Theorem 3.2.2, due to Zerner [83], actually shows that if a 0-1 law holds true, a LLN holds, at least for i.i.d. environments. More precisely, one has the following:

Theorem 3.2.11 *There exist deterministic $v_\ell, v_{-\ell}$ (possibly zero) such that*

$$\lim_{n \rightarrow \infty} \frac{X_n \cdot \ell}{n} = v_\ell \mathbf{1}_{A_\ell} + v_{-\ell} \mathbf{1}_{A_{-\ell}}, \quad \mathbb{P}^o\text{-a.s.} \tag{3.2.12}$$

An immediate corollary, obtained by applying Theorem 3.2.11 d times with respect to the basis $\ell = e_i, i = 1, \dots, d$, is the following:

Corollary 3.2.13 *Assume that $\mathbb{P}^o(A_\ell) \in \{0, 1\}$ for every ℓ . Then there exists a deterministic v (possibly zero) such that*

$$\lim_{n \rightarrow \infty} \frac{X_n}{n} = v, \quad \mathbb{P}^o\text{-a.s.}$$

Proof of Theorem 3.2.11: (sketch) In view of the 0-1 law Theorem 3.1.2 and of Theorem 3.2.2, all that remains to prove is that if $\mathbb{P}^o(A_\ell \cup A_{-\ell}) = 0$ then $X_n \cdot \ell/n \rightarrow 0$, \mathbb{P}^o -a.s. The complete proof for that is given in [83], and we provide next a brief description.

Consider the set of visits to the hyperplane $\mathcal{H}_m := \{z : z \cdot \ell = m\}$, defining $\tau_m^0 = T_m$ and $\tau_m^i = \min\{n > \tau_m^{i-1} : X_n \cdot \ell = m\}$. Fixing an integer L , let

$$h_{m,L} = \sup_{i \geq 0} \{ \tau_m^i - \tau_m^0 : \tau_m^i < T_{m+L} \}$$

be the diameter of the set of visits to \mathcal{H}_m before T_{m+L} . For any constant $c > 0$, let

$$F_{M,L}(c) = \frac{\#\{0 \leq m \leq M : h_{m,L} \leq c\}}{M + 1}$$

denote the fraction of m 's smaller than M such that the time between the first and last visit to \mathcal{H}_m before T_{m+L} is smaller than c . The first observation, which is a deterministic (combinatorial) computation that we skip, is that for any path with $\liminf_{n \rightarrow \infty} X_n \cdot \ell/n > 0$ there exists a constant c such that

$$\inf_{L \geq 1} \limsup_{M \rightarrow \infty} F_{M,L}(c) > 0,$$

that is, roughly, there is a fraction of m 's for which the time between first and last visits of \mathcal{H}_m (before hitting \mathcal{H}_{m+L}) is not too large.

Assume now that $\mathbb{P}^o(\limsup X_n \cdot \ell/n > 0) > 0$. Then, by the above observation, there is some $c > 0$ such that

$$\mathbb{P}^o(\limsup_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} F_{M,L}(c) > 0) > 0. \tag{3.2.14}$$

But on the event $\{h_{m,L} \leq c\}$, the last point visited in \mathcal{H}_m before hitting \mathcal{H}_{m+L} is at most at distance c from X_{T_m} and has been visited at most c times before T_{m+L} . Thus, there is a $z \in \mathcal{H}_0$ with $|z|_1 \leq c$, and an $1 \leq r \leq c$ such that the r -th visit to $X_{T_m} + z$ occurs before T_{m+L} and the walk does not backtrack from \mathcal{H}_m after this r -th visit. Denoting the last event by $B_{m,L}^1(z, r)$, it follows that

$$F_{M,L}(c) \leq \frac{1}{M + 1} \sum_{z \in \mathcal{H}_0, |z|_1 \leq c} \sum_{r=1}^c \sum_{m=0}^M \mathbf{1}_{B_{m,L}^1(z,r)}.$$

Noting that the summation over r and z is over a finite set, and combining the last inequality with (3.2.14), it follows that for some z and r ,

$$\mathbb{P}^o(\limsup_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} \frac{1}{M + 1} \sum_{m=0}^M \mathbf{1}_{B_{m,L}^1(z,r)} > 0) > 0. \tag{3.2.15}$$

While the events $\{B_{m,L}^1(z, r)\}_m$ are not independent, some independence can be restored in the following way: construct independent (given the environment) copies Y^y of the RWRE, starting at y . Define the event $B_{m,L}(z, r)$ as the union of $B_{m,L}^1(z, r)$ with the event that X does not hit $X_{T_m} + z$ for the r -th time before T_{m+L} , but $Y^{X_{T_m} + z}$ does not backtrack from \mathcal{H}_m before it hits \mathcal{H}_{m+L} . An easy computation involving the Markov property shows that for each fixed $i = 0, 1, \dots, L - 1$, the events $\{B_{jL+i,L}(z, r)\}_j$ are independent, with

$$\mathbb{P}^\circ(B_{jL+i,L}(z,r)) = \mathbb{P}^\circ(D \geq T_L).$$

(Here and in the sequel, we abuse notations by still using \mathbb{P}° to denote the annealed law on the enlarged probability space that supports the extra Y^y walks). Hence, since we have from (3.2.15) that

$$\mathbb{P}^\circ(\limsup_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} \frac{1}{M+1} \sum_{i=0}^{L-1} \sum_{j=0}^{[M/L]} \mathbf{1}_{B_{jL+i,L}^1(z,r)} > 0) > 0, \tag{3.2.16}$$

it follows, by the standard law of large numbers, that

$$\mathbb{P}^\circ(D = \infty) = \limsup_{L \rightarrow \infty} \mathbb{P}^\circ(D \geq T_L) > 0.$$

But from (3.1.4), we have that $\mathbb{P}^\circ(A_\ell) \geq \mathbb{P}^\circ(D = \infty) > 0$. In particular, this shows that $\mathbb{P}^\circ(A_\ell) = 0$ implies that $\limsup X_n \cdot \ell/n \leq 0$, \mathbb{P}° -a.s. Repeating this argument with $-\ell$ instead of ℓ completes the proof of the theorem. \square

Bibliographical notes:

The proof here follows closely [76], except that Lemma 3.2.5 is due to private communication with Martin Zerner. The improvement Theorem 3.2.11 is based on [83].

The ballistic LLN has been proved for certain non iid environments in [13]. Alternative approaches to ballistic LLN's using the environment viewed from the particle were developed in [44] and in great generality in [62].

There are only a few LLN results in the non-ballistic case, see the bibliographical notes of Section 3.3.

3.3 CLT for walks in balanced environments

The setup in this section is the following:

Assumption 3.3.1

- (B1) *P is stationary and ergodic.*
- (B2) *P is balanced: for $i = 1, \dots, d$, $P(\omega(x, x + e_i) = \omega(x, x - e_i)) = 1$.*
- (B3) *P is uniformly elliptic: there exists an $\varepsilon > 0$ such that for $i = 1, \dots, d$,*

$$P(\omega(x, x + e_i) > \varepsilon) = 1.$$

Unlike the situation in Section 2.1, we do not have an explicit construction of invariant measures at our disposal. The approach toward the LLN and CLT uses however (B2) in an essential way: indeed, note that in the notations of (2.1.28),

$$d(x, \omega) = \sum_{i=1}^d e_i \left[\omega(x, x + e_i) - \omega(x, x - e_i) \right] = 0.$$

Hence, the processes $(X_n(i))_{n \geq 0}, i = 1, \dots, d$, are martingales, with, denoting $\mathcal{F}_n = \sigma(X_0, \dots, X_n)$,

$$E_\omega^o((X_n(i) - X_{n-1}(i))(X_n(j) - X_{n-1}(j)) | \mathcal{F}_{n-1}) = 2\delta_{ij}\omega(X_{n-1}, X_{n-1} + e_i).$$

Since $|\omega(\cdot, \cdot)| \leq 1$ P -a.s., it immediately follows that $X_n/n \xrightarrow[n \rightarrow \infty]{} 0$, \mathbb{P}^o -a.s. Further, the multi-dimensional CLT (compare with Lemma 2.2.4) yields that if there exists a deterministic vector $\mathbf{a} = (a_1, \dots, a_d)$ such that

$$\frac{1}{n} \sum_{k=1}^n \omega(X_{k-1}, X_{k-1} + e_i) \xrightarrow[n \rightarrow \infty]{} \frac{a_i}{2} > 0, \quad \mathbb{P}^o\text{-a.s.}, \tag{3.3.2}$$

then, for any bounded continuous function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and any $y \in \mathbb{R}$,

$$\begin{aligned} \lim_{n \rightarrow \infty} P_\omega^o \left(f \left(\frac{X_n}{\sqrt{n}} \right) \leq y \right) & \tag{3.3.3} \\ &= \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sqrt{a_i}} \int_{\mathbb{R}^d} \mathbf{1}_{\{f(\mathbf{x}) \leq y\}} \exp \left(- \sum_{i=1}^d \frac{x_i^2}{2a_i} \right) \prod_{i=1}^d dx_i, \quad P\text{-a.s.} \end{aligned}$$

Our goal in this section is to demonstrate such a CLT, and to study transience and recurrent questions for the RWRE.

Central limit theorems

Theorem 3.3.4 *Assume Assumption 3.3.1. Then, there exists a deterministic vector \mathbf{a} such that (3.3.2) holds true. Consequently, the quenched CLT (3.3.3) holds true.*

Remark 3.3.5 *In fact, the above observations yield not only a CLT in the form of (3.3.3) but also a trajectorial CLT for the process $\{X_{[nt]}/\sqrt{n}, t \in [0, 1]\}$.*

Proof of Theorem 3.3.4

As in Section 2.1, the key to the proof of (3.3.2) is to consider the environment viewed from the particle. Define $\bar{\omega}(n) = \theta^{X_n}\omega$, and the Markov transition kernel

$$M(\omega, d\omega') = \sum_{e_i} \left[\omega(0, e_i) \delta_{\theta^{e_i}\omega = \omega'} + \omega(0, -e_i) \delta_{\theta^{-e_i}\omega = \omega'} \right]. \tag{3.3.6}$$

As in Lemma 2.1.18, the process $\bar{\omega}(n)$ is Markov under either P_ω^o or \mathbb{P}^o . Mimicking the proof of Corollary 2.1.25, if we can construct a measure Q on Ω which is absolutely continuous with respect to P and such that it is invariant under the Markov transition M , we will conclude, as in Corollary 2.1.25, that $\bar{\omega}(n)$ is stationary and ergodic and hence

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \omega(X_{n-1}, X_{n-1} + e_i) &= \frac{1}{n} \sum_{i=1}^n \bar{\omega}(n)(0, e_i) \xrightarrow[n \rightarrow \infty]{} \frac{a_i}{2} \\ &:= E_Q \bar{\omega}(0, e_i) \geq \varepsilon, \mathbb{P}^o\text{-a.s.}, \end{aligned} \tag{3.3.7}$$

yielding (3.3.2). Our effort therefore is directed towards the construction of such a measure. Naturally, such measures will be constructed from periodic modifications of the RWRE, and require certain a-priori estimates on harmonic functions. We state these now, and defer their proof to the end of the section. The estimates we state are slightly more general than needed, but will be useful also in the study of transience and recurrence.

We let $|x|_\infty := \max_{i=1}^d |x_i|$ and define $D = D_R(x_0) = \{x \in \mathbb{Z}^d : |x-x_0|_\infty < R\}$. The generator of the RWRE, under P_ω , is the operator

$$(L_\omega f)(x) = \sum_{i=1}^d \omega(x, x + e_i) \left[f(x + e_i) + f(x - e_i) - 2f(x) \right].$$

For any bounded $E \subset \mathbb{Z}^d$ of cardinality $|E|$, set $\partial E = \{y \in E^c : \exists x \in$

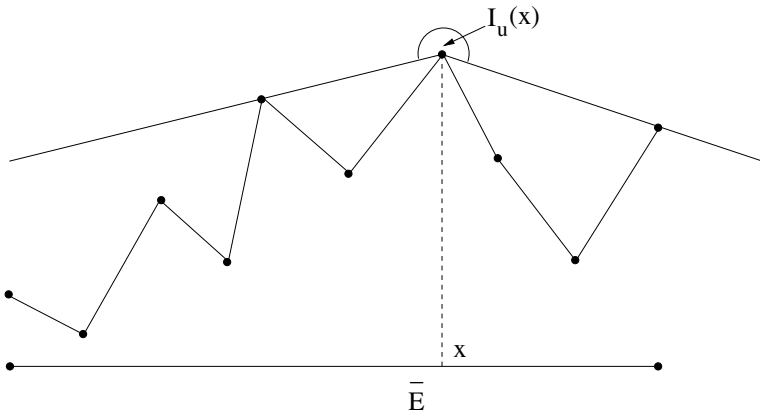


Fig. 3.3.1. The normal set at $x \in E$

$E, |x - y|_\infty = 1\}$, $\overline{E} = E \cup \partial E$, and $\text{diam}(E) = \max\{|x - y|_\infty : x, y \in \overline{E}\}$. For any function $u : \mathbb{Z}^d \rightarrow \mathbb{R}$, we define the *normal set* at a point $x \in E$ as

$$I_u(x) = \{s \in \mathbb{R}^d : u(z) \leq u(x) + s \cdot (z - x), \forall z \in \overline{E}\}.$$

Finally, for any $q > 0$, E and u as above, define

$$\|g\|_{E,q,u} := \left(\frac{1}{|E|} \sum_{x \in E} 1_{\{I_u(x) \neq \emptyset\}} |g(x)|^q \right)^{1/q}, \quad \|g\|_{E,q} := \left(\frac{1}{|E|} \sum_{x \in E} |g(x)|^q \right)^{1/q}.$$

Then we have the following:

Lemma 3.3.8 *There exists a constant $C = C(\varepsilon, d)$ such that*

(a) (*maximum principle*) For any $E \subset \mathbb{Z}^d$ bounded, any functions u and g such that

$$L_\omega u(x) \geq -g(x), \quad x \in E$$

satisfy

$$\max_{x \in E} u(x) \leq C \text{diam}(E) |E|^{1/d} \|g\|_{E,d,u} + \max_{x \in \partial E} u^+(x).$$

(b) (*Harnack inequality*) Any function $u \geq 0$ such that

$$L_\omega u(x) = 0, \quad x \in D_R(x_0), \tag{3.3.9}$$

satisfies

$$\frac{1}{C} u(x_0) \leq u(x) \leq C u(x_0), \quad x \in D_{R/2}(x_0).$$

We now introduce a periodic structure. Set $\Delta_N = \{-N, \dots, N\}^d \subset \mathbb{Z}^d$ and identify elements of $T_N = \mathbb{Z}^d / (2N + 1)\mathbb{Z}^d$ with a point of Δ_N , setting $\pi_N : \mathbb{Z}^d \rightarrow T_N$ and $\hat{\pi}_N : \mathbb{Z}^d \rightarrow \Delta_N$ to be the canonical projections. Set $\Omega^N = \{\omega \in \Omega : \theta^x \omega = \omega, \forall x \in (2N + 1)\mathbb{Z}^d\}$. For any $\omega \in \Omega$, define $\omega^N \in \Omega^N$ by $\omega^N(x) = \omega(\hat{\pi}_N x)$. Note that ω^N is then a well defined function on T_N too.

Due to the ergodicity of P , it holds that in the sense of weak convergence,

$$P_N := \frac{1}{(2N + 1)^d} \sum_{x \in \Delta_N} \delta_{\theta^x \omega^N} \xrightarrow[N \rightarrow \infty]{} P, \quad P\text{-a.s.} \tag{3.3.10}$$

Let $\Omega_0 \subset \Omega$ denote those environments ω for which the convergence holds in (3.3.10) (clearly, $P(\Omega_0) = 1$).

Fixing $\omega \in \Omega_0$, let $(X_{n,N})_{n \geq 0}$ denote the RWRE on \mathbb{Z}^d with law $P_{\omega^N}^{X_{0,N}}$. Then, $\bar{X}_{n,N} := \pi_N X_{0,N}$ is an irreducible Markov chain with finite state space T_N , and hence it possesses a unique invariant measure $\mu_N = \frac{1}{(2N+1)^d} \sum_{x \in T_N} \phi_N(x) \delta_x$. Setting $\bar{\omega}^N(n) := \theta^{X_{n,N}} \omega^N$, it follows that $\bar{\omega}^N(n)$ is an irreducible Markov chain with finite state space $S_N := \{\theta^x \omega^N\}_{x \in \Delta_N}$ and transition kernel M . Its unique invariant measure, supported on Ω^N , is then easily checked to be of the form

$$Q_N = \frac{1}{(2N + 1)^d} \sum_{x \in \Delta_N} \phi_N(\pi_N x) \delta_{\theta^x \omega^N}.$$

Partitioning the state space S_N into finitely many *disjoint* states $\{\omega_\alpha^N\}_{\alpha=1}^K$, set $C_N(\alpha) = \{x \in \Delta_N : \theta^x \omega^N = \omega_\alpha^N\}$. Then,

$$f_N := \frac{dQ_N}{dP_N} = \sum_{\alpha=1}^K \mathbf{1}_{\{\omega = \omega_\alpha^N\}} \frac{1}{|C_N(\alpha)|} \sum_{x \in C_N(\alpha)} \phi_N(\pi_N x).$$

We show below, as a consequence of part (a) of Lemma 3.3.8, that there exists a constant $C_2 = C_2(\varepsilon, d)$, independent of N , such that

$$\|\phi_N(\pi_N \cdot)\|_{D_{N+1}(0),d/d-1} \leq C_2. \tag{3.3.11}$$

Thus, using Jensen’s inequality in the first inequality and (3.3.11) in the second,

$$\begin{aligned} \int f_N^{d/d-1} dP_N &= \sum_{\alpha=1}^K \left[\frac{1}{|C_N(\alpha)|} \sum_{x \in C_N(\alpha)} \phi_N(\pi_N x) \right]^{d/d-1} \frac{|C_N(\alpha)|}{(2N+1)^d} \\ &\leq \sum_{\alpha=1}^K \sum_{x \in C_N(\alpha)} \phi_N(\pi_N(x))^{d/d-1} \frac{1}{(2N+1)^d} \\ &= \frac{1}{(2N+1)^d} \sum_{x \in \Delta_N} \phi_N(\pi_N(x))^{d/d-1} \leq C_2^{(d-1)/d}. \end{aligned} \tag{3.3.12}$$

Note that f_N extends to a measurable function on Ω , and the latter is, due to (3.3.12), uniformly integrable with respect to P_N . Thus, any weak limit of Q_N is absolutely continuous with respect to P , and further it is invariant with respect to the Markov kernel M .

Let $E = \{\omega : \frac{dQ}{dP} = 0\}$. By invariance, $E_Q M \mathbf{1}_E = E_Q \mathbf{1}_E = 0$, and hence $M \mathbf{1}_E \leq \mathbf{1}_E$, P -a.s. But, $M \mathbf{1}_E \geq \varepsilon \sum_{i=1}^d (\mathbf{1}_E \circ \theta^{e_i} + \mathbf{1}_E \circ \theta^{-e_i})$. Hence, $\mathbf{1}_E \geq \mathbf{1}_E \circ \theta^{\pm e_i}$, P -a.s. Since P is stationary, $\mathbf{1}_E = \mathbf{1}_E \circ \theta^{\pm e_i}$, P -a.s., and hence by ergodicity (considering the invariant event $\cap_{x \in \mathbb{Z}^d} (\theta^x)^{-1} E$) $P(E) \in \{0, 1\}$. But $Q \ll P$ implies $P(E) = 0$. Hence, $Q \sim P$, as claimed (further, by (3.3.7), Q is then uniquely defined).

It thus only remains to prove (3.3.11). Fix a function g on T_N , and define the resolvent

$$\begin{aligned} R^{\omega_N} g(x) &:= \sum_{j=0}^{\infty} \left(1 - \frac{1}{N^2}\right)^j E_{\omega_N}^x g(\bar{X}_{j,N}) \\ &= \sum_{j=0}^{\infty} \left(1 - \frac{1}{N^2}\right)^j E_{\omega_N}^x g \circ \pi_N(X_{j,N}), \quad x \in T_N \end{aligned}$$

and the stopping times $\tau_0 = 0, \tau_1 = \tau := \min\{k \geq 1 : |X_{k,N} - X_{0,N}| \geq N\}$ and $\tau_{k+1} = \tau \circ \theta^k + \tau_k$. Since for $x \in \mathbb{Z}^d$ with $|x - X_{0,N}| < N$ it holds that $L_{\omega_N} E_{\omega_N}^x \left(\sum_{j=0}^{\tau-1} g \circ \pi_N(X_{j,N})\right) = -g(x)$, we have by Lemma 3.3.8(a) that for some constant $C = C(\varepsilon, d)$,

$$\sup_{|x - X_{0,N}| < N} \left| E_{\omega_N}^x \left(\sum_{j=0}^{\tau-1} g \circ \pi_N(X_{j,N}) \right) \right| \leq CN^2 \|g\|_{D_{N+1}(0),d}. \tag{3.3.13}$$

Since $(X_{n,N})_{n \geq 0}$ is a martingale, it follows from Doob’s inequality that, for any $K \geq 1$,

$$\begin{aligned}
 P_{\theta^x \omega^N}^o[\tau \leq K] &\leq 2 \sum_{i=1}^d P_{\theta^x \omega^N}^o \left[\sup_{n \leq K} X_n(i) \geq N \right] \\
 &\leq \frac{2}{N} \sum_{i=1}^d E_{\theta^x \omega^N}^o \left((X_K(i))_+ \right) \leq \frac{2d}{N} \sqrt{K}.
 \end{aligned}$$

Hence, using $K = N^2/8d^2$,

$$E_{\theta^x \omega^N}^o \left(\left(1 - \frac{1}{N^2} \right)^\tau \right) \leq \frac{2d}{N} \sqrt{K} + \left(1 - \frac{1}{N^2} \right)^K \leq C_3 \tag{3.3.14}$$

where $C_3 = C_3(d) < 1$ is independent of N . Thus, using the strong Markov property, (3.3.13) and (3.3.14),

$$\begin{aligned}
 |R^{\omega^N} g(x)| &= \sum_{m \geq 0} E_{\omega^N}^x \left(\sum_{\tau_m \leq j < \tau_{m+1}} \left(1 - \frac{1}{N^2} \right)^j g \circ \pi_N(X_{j,N}) \right) \\
 &\leq \sum_{m \geq 0} E_{\omega^N}^x \left(\left(1 - \frac{1}{N^2} \right)^{\tau_m} E_{\omega^N}^{X_{\tau_m,N}} \sum_{j=0}^{\tau-1} g \circ \pi_N(X_{j,N}) \right) \\
 &\leq \sum_{m \geq 0} \left(\sup_{x \in \mathbb{Z}^d} E_{\omega^N}^x \left(\left(1 - \frac{1}{N^2} \right)^\tau \right) \right)^m \cdot \sup_{x \in \mathbb{Z}^d} E_{\omega^N}^x \left(\sum_{j=0}^{\tau-1} g \circ \pi_N(X_{j,N}) \right) \\
 &\leq C_4 N^2 \|g\|_{D_{N+1}(0),d}
 \end{aligned}$$

where $C_4 = C_4(d, \varepsilon)$. Using the invariance of ϕ_N , we now get

$$\begin{aligned}
 \|\phi_N(\pi_N \cdot)\|_{D_{N+1}(x_0),d/d-1} &= \|\phi_N(\pi_N \cdot)\|_{D_{N+1}(0),d/d-1} \\
 &= \sup_{g: \|g\|_{D_{N+1}(0),d} \leq 1} \frac{1}{|D_{N+1}(0)|} \sum_{y \in D_{N+1}(0)} \phi_N(\pi_N y) g(y) \\
 &= \frac{1}{N^2} \sup_{g: \|g\|_{D_{N+1}(0),d} \leq 1} \sum_{k \geq 0} \left(1 - \frac{1}{N^2} \right)^k \frac{1}{(2N+1)^d} \sum_{x \in \Delta_N} \phi_N(x) E_{\omega^N}^x (g \circ \pi_N(X_{k,N})) \\
 &\leq C_2
 \end{aligned}$$

with $C_2 = C_2(d, \varepsilon)$, proving (3.3.11). □

Proof of Lemma 3.3.8

(a) We may assume without loss of generality that $\max_{x \in \partial E} u(x) \leq 0$, $g \geq 0$, $g \neq 0$ and that $u \geq 0$ is not identically 0. Let $\bar{u} = \max_{x \in \bar{E}} u = u(x_0)$, some $x_0 \in E$. Then, for s satisfying $|s|_\infty < \bar{u}/\text{diam}(\bar{E})$, it holds that

$$u(x_0) + s \cdot (x - x_0) > 0, \quad \forall x \in \bar{E}.$$

Hence, with $t = \inf\{\rho \geq 0 : u(x_0) + s(x - x_0) + \rho > u(x), \forall x \in \overline{E}\}$, we have that $u(x) = u(x_0) + s \cdot (x - x_0) + t$, some $x \in \overline{E}$, and hence $u(x) + s \cdot (z - x) = u(x_0) + s \cdot (z - x_0) + t \geq u(z)$, $\forall z \in \overline{E}$. Hence,

$$s \in I_u(x) \subset \bigcup_{x \in E} I_u(x), \text{ for all } s \text{ with } |s|_\infty < \frac{\overline{u}}{\text{diam}(\overline{E})}. \tag{3.3.15}$$

Assume $s \in I_u(x)$. Then, with $e \in \{\pm e_i\}$, and $v(y) = u(x) + s \cdot (y - x)$,

$$\begin{aligned} 0 &= \omega(x, x + e) \left(2v(x) - v(x + e) - v(x - e) \right) \\ &\leq \omega(x, x + e) (2u(x) - u(x + e) - u(x - e)), \end{aligned}$$

and hence,

$$\begin{aligned} 0 &\leq \omega(x, x + e) (2u(x) - u(x + e) - u(x - e)) \\ &\leq \sum_{i=1}^d \omega(x, x + e_i) (2u(x) - u(x + e_i) - u(x - e_i)) = -L_\omega u(x) \leq g(x). \end{aligned}$$

Hence,

$$\left(u(x) - u(x - e) \right) - \left(u(x + e) - u(x) \right) \leq \frac{g(x)}{\omega(x, x + e)} \leq \frac{g(x)}{\varepsilon}.$$

Because $s \in I_u(x)$, it holds that

$$u(x + e) - u(x) \leq s \cdot e \leq u(x) - u(x - e)$$

and hence

$$u(x) - u(x - e) - \frac{g(x)}{\varepsilon} \leq s \cdot e \leq u(x) - u(x - e), \forall s \in I_u(x). \tag{3.3.16}$$

Using (3.3.15) in the first inequality and (3.3.16) in the second, we have that

$$\left(\frac{2\overline{u}}{\text{diam}(\overline{E})} \right)^d \leq \left| \bigcup_{x \in E} I_u(x) \right| \leq \sum_{x \in E} \left(\frac{g(x)}{\varepsilon} \right)^d \mathbf{1}_{\{I_u(x) \neq \emptyset\}}.$$

Hence,

$$\overline{u} \leq C_0(d, \varepsilon) \text{diam}(\overline{E}) |E|^{1/d} \left(\frac{1}{|E|} \sum_{x \in E} |g(x)|^d \mathbf{1}_{\{I_u(x) \neq \emptyset\}} \right)^{\frac{1}{d}},$$

completing the proof of part (a).

(b) It is enough to consider $x_0 = 0$. We begin with some estimates. For parts of the proof, it is easier to work with L_2 (instead of L_∞) balls. Set

$B_R = \{x \in \mathbb{Z}^d : |x|_2 < R\}$. We first deduce from part (a) that for any $p \leq d$ and $\sigma < 1$, there exists a constant $C_1 = C_1(p, \sigma, d)$ such that

$$\max_{x \in B_{\sigma R}} u(x) \leq C_1 \left(\frac{1}{|B_R|} \sum_{x \in B_R} |u^+(x)|^p \right)^{\frac{1}{p}}. \tag{3.3.17}$$

Indeed, define $\eta(x) = \left(1 - \frac{|x|^2}{R^2}\right)^{2d/p}$. A Taylor expansion reveals that for some $C_2 = C_2(p, d)$, it holds that

$$|\eta(x \pm e_i) - \eta(x)| < \frac{C_2}{R}, \quad |\eta(x + e_1) + \eta(x - e_1) - 2\eta(x)| \leq \frac{C_2}{R^2}. \tag{3.3.18}$$

Fix $\kappa_i = \kappa_i(x) \in [0, 1]$, $i = 1, \dots, d$, set $\nu(x) = \eta(x)u(x)$, $x \in B_R$, and

$$\hat{L}_\omega \nu(x) = \sum_{i=1}^d \hat{\omega}(x, x + e_i)(\nu(x + e_i) + \nu(x - e_i) - 2\nu(x))$$

where

$$\hat{\omega}(x, x + e_i) = \begin{cases} \omega(x, x + e_i) \left[\frac{\kappa_i}{\eta(x - e_i)} + \frac{1 - \kappa_i}{\eta(x + e_i)} \right], & |x|^2 \leq R^2 - 4R \\ \omega(x, x + e_i), & R^2 \geq |x|^2 > R^2 - 4R \end{cases}.$$

Then, a tedious computation reveals that, on the set $|x|^2 \leq R^2 - 4R$,

$$\begin{aligned} & -\hat{L}_\omega \nu(x) = -L_\omega u(x) \\ & - 2 \sum_i \frac{\kappa_i(\nu(x + e_i) - \nu(x)) + (1 - \kappa_i)(\nu(x) - \nu(x - e_i))}{\eta(x + e_i)\eta(x - e_i)} [\eta(x + e_i) - \eta(x - e_i)] \\ & + \sum_i \frac{u(x)}{\eta(x + e_i)\eta(x - e_i)} \\ & \left[2(\eta(x + e_i) - \eta(x))(\eta(x) - \eta(x - e_i)) - \eta(x)(\eta(x + e_i) + \eta(x - e_i) - 2\eta(x)) \right] \\ & \leq C_3(d, p) \left[\sum_i \frac{|\kappa_i(\nu(x + e_i) - \nu(x)) + (1 - \kappa_i)(\nu(x) - \nu(x - e_i))|}{R} + \frac{u(x)}{R^2} \right] \end{aligned}$$

where we used (3.3.9) in the first equality.

If for such x , $I_\nu(x) \neq \phi$, then by the proof in part (a), there exists a vector $q \in I_\nu(x)$ with $|q| \leq \frac{\nu(x)}{R - |x|_\infty}$, and one may find a $\kappa_i \in [0, 1]$ such that

$$\kappa_i(\nu(x + e_i) - \nu(x)) + (1 - \kappa_i)(\nu(x) - \nu(x - e_i)) = q_i.$$

Thus, on $\{I_\nu(x) \neq \phi\} \cap \{x : |x|^2 \leq R^2 - 4R\}$, it holds that $-\hat{L}_\omega \nu(x) \leq C_4(d, p) \frac{u(x)}{R^2}$. On the other hand, when $|x|^2 \geq R^2 - 4R$, recalling that $u \geq 0$, it holds that

$$-\left(\nu(x + e_i) + \nu(x - e_i) - 2\nu(x)\right) \leq 2\eta(x)u(x) \leq C_5(d, p) \frac{u(x)}{R^2} \eta(x)^{\frac{d-p}{d}}$$

and in conclusion,

$$-\hat{L}_\omega \nu(x) \leq g(x), \quad x \in B_R$$

where

$$\left|g(x)\mathbf{1}_{I_\nu(x) \neq \phi}\right| \leq \frac{C_6(d, p)u(x)}{R^2}.$$

Applying part (a) of the lemma, we get (3.3.17).

Next, let $\sigma < \tau < 1$, and set

$$\mathbf{u}_\sigma = \min_{x \in B_{\sigma R}} u(x), \quad \mathbf{u}_\tau = \min_{x \in B_{\tau R}} u(x).$$

We claim that (3.3.17) implies the existence of a constant $\gamma = \gamma(d, \sigma, \tau, \varepsilon)$ such that

$$\mathbf{u}_\tau \geq \gamma \mathbf{u}_\sigma. \tag{3.3.19}$$

Indeed, set $\bar{\eta}(x) = (R^2 - |x|^2)^\beta$ with $\beta > 2 \vee 1/\sigma$ and $w(x) = \mathbf{u}_\sigma R^{-2\beta} \bar{\eta}(x) - u(x)$. Then, $w(x) \leq 0$ on $B_{\sigma R} \cup B_R^c$, and $L_\omega w = \mathbf{u}_\sigma R^{-2\beta} L_\omega \bar{\eta}$ on B_R . But, there is an $R_1(\beta)$ such that on $B_R \setminus B_{\sigma R}$, $R > R_1$,

$$L_\omega \bar{\eta}(x) \geq \begin{cases} 0, & |x| < R \\ -C(\beta, d, \varepsilon)R^{2(\beta-1)}, & |x| = R \end{cases},$$

implying by part (a) that on $B_R \setminus B_{\sigma R}$, $R > R_1(\beta)$,

$$w(x) \leq C(\beta, d, \varepsilon) \mathbf{u}_\sigma R^2 R^{-2\beta} \left(\frac{1}{R^d} \sum_{|x|=R} R^{2(\beta-1)d}\right)^{\frac{1}{d}} \leq \frac{C(\beta, d, \varepsilon)}{R^{1/d}} \mathbf{u}_\sigma.$$

Thus,

$$\mathbf{u}_\tau \geq \mathbf{u}_\sigma \left[(1 - \tau^2)^\beta - \frac{C(\beta, d, \varepsilon)}{R^{1/d}} \right].$$

We conclude that there exists an $R_0 = R_0(\sigma, \tau, d, \varepsilon)$ and $\gamma = \gamma(d, \sigma, \tau, \varepsilon)$ such that for all $R > R_0$, (3.3.19) holds. On the other hand, for $R < R_0$ (but $(1 - \tau)R > 1!$), (3.3.19) is trivial by finitely many applications of the equality $L_\omega u = 0$. Thus, (3.3.19) is always satisfied.

A conclusion of (3.3.19) is that if $L_\omega u = 0$ on B_R , $\sigma < 1$, and $\Gamma \subset B_{\sigma R} \subset B_{\tau R} \subset B_R$, letting $\mathbf{u}_\Gamma = \min_{x \in \Gamma} u(x)$, we have that for some $\delta = \delta(\varepsilon, d)$,

$$|\Gamma| \geq \delta |B_{\sigma R}| \implies \mathbf{u}_\tau \geq \gamma \mathbf{u}_\Gamma. \tag{3.3.20}$$

Indeed, define $\nu = \mathbf{u}_\Gamma - u$ and conclude from (3.3.17) that

$$\max_{x \in B_{\sigma R/2}} \nu(x) \leq C_1 \left(\frac{1}{|B_{\sigma R}|} \sum_{x \in B_{\sigma R}} \nu^+(x) \right) \leq C_1(1 - \delta) \max_{x \in B_{\sigma R}} \nu(x)$$

and hence, taking $\delta < 1$ such that $C_1(1 - \delta) < 1/2$,

$$\mathbf{u}_\Gamma - \min_{x \in B_{\sigma R/2}} u(x) \leq C_1(1 - \delta)(\mathbf{u}_\Gamma - \mathbf{u}_\sigma) \leq \frac{1}{2}(\mathbf{u}_\Gamma - \mathbf{u}_\sigma),$$

from which one concludes that $\mathbf{u}_\Gamma \leq \mathbf{u}_{\sigma/2}$. (3.3.20) follows from combining this and (3.3.19).

We finally use the following covering argument. Fix a cube $Q \subset \mathbb{Z}^d$. For $t > 0$, set

$$\Gamma_t = \{x \in Q : u(x) > t\}.$$

Note that if $Q' = Q'(z, r)$ is any cube in \mathbb{Z}^d , centered at z and of side r , (3.3.20) implies that

$$|\Gamma_t \cap Q'| \geq \delta|Q'| \Rightarrow u(x) \geq \gamma t, \text{ some } \gamma = \gamma(\delta, d, \varepsilon). \tag{3.3.21}$$

Define, for any $A \subset Q$,

$$A_\delta = \bigcup_{\substack{\{r, z\} \\ z \in (\frac{1}{2}\mathbb{Z})^d}} \{Q'(z, 3r) \cap Q : |A \cap Q'(z, r)| \geq |Q'(z, r)|\}.$$

Then, cf. [78, Lemma 3] for a proof, either $A_\delta = Q$ or $|A_\delta| \geq |A|/\delta$. Thus, if $|\Gamma_t| \geq \delta^s|Q|$, then iterating (3.3.21) and the above, $\inf_{x \in Q} u(x) \geq \gamma^s t$. Choosing s such that $\delta^s \leq \frac{|\Gamma_t|}{|Q|} \leq \delta^{s-1}$, we conclude that $\inf_{x \in D_R} u(x) \geq \gamma t \left(\frac{|\Gamma_t|}{|D_R|}\right)^{\log \gamma / \log \delta}$. Hence, with $p < \log \delta / \log \gamma := p'$, and $\mathbf{u} = \min_{D_R} u$, we have

$$\begin{aligned} \frac{1}{|D_R|} \sum_{x \in D_R} |u(x)|^p &= p \int_{\mathbf{u}}^\infty t^{p-1} \left(\frac{1}{|D_R|} \sum_{x \in D_R} \mathbf{1}_{u(x) \geq t} \right) dt \\ &= p \int_{\mathbf{u}}^\infty t^{p-1} \left(\frac{|\Gamma_t|}{|D_R|} \right) dt \\ &\leq c(p) \mathbf{u}^{p'} \int_{\mathbf{u}}^\infty \frac{t^{p-1}}{t^{p'}} dt = c(p, p') \mathbf{u}^p, \end{aligned}$$

for some constants $c(p), c(p, p')$, since $p' + 1 - p > 1$. Combining this and (3.3.17) yields the lemma. □

Transience and recurrence of balanced walks

The main result in this section is the following:

Theorem 3.3.22 *Assume Assumption 3.3.1. Then the RWRE $(X_n)_{n \geq 0}$ is transient if $d \geq 3$ and recurrent if $d = 2$.*

Proof. We begin with the transience statement. Fix $d \geq 3$, K large, and define $r_i = K^i$, with $B_i = \{x : |x|_\infty \leq r_i\}$. Set $\tau_0 = 1$ and

$$\tau_i = \min\{n > \tau_{i-1} : X_n \in \partial B_i\}.$$

We use the following uniform estimate on exit probabilities, that actually is stronger than needed: there exists some constant $C = C(\delta, \varepsilon, d) > 0$ such that, if $\Omega_0 = \{\omega : \omega(z, z + e_i) = \omega(z, z - e_i) > \varepsilon, i = 1, \dots, d, \forall z \in \mathbb{Z}^d\}$,

$$\sup_{\omega \in \Omega_0} P_\omega^o(|X_n| < L, n = 1, \dots, L^{2(1+\delta)}) \leq C e^{-CL^{2\delta}}. \tag{3.3.23}$$

There are many ways to prove (3.3.23), including a coupling argument. We use here an optimal control trick. Let $\{B_n\}_{n \geq 0}$ denote a sequence of i.i.d. Bernoulli(1/2) random variables, independent of the environment, of law Q . Then, X_n can be constructed as follows:

$$P_{\omega, B}^o(X_{n+1} = X_n + e_i | X_n = x, X_{n-1}, \dots, X_0) = 2\omega(x, e_i) \mathbf{1}_{2B_{n+1}-1=\pm 1}.$$

(As in Section 3.1, $Q \times P_{\omega, B}^o$, when restricted to $(\mathbb{Z}^d)^\mathbb{N}$, equals P_ω^o .) Set $\mathcal{G}_n = \sigma(B_0, B_1, \dots, B_n, X_0, \dots, X_n, (\omega_z)_{z \in \mathbb{Z}^d})$. An admissible control $\alpha = (\alpha_n)_{n \geq 0}$ is a sequence of \mathcal{G}_n measurable function taking values in $\mathcal{A} := [2\varepsilon, \frac{1}{2} - \varepsilon(d-1)]$. Then define the \mathbb{Z} -valued controlled process (Y_n^α) by $Y_0 = 0$ and

$$P(Y_{n+1}^\alpha = Y_n^\alpha \pm 1 | \mathcal{G}_n, Y_0^\alpha, \dots, Y_n^\alpha) = \alpha_n \mathbf{1}_{2B_{n+1}-1=\pm 1}.$$

Note that, by taking $\hat{\alpha}_n = 2\omega(X_n, e_1)$, we may construct $(Y_n^{\hat{\alpha}})$ and X_n on the same probability space such that $Y_n^{\hat{\alpha}} = X_n$, $Q \times P_{\omega, B}^o$ -a.s. Thus,

$$\begin{aligned} \sup_{\omega \in \Omega_0} P_\omega^o(|X_n|_\infty < L, n = 1, \dots, L^{2(1+\delta)}) \\ \leq \sup_{\omega \in \Omega_0} \sup_{\alpha} Q \times P_{\omega, B}^o(|Y_n^\alpha| < L, n = 1, \dots, L^{2(1+\delta)}). \end{aligned} \tag{3.3.24}$$

Let $g_{n,\omega}(x) = \sup_{\alpha} Q \times P_{\omega, B}(|Y_i^\alpha| < L, i = 1, \dots, n | Y_0 = x)$ (it turns out eventually that $g_{n,\omega}$ does not depend on ω !) Then, due to the Markov property, $g_{n,\omega}(\cdot)$ must satisfy the dynamic programming equation

$$g_{n,\omega}(x) = \begin{cases} \max_{\alpha \in \mathcal{A}} \left(\frac{\alpha}{2} (g_{n+1,\omega}(x+1) + g_{n-1,\omega}(x-1)) + (1-\alpha)g_{n-1,\omega}(x) \right), & |x| < L \\ 0, & |x| \geq L \end{cases}$$

and $g_{0,\omega}(x) = \mathbf{1}_{|x| < L}$. Next, we note that $g_{n,\omega}(\cdot)$ satisfies

$$g_{n,\omega}(x+1) + g_{n,\omega}(x-1) - 2g_{n,\omega}(x) \leq 0. \tag{3.3.25}$$

For $n = 0$ this is immediate, and hence

$$g_{1,\omega}(x) = g_{0,\omega}(x) + \varepsilon \left(g_{0,\omega}(x+1) + g_{0,\omega}(x-1) - 2g_{0,\omega}(x) \right).$$

We then have that $g_{1,\omega}(x)$ satisfies (3.3.25), and the argument can be iterated. We further conclude that

$$g_{n,\omega}(x) = g_{n-1,\omega}(x) + \varepsilon \left(g_{n-1,\omega}(x+1) + g_{n-1,\omega}(x-1) - 2g_{n-1,\omega}(x) \right). \quad (3.3.26)$$

Thus, $g_{n,\omega}(x)$ is nothing but the probability that a simple random walk on \mathbb{Z} with geometric $(1 - 2\varepsilon)$ holding times, stays confined in a strip of size L for $L^{2(1+\delta)}$ units of time (note that (3.3.26) possesses a unique solution, which does not depend on $\omega \in \Omega_0!$). The conclusion (3.3.23) follows from solving (3.3.26) and combining it with (3.3.24).

From (3.3.23), we conclude that $E_\omega^o(\tau_{i+2}) \leq Cr_{i+2}^{2(1+\delta)}$, for all i large enough, all $\omega \in \Omega_0$, where $\mathbb{P}(\Omega_0) = 1$. Thus,

$$\begin{aligned} Cr_{i+2}^{2(1+\delta)} &\geq E_\omega^o \left(E_\omega^o(\# \text{ visits of } X_n \text{ at } B_{i-1} \text{ for } n \in (\tau_i + 1, \dots, \tau_{i+2}) | X_{\tau_i}) \right) \\ &= E_\omega^o \left(\sum_{y \in B_{i-1}} E_\omega^{X_{\tau_i}}(\# \text{ visits at } y \text{ before } \tau_{i+2}) \right) \\ &\geq \sum_{y \in B_{i-1}} E_\omega^o \left(E_{\theta^{-y}\omega}^{X_{\tau_i-y}}(\# \text{ visits at } 0 \text{ before } \tau_{i+1}) \right) \\ &\geq C \sum_{y \in B_{i-1}} \max_{z \in E_i} \left(E_{\theta^{-y}\omega}^z(\# \text{ of visits at } 0 \text{ before } \tau_{i+1}) \right) \end{aligned}$$

where $E_i = \{x : \frac{r_i}{2} < |x|_\infty < \frac{3r_i}{2}\}$, and Harnack’s inequality (Lemma 3.3.8) was used in the last step. Taking P -expectations, we conclude that

$$\begin{aligned} Cr_{i+1}^{2(1+\delta)} &\geq C \sum_{y \in B_{i-1}} \mathbb{E}^o \left(\max_{z \in E_i} E_{\theta^{-y}\omega}^z(\# \text{ of visits at } 0 \text{ before } \tau_{i+1}) \right) \\ &\geq C \sum_{y \in B_{i-1}} \mathbb{E}^o \left(E_{\theta^{-y}\omega}^{X_{\tau_i}}(\# \text{ of visits at } 0 \text{ before } \tau_{i+1}) \right) \\ &= C \sum_{y \in B_{i-1}} \mathbb{E}^o \left(E_\omega^{X_{\tau_i}}(\# \text{ of visits at } 0 \text{ before } \tau_{i+1}) \right) \\ &= C'(r_{i-1})^d \mathbb{E}^o \left(E_\omega^{X_{\tau_i}}(\# \text{ of visits at } 0 \text{ before } \tau_{i+1}) \right), \end{aligned}$$

where the shift invariance of P was used in the next to last equality. Therefore,

$$\mathbb{E}^o(\# \text{ of visits at } 0 \text{ between } \tau_i + 1 \text{ and } \tau_{i+1}) \leq C''r_i^{2+\delta-d}.$$

Hence, for $d \geq 3$,

$$\mathbb{E}^o(\# \text{ of visits at } 0) \leq C'' \sum_{i=1}^{\infty} r_i^{2+\delta-d} < \infty,$$

implying that P -a.s., $E_{\omega}^o(\# \text{ of visits at } 0) < \infty$, i.e. (X_n) is transient if $d \geq 3$.

Turning to $d = 2$, we recall the following lemma:

Lemma 3.3.27 (Derrienic[20]) *Let (Y_i) be a stationary and ergodic lattice valued sequence, and set $S_n = \sum_{i=1}^n Y_i$. Define*

$$R_n = \{ \# \text{ of sites visited up to time } n \}.$$

Then,

$$\frac{R_n}{n} \xrightarrow[n \rightarrow \infty]{} \text{Prob}(S_i \neq 0, i \geq 1).$$

Proof. The sequence R_n is sub-additive and hence, by Kingman’s ergodic sub-additive theorem, $R_n/n \rightarrow_{n \rightarrow \infty} a$, a.s. and in L^1 , for some constant a . Noting that $R_{n+1} = R_n \circ \theta + \mathbf{1}_{Y_1 \notin \{\cup_{i=2}^{n+1} S_i\}}$, it holds that $(R_{n+1} - R_n \circ \theta) \xrightarrow[n \rightarrow \infty]{} \mathbf{1}_A \circ \theta$, where $A = \{S_i \neq 0, i \geq 1\}$. Thus, $ER_n/n \xrightarrow[n \rightarrow \infty]{} E\mathbf{1}_A =: a$. \square

Under the measure on the environment Q introduced in this section, the increments $\{X_{n+1} - X_n\}$ are stationary and ergodic. Letting R_n denote the range of the RWRE up to time n , we have that

$$\frac{R_n}{n} \xrightarrow[n \rightarrow \infty]{} Q \times P_{\omega}^o(\text{no return to } 0), \quad Q\text{-a.s.}$$

But, due to the CLT (Theorem 3.3.4 and Remark 3.3.5), for any $\delta > 0$,

$$\liminf_{n \rightarrow \infty} P_{\omega}^o \left(\frac{R_n}{n} < \delta \right) > 0, \quad Q\text{-a.s.}$$

Hence, for any $\delta > 0$,

$$P_{\omega}^o(\text{no return to } 0) < \delta, \quad Q\text{-a.s.}$$

and hence also P -a.s. This concludes the recurrence proof. \square

Remark: It is interesting to note that the transience (for $d \geq 3$) and recurrence (for $d = 2$) results are *false* for certain balanced, elliptic environments in Ω_0 (however, the P -probability of these environments is, of course, null). A simple example that exhibits the failure of recurrence for $d = 2$ was suggested by N. Gantert: fix $0.25 < p < 0.5$ and $q = 0.5 - p$. With $x = (x_1, x_2) \in \mathbb{Z}^2$, define

$$\omega(x, e) = \begin{cases} \frac{1}{4}, & x_1 = x_2, |e| = 1 \\ p, & \begin{cases} e = \pm e_2, |x_1| > |x_2| \\ \text{or} \\ e = \pm e_1, |x_1| < |x_2| \end{cases} \\ q, & \begin{cases} e = \pm e_1, |x_1| > |x_2| \\ \text{or} \\ e = \pm e_2, |x_1| < |x_2| \end{cases} \end{cases}$$

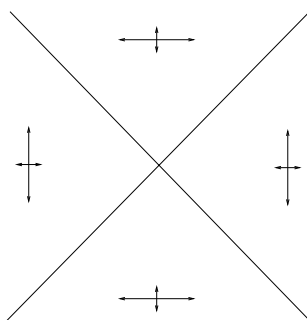


Fig. 3.3.2. Transient balanced environment, $d = 2$

Define

$$\nu(x) = \begin{cases} 1, & x \neq 0 \\ 4q, & x = 0. \end{cases}$$

Then, $\nu(\cdot)$ is an excessive measure, i.e.

$$(L_\omega^* \nu)(x) := \sum_{e:|e|=1} \omega(x - e, e) \nu(x - e) \leq \nu(x), \quad x \in \mathbb{Z}^2.$$

If $\{X_n\}$ was recurrent, then every excessive measure needs to equal the (unique) invariant measure. But, with

$$\nu((1, 0)) = 1 > (\nu L_\omega)((1, 0)) = 2q + 0.5.$$

Thus, $\nu(\cdot)$ is not invariant, contradicting the recurrence of the chain.

The intuitive idea behind the example above is that for points far from the origin, the “radial component” of the walk behaves roughly like a Bessel process of dimension $2 + \delta$, some $\delta > 0$, implying the transience. A similar argument, only more complicated, allows one to construct environments in $d \geq 3$ where the radial component behaves like a Bessel process of dimension $2 - \delta$, some $\delta > 0$. It is not hard to prove, using Lyapunov functions techniques, that there exists a $\kappa(d) < 1/2d$ such that if $d \geq 3$ and the balanced environment is such that $\min_{e:|e|=1} \omega(x, e) > \kappa(d)$ then the walk is transient.

Bibliographical notes: The basic CLT under Assumption 3.3.1 is due to Lawler [47], who transferred to the discrete setting some results of Papanicolau and Varadhan. An extension to the case of non nearest neighbour walks appears in [48]. The Harnack principle (Lemma 3.3.8) was provided in [49], and in greater generality in [46], whose approach we follow, after a suggestion by Sznitman (see also [69]).

The proof of the transience part in Theorem 3.3.22 was suggested by G. Lawler in private communication. The proof of the recurrence part is due to H. Kesten, also in private communication. A recent independent proof appears

in [8]. Finally, the examples mentioned at the end of the section go back to Krylov (in the context of diffusions), with this version based on discussions with Comets and Gantert.

We comment that there are very few results on LLN's and CLT's for non balanced, non ballistic walks. One exception is the result in [9], where renormalization techniques are used to prove a (quenched) CLT in symmetric (not-balanced!) environments with small disorder. Another case, in which some of the RWRE coordinates perform a simple random walk, is analysed in details in [4], using cut-times of the random walk instead of the regeneration times used in Section 3.5.

3.4 Large deviations for nestling walks

In this section, we derive an LDP for a class of nearest neighbour random walks in random environment, in \mathbb{Z}^d . For reasons that will become clearer below, we need to restrict attention to environments which satisfy a condition on the support of P , which we call, after M. Zerner, “nestling environments”. For technical reasons, we also need to make an independence assumption (see however the remark at the end of this section).

Define $d(\omega) := \sum_{e:|e|=1} \omega(0, e)e$, and let $P_d := P \circ d^{-1}$ denote the law of $d(\omega)$ under P .

Assumption 3.4.1

(C1) P is i.i.d.

(C2) P is elliptic: there exists an $\varepsilon > 0$ such that $P(\omega(z, z + e_i) \geq \varepsilon) = P(\omega(z, z - e_i) \geq \varepsilon) = 1, i = 1, \dots, d$.

(C3) (Nestling property): $0 \in \text{conv}(\text{supp}(P_d))$.

We elaborate below on the nestling assumption. Clearly, balanced walks are nestling, but one may construct examples, as in $d = 1$, of nestling environments with ballistic behaviour.

For any $y \in \mathbb{R}^d$, we denote by $[y]$ the point in \mathbb{Z}^d with $1 > y_i - [y]_i \geq 0$. For $z \in \mathbb{Z}^d$, we let $T_z = \inf\{n \geq 0 : X_n = z\}$. As in Section 2.3, the key to our approach to large deviation results for (X_n) is a large deviation principle for $T_{[nz]}, z \in \mathbb{R}^d$, stated next.

Theorem 3.4.2 (a) Assume P is ergodic and elliptic. For any $z \in \mathbb{R}^d, |z|_1 = 1$, any $\lambda \leq 0$, the following deterministic limit exists P -a.s.

$$a(\lambda, z) := \lim_{n \rightarrow \infty} \frac{1}{n} \log E_\omega^o(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty}).$$

(b) Further assume Assumption 3.4.1, and define

$$I_{T,z}(s) = \sup_{\lambda < 0} (\lambda s - a(\lambda, z)).$$

Then, $T_{[nz]}/n$ satisfies, P -a.s., under P_ω^o , a (weak) LDP, with rate function $I_{T,z}(s)$. That is,

$$\begin{aligned} & \lim_{\delta \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(T_{[nz]}/n \in (s - \delta, s + \delta)) \\ &= \lim_{\delta \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(T_{[nz]}/n \in (s - \delta, s + \delta)) = -I_{T,z}(s), \quad P - a.s. \end{aligned} \tag{3.4.3}$$

(Note that $I_{T,z}(s) = \infty$ for $s < 1$).

With Theorem 3.4.2 at hand, we may state the LDP for X_n/n . Define, for $x \in \mathbb{R}^d$,

$$I(x) = \begin{cases} |x|_1 I_{T,x/|x|_1}(1/|x|_1), & |x|_1 \leq 1 \\ \infty, & \text{otherwise} \end{cases}.$$

Obviously, $a(\lambda, z)$ is defined for any $z \in \mathbb{R}^d \setminus \{0\}$, and is by definition homogeneous in $|z|_1$. An easy computation then reveals that $I(x) = \sup_{\lambda < 0} (\lambda - a(\lambda, x))$. We have the

Theorem 3.4.4 *Assume Assumption 3.4.1. Then, P -a.s., the random variables X_n/n satisfy the LDP in \mathbb{R}^d with good, convex rate function $I(\cdot)$. That is,*

$$\begin{aligned} \lim_{\delta \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o \left(\frac{X_n}{n} \in B_x(\delta) \right) &= \lim_{\delta \rightarrow \infty} \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o \left(\frac{X_n}{n} \in B_x(\delta) \right) \\ &= -I(x), \quad P - a.s. \end{aligned}$$

Proof of Theorem 3.4.2

The idea behind the proof is relatively simple, and is related to our proof of large deviations for $d = 1$. However, there are certain complications in the proof of the lower bound, which can be overcome at present only under the nestling assumption.

a) We begin by defining, for $\lambda \leq 0$,

$$a_{n,m}(\lambda, z) := \log E_\omega^{[mz]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \right).$$

We then have (since the time to reach $[nz]$ is not larger than the time to reach $[nz]$, when one is forced also to first visit $[mz]$), that

$$a_{n,0}(\lambda, z) \geq a_{m,0}(\lambda, z) + a_{n,m}(\lambda, z).$$

Further, we note that due to **C2**, there exists a constant $C(\lambda, \varepsilon)$ such that $n^{-1}|a_{n,0}(\lambda, z)| \leq C(\lambda, \varepsilon)$, for all ω with $\omega(x, x + e) \geq \varepsilon$, all $x \in \mathbb{Z}^d$ and e such that $|e| = 1$. Thus, by Kingman’s subadditive ergodic theorem,

$$\frac{a_{n,0}(\lambda, z)}{n} \xrightarrow[n \rightarrow \infty]{} a(\lambda, z), \quad P\text{-a.s.} \tag{3.4.5}$$

b) By Chebycheff’s inequality, (3.4.5) immediately implies the upper bound in (3.4.3). Thus, all our effort is now concentrated in proving the lower bound.

Toward this end, note that by Jensen’s inequality, the deterministic function $a(\cdot, z)$ is convex, and thus differentiable a.e. We denote by \mathcal{D} the set of $\lambda < 0$ such that $a(\cdot, z)$ is differentiable at λ . Recall that a point $s \in \mathbb{R}_+$ is an exposed point of $I_{T,z}(\cdot)$ if for some $\lambda < 0$ (“the exposing plane”) and all $t \neq s$,

$$\lambda t - I_{T,z}(t) > \lambda s - I_{T,z}(s).$$

It is straightforward to check, see e.g., [19, Lemma 2.3.9(b)] that if $y = a'(\lambda, z)$ for some $\lambda \in \mathcal{D}$, then $I_{T,z}(y) = \lambda y - a(\lambda, z)$, and further y is an exposed point of $I_{T,z}(\cdot)$, with exposing plane λ .

As we already saw, it is then standard (see, e.g., [19, Theorem 2.3.6(b)]) that the lower bound in (3.4.3) holds for any exposed point. Thus, it only remains to handle points which are not exposed. Toward this end, define (using the monotonicity to ensure the existence of the limit!)

$$s_+ := \lim_{\lambda \rightarrow 0, \lambda \in \mathcal{D}} a'(\lambda, z) \leq \infty.$$

Note that, for any $s \geq s_+$, $I_{T,z}(s) = -\lim_{\lambda \rightarrow 0} a(\lambda, z)$.

The approach toward the lower bound is different when $s \geq s_+$ (case a) and $s < s_+$ (case b): in case a, a strategy which will achieve a lower bound consists of spending first some time in a “trap” at the neighborhood of the origin, returning to the origin and then getting to $[nz]$ within time roughly ns_-^n , where $s_-^n < s_+$ is an exposed point with $|I_{T,z}(s_-^n) - I_{T,z}(s_+)| \leq \eta$. The nestling assumption is crucial to create the trap. In case b, the achieving strategy consists of finding an intermediate point, progressing faster than needed toward the intermediate point, and then progressing slower than expected toward $[nz]$. To control the behavior of the walk starting at intermediate points, the independence assumption comes in handy.

Turning to case a, the role of the nestling assumption is evident in the following lemma:

Lemma 3.4.6 *Assume Assumption 3.4.1. Then, there exists an $\Omega_0 \subset \Omega$ with $P(\Omega_0) = 1$ with the following property: for each $\delta > 0$ and each $\omega \in \Omega_0$, there exists an $R(\delta, \omega)$ and an $n_0 = n_0(\delta, \omega)$ such that, for any $n > n_0$ even,*

$$P_\omega^o(|X_m|_2 \leq R(\delta, \omega), \quad m = 1, \dots, n - 1, X_n = 0) \geq e^{-\delta n}.$$

Proof of Lemma 3.4.6

We begin by constructing a “trap”. As a preliminary, with $x \in \mathbb{Z}^d$, and $\bar{\tau}$ such that $|x_{\bar{\tau}}| \geq |x_j|$, $j = 1, \dots, d$ (and hence $|x|_2 \leq \sqrt{d}|x_{\bar{\tau}}|$) we have, defining $y_x = x - \text{sign}(x_i)e_i$, that

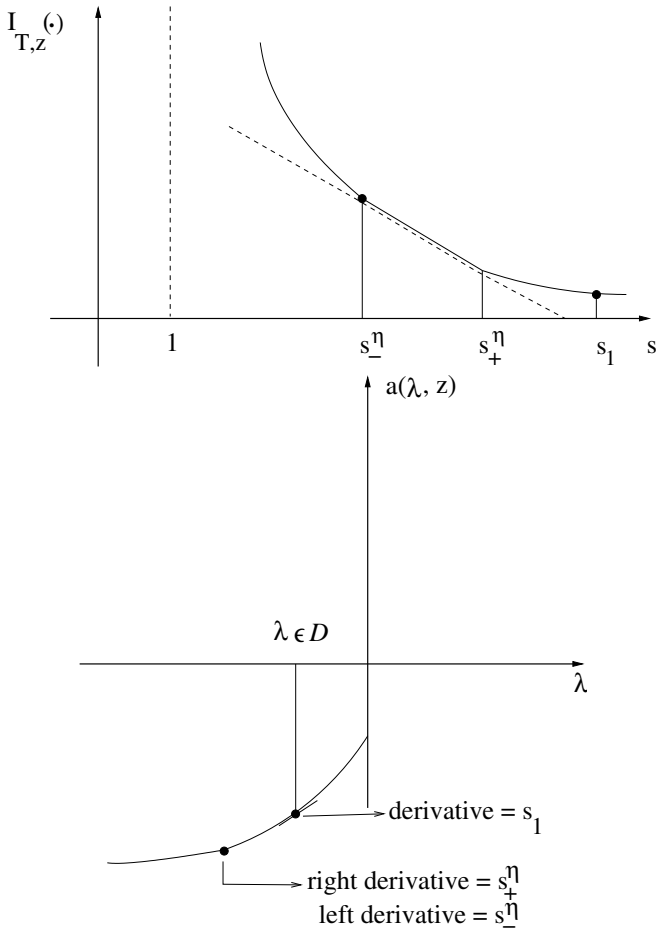


Fig. 3.4.1. exposed points and differentiability

$$|x|_2 - |y_x|_2 = \frac{|x_{\bar{t}}|^2 - (|x_{\bar{t}}| - 1)^2}{|x|_2 + |y_x|_2} \geq \frac{1}{2\sqrt{d}}.$$

Fix $\kappa = \varepsilon\delta/32\sqrt{d}$ and $F(x) = (1 - \kappa^2|x|_2^2) \vee 0$. Call a site $x \in \mathbb{Z}^d$ “successful” if

$$x \cdot \sum_{i=1}^d (\omega(x, x + e_i)e_i - \omega(x, x - e_i)e_i) \leq 1.$$

Due to **(C3)**,

$$P(x \text{ is successful}) > 0.$$

and hence, by the independence assumption **(C1)**,

$$P(\text{all sites } x \in B_{1/\kappa}(0) \text{ are successful}) > 0. \tag{3.4.7}$$

Fix now $\omega \in \Omega$ such that all sites $x \in B_{1/\kappa}(0)$ are successful. We next claim that for such ω ,

$$\sum_i \omega(x, x \pm e_i) F(x \pm e_i) \geq e^{-\delta/3} F(x). \tag{3.4.8}$$

Indeed, for $|x|_2 \geq 1/\kappa$ this is obvious, for $1/\kappa - \varepsilon/4\sqrt{d} < |x|_2 < 1/\kappa$ this follows from the ellipticity assumption **(C2)** while for $|x|_2 < 1/\kappa - \varepsilon/4\sqrt{d}$ this follows from a Taylor expansion. Thus, $e^{\delta n/3} F(X_n)$ is, for such ω , a submartingale under P_ω^o , and we have that for all $n \geq 1$,

$$e^{-\delta n/3} = e^{-\delta n/3} E_\omega^o F(X_0) \leq e^{-\delta n/3} E_\omega^o \left(e^{\delta n/3} F(X_n) \right) \leq P_\omega^o \left(|X_n|_2 < \frac{1}{\kappa} \right). \tag{3.4.9}$$

Fixing n_1 even large enough such that $e^{-\delta n_1/3} \varepsilon \sqrt{d}/\kappa \geq e^{-2\delta n_1/3}$, we conclude that for such ω ,

$$P_\omega^o \left(|X_m|_\infty \leq n_1, m = 1, \dots, n_1 - 1, X_{n_1} = 0 \right) \geq e^{-2\delta n_1/3}.$$

Due to (3.4.7) and **(C1)**, there exists (P -a.s.) an $x_0 = x_0(\omega, \delta)$ such that all sites in $B_{1/\kappa}(x_0)$ are successful. Set $m_0 = m_0(\omega, \delta) := \sum_{i=1}^d |x_0(\omega, \delta)(i)|$. Due to the ellipticity assumption **(C2)**, we have

$$P_\omega^o \left(X_{m_0} = x_0(\omega, \delta) \right) \geq \varepsilon^{m_0}, \quad P_\omega^{x_0} (X_{m_0} = 0) \geq \varepsilon^{m_0}.$$

Next set $R(\delta, \omega) := n_1 + 2m_0 + 1$. Define $K = \lfloor (n - 2m_0)/n_1 \rfloor$. We then have, using the Markov property, that

$$\begin{aligned} & P_\omega^o \left(|X_m|_2 \leq R(\delta, \omega), \quad m = 1, \dots, n, \quad X_n = 0 \right) \\ & \geq P_\omega^o \left(X_{m_0} = x_0(\omega, \delta) \right) P_\omega^{x_0} \left(|X_m - x_0|_\infty \leq n_1, X_{n_1} = 0 \right)^K \\ & \quad P_\omega^{x_0} (X_{m_0} = 0) P_\omega^o (X_{n - Kn_1 - 2m_0} = 0) \\ & \geq \varepsilon^{2m_0} \varepsilon^{n_1} \cdot e^{-\frac{2\delta}{3} Kn_1} \geq e^{-\delta n}, \end{aligned}$$

for all $n > n_0(\delta, \varepsilon, \omega)$. □

Equipped with Lemma 3.4.6 we may complete the proof in case a. Indeed, all we need to prove is that for any $\delta > 0$,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o \left(T_{\lfloor nz \rfloor} / n \in (s - \delta, s + \delta) \right) = -I_{T,z}(s_+), \quad P - a.s.$$

Fix $\eta > 0$ and an exposed point s_η^- with $|I_{T,z}(s_\eta^-) - I_{T,z}(s_+)| \leq \eta$. Due to the Markov property, for all n such that $\lfloor nz \rfloor_\infty > R(\delta, \omega)$,

$$\begin{aligned}
 & P_\omega^o \left(T_{[nz]}/n \in (s - \delta, s + \delta) \right) \\
 & \geq P_\omega^o \left(|X_m|_2 \leq R(\delta, \omega), m = 1, \dots, \lceil n(s - s_-^\eta) \rceil, X_{\lceil n(s - s_-^\eta) \rceil} = 0 \right) \\
 & \quad P_\omega^o \left(T_{[nz]}/n \in (s_-^\eta - \delta', s_-^\eta + \delta') \right)
 \end{aligned}$$

where $\delta' = \delta \cdot s_-^\eta / 2s$, and hence,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o \left(T_{[nz]}/n \in (s - \delta, s + \delta) \right) \geq -\delta - I_{T,z}(s_-^\eta).$$

Since δ is arbitrary and $I_{T,z}(s_-^\eta) \xrightarrow{\eta \rightarrow 0} I_{T,z}(s_+)$, the proof is concluded for $s \geq s_+$.

Turning to case b, recall that our plan is to consider intermediate points. This requires a slight strengthening of the convergence of $a(\lambda, z)$. We state this in the following Lemma, whose proof is deferred.

Lemma 3.4.10 *Assume Assumption 3.4.1 and set $\nu \in (0, 1)$. Then, for $z \in \mathbb{R}^d$, $|z|_1 = 1$, and any $\lambda < 0$, we have*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log E_\omega^{\lfloor \nu n z \rfloor} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \right) = (1 - \nu)a(\lambda, z), \quad P - a.s.$$

Assuming Lemma 3.4.10, we complete the proof of part b. Note the existence, for any $\eta \geq 0$, of $s_-^\eta < s < s_+^\eta$ such that s_-^η, s_+^η are exposed, and further

$$\left| I_{T,z}(s) - \left(\frac{s - s_-^\eta}{s_+^\eta - s_-^\eta} \right) I_{T,z}(s_-^\eta) - \left(\frac{s_+^\eta - s}{s_+^\eta - s_-^\eta} \right) I_{T,z}(s_+^\eta) \right| < \eta. \quad (3.4.11)$$

Set $\nu := (s_+^\eta - s) / (s_+^\eta - s_-^\eta)$. By the Markov property,

$$\begin{aligned}
 & P_\omega^o \left(T_{[nz]}/n \in (s - \delta, s + \delta) \right) \\
 & \geq P_\omega^o \left(T_{\lfloor \nu n z \rfloor} / n \in (s_-^\eta - \delta', s_-^\eta + \delta') \right) P_\omega^{\lfloor \nu n z \rfloor} \left(T_{[nz]}/n \in (s_+^\eta - \delta', s_+^\eta + \delta') \right)
 \end{aligned}$$

where $\delta' = \frac{\min(\nu, 1-\nu)\delta}{2}$. Due to Lemma 3.4.10, and the fact that s_+^η, s_-^η are exposed points of $I_{T,z}(\cdot)$, one concludes that

$$\begin{aligned}
 & \lim_{\delta \rightarrow \infty} \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o \left(T_{[nz]}/n \in (s - \delta, s + \delta) \right) \\
 & \geq - \left[\left(\frac{s - s_-^\eta}{s_+^\eta - s_-^\eta} \right) I_{T,z}(s_-^\eta) + \left(\frac{s_+^\eta - s}{s_+^\eta - s_-^\eta} \right) I_{T,z}(s_+^\eta) \right].
 \end{aligned}$$

Using (3.4.11), this completes the proof of the theorem. □

Proof of Theorem 3.4.4

Fix x and δ as in the statement of the theorem. Then, using the ellipticity assumption **(C2)**, for any n large enough,

$$P_\omega^o\left(\frac{X_n}{n} \in B_x(\delta)\right) \geq P_\omega^o\left(T_{[nx]} \in n\left(1 - \frac{\delta}{2}, 1 + \frac{\delta}{2}\right)\right) \varepsilon^{n\delta/2}$$

and the lower bound follows from Theorem 3.4.2.

Turning to the upper bound, note that $|nB_x(\delta) \cap \mathbb{Z}^d| \leq C_\delta n^d$, and that

$$P_\omega^o\left(\frac{X_n}{n} \in B_x(\delta)\right) = \sum_{y \in nB_x(\delta) \cap \mathbb{Z}^d} P_\omega^o(X_n = y).$$

Further, note that $P_\omega^o(X_n = y) \leq P_\omega^o(T_{[y]} \leq n)$, and that due to the ellipticity **(C2)**,

$$\sup_{y \in nB_x(\delta)} P_\omega^o(X_n = y) \leq \varepsilon^{-n\sqrt{d}\delta} P_\omega^o\left(T_{[nx]} \leq n(1 + \delta)\right)$$

and hence,

$$\begin{aligned} & \lim_{\delta \rightarrow \infty} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(\frac{X_n}{n} \in B_x(\delta)\right) \\ & \leq \lim_{\delta \rightarrow \infty} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o\left(T_{[nx]} \leq n(1 + \delta)\right) \\ & \leq - \inf_{0 \leq \eta \leq 1} I(\eta x), \quad P - a.s. \end{aligned}$$

The monotonicity of $I(\eta \cdot)$ in η , which is induced from that of $I_{T,z}(\cdot)$, completes the proof. \square

Proof of Lemma 3.4.10

By the ellipticity assumption **(C2)**, $\frac{1}{n} \log E_\omega^{[\nu n z]}(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty})$ is uniformly bounded. Further, it possesses the same law as $\frac{1}{n} \log E_\omega^o(e^{\lambda T_{[n(1-\eta)z]}} \mathbf{1}_{T_{[n(1-\eta)z]} < \infty})$. Thus,

$$\frac{1}{n} \log E_\omega^{[\nu n z]}(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty}) \xrightarrow[n \rightarrow \infty]{P} (1 - \nu)a(\lambda, z). \tag{3.4.12}$$

Our goal is thus to prove that the convergence in (3.4.12) is in fact a.s.

Toward this end, as a first step we truncate appropriately the expectation.

Set

$$N_x = \#\{\text{visits at } x \text{ before } T_{[nz]}\},$$

and $N = \sup_{x \in \mathbb{Z}^d} N_x$. We show that, for some $\delta < 1$,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \frac{E_\omega^{[\nu n z]}(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty})}{E_\omega^{[\nu n z]}(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \mathbf{1}_{N < n^\delta})} = 0, \quad P - a.s. \tag{3.4.13}$$

Indeed, note first that

$$E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \right) \leq E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \mathbf{1}_{N < n^\delta} \right) + \sum_{x \in \mathbb{Z}^d} E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \mathbf{1}_{N_x > n^\delta} \right).$$

But, due to the Markov property,

$$\begin{aligned} & E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \mathbf{1}_{N_x > n^\delta} \right) \\ & \leq \sum_{k=[n^\delta]+1}^\infty E_\omega(e^{\lambda T_x} \mathbf{1}_{T_x < T_{[nz]}})^k E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_x < T_{[nz]} < \infty} \right) \\ & \leq \frac{e^{\lambda n^\delta}}{1 - e^\lambda} E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \right), \end{aligned}$$

and hence,

$$E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \right) \leq \frac{E_\omega^{[\nu n z]} \left(e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \mathbf{1}_{N < n^\delta} \right)}{1 - n^d e^{\lambda n^\delta} / (1 - e^\lambda)},$$

yielding (3.4.13). Further, due to the ellipticity assumption **(C2)**, it holds that for some constant $K = K(\lambda)$ large enough,

$$E_\omega^{[\nu n z]} (e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < Kn} \mathbf{1}_{N < n^\delta}) \geq E_\omega^{[\nu n z]} (e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < \infty} \mathbf{1}_{N < n^\delta}) / 2.$$

Thus, it suffices to consider

$$g_\omega^\delta = \log E_\omega^{[\nu n z]} (e^{\lambda T_{[nz]}} \mathbf{1}_{T_{[nz]} < Kn} \mathbf{1}_{N < n^\delta}).$$

Denote by $\mathcal{P}_{k,\delta}$ the set of nearest neighbour paths (γ_n) on \mathbb{Z}^d with $\gamma_0 = [nz]$, $\gamma_k = [nz]$, and $N(\gamma) \leq n^\delta$. For $e \in \{\pm e_i\}_{i=1}^d =: \mathcal{E}$, set

$$N_{x,e}(\gamma) = \#\{\text{steps from } x \text{ to } x + e \text{ of } \gamma \text{ before } T_{[nz]}(\gamma)\}.$$

Then, with $\beta(x, x + e) = \log \omega(x, x + e)$ and $D_{Kn} = \{x \in \mathbb{Z}^d : |x|_\infty \leq Kn\}$,

$$g_\omega^\delta = \log \sum_{k \leq Kn} e^{\lambda k} \sum_{\gamma \in \mathcal{P}_{k,\delta}} \prod_{x \in D_{Kn}} \prod_{e \in \mathcal{E}} e^{\beta(x, x+e) N_{x,e}(\gamma)}.$$

We use the following concentration inequality, which is a slight variant of [77, Theorem 6.6].

Lemma 3.4.14 (Talagrand) *Let $\mathcal{K} \subset \mathbb{R}^{d_1}$ be compact and convex. Let μ be a law supported on \mathcal{K} , and let $f : \mathcal{K}^N \rightarrow \mathbb{R}$ be convex and of Lipschitz constant L . Finally, let M_N denote the median of f with respect to $\mu^{\otimes N}$, i.e. M_N is the smallest number such that*

$$\mu^{\otimes N}(f \leq M_N) \geq \frac{1}{2}, \quad \mu(f \geq M_N) \geq \frac{1}{2}.$$

Then, there exists a constant $C = C(\mathcal{K})$, independent of f, μ , such that for all $t > 0$,

$$\mu^{\otimes N}(|f - M_N| \geq t) \leq C \exp(-Ct^2/L^2).$$

To apply Lemma 3.4.14, note that

$$\begin{aligned} & \left| \frac{\partial g_\omega^\delta}{\partial \beta(x, x + e)} \right| \\ & \leq \frac{1}{\varepsilon} \frac{\sum_{k \leq K_n} e^{\lambda k} \sum_{\gamma \in \mathcal{P}_{k, \delta}} N_{x, e}(\gamma) \prod_{x' \in D_{K_n}} \prod_{e \in \mathcal{E}} e^{\beta(x', x' + e) N_{x', e}(\gamma)}}{\sum_{k \leq K_n} e^{\lambda k} \sum_{\gamma \in \mathcal{P}_{k, \delta}} \prod_{x' \in D_{K_n}} \prod_{e \in \mathcal{E}} e^{\beta(x', x' + e) N_{x', e}(\gamma)}}. \end{aligned} \tag{3.4.15}$$

Thus, using Jensen’s inequality in the first inequality,

$$\begin{aligned} \sum_{x \in D_{K_n}} \left| \frac{\partial g_\omega^\delta}{\partial \beta(x, x + e)} \right|^2 & \leq \frac{1}{\varepsilon} \sum_{x \in D_{K_n}} \\ & \frac{\sum_{k \leq K_n} e^{\lambda k} \sum_{\gamma \in \mathcal{P}_{k, \delta}} N_{x, e}(\gamma)^2 \prod_{x' \in D_{K_n}} \prod_{e \in \mathcal{E}} e^{\beta(x', x' + e) N_{x', e}(\gamma)}}{\sum_{k \leq K_n} e^{\lambda k} \sum_{\gamma \in \mathcal{P}_{k, \delta}} \prod_{x' \in D_{K_n}} \prod_{e \in \mathcal{E}} e^{\beta(x', x' + e) N_{x', e}(\gamma)}} \\ & \leq \frac{K_n^{1+\delta}}{\varepsilon}. \end{aligned}$$

It is immediate to see that on the other hand g_ω^δ is a convex function of $\{\beta(x, x + e)\}$. Hence, by Lemma 3.4.14 and the above,

$$P(|g_\omega^\delta - E g_\omega^\delta| > tn) \leq C_1 e^{-C_1 n^{1-\delta}},$$

where $C_1 = C_1(\varepsilon, \delta)$. The Borel-Cantelli lemma then completes the proof of Lemma 3.4.10. \square

Remarks: 1. In the proof above, the independence assumption **(C1)** was used in two places. The first is the construction of traps (Lemma 3.4.6), where the independence assumption may be replaced by the requirement that P , when restricted to finite subsets, be equivalent to a product measure. More seriously, the product structure was used in the application of Talagrand’s Lemma 3.4.14. It is plausible that this can be bypassed, e.g. using the techniques in [68].

2. S. R. S. Varadhan has kindly indicated to me a direct argument which gives the quenched LDP for the position, for ergodic environments, without passing through hitting times. Fix $\varepsilon > 0$, and define X_n^ε to be the RWRE with geometric holding times of parameter $1/\varepsilon$. Fix a deterministic v , with $|v|_1 < 1$, and define

$$g(m, n) = P_{\theta^{[mv]_\omega}}^0(X_{m-n}^\varepsilon - X_0^\varepsilon = [(n - m)v]).$$

Then, $g(0, n + m) \geq g(0, m)g(m, n + m) > 0$ for all $n, m \geq 1$. Consequently, by Kingman’s ergodic sub-additive theorem,

$$\frac{1}{n} \log g(0, n) \rightarrow_{n \rightarrow \infty} -I^\epsilon(v), P - a.s.,$$

for some deterministic $I^\epsilon(v)$. From this it follows (see e.g. [19, Theorem 4.1.11]) that X_n^ϵ/n satisfies the (quenched) LDP with convex, good rate function $I^\epsilon(\cdot)$. Finally, it is easy to check that $I^\epsilon(\cdot) \rightarrow_{\epsilon \rightarrow 0} I(\cdot)$ (even uniformly on compacts) and that

$$\limsup_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_\omega^o(|X_n - X_n^\epsilon| > \delta n) = -\infty, P - a.s.,$$

from which it follows that X_n/n satisfies the quenched LDP with deterministic, convex, good rate function $I(\cdot)$.

3. Returning to the i.i.d. nestling setup, a natural question is whether one may prove an annealed large deviations principle for the position. A partial answer is given by the following. Fix a direction ℓ and recall the time $D = D(\ell)$ introduced in Section 3.2. Define $T_k^\ell = \min\{n : (X_n - X_0) \cdot \ell \geq k\}$. Then, for any $\lambda \in \mathbb{R}$,

$$\begin{aligned} \mathbb{E}^o(e^{\lambda T_{k+m}^\ell} \mathbf{1}_{\{D(\ell)=\infty\}}) &\geq E_P \left(E_\omega^o \left(e^{\lambda T_k^\ell} \mathbf{1}_{\{D(\ell) > T_k^\ell\}} \right) E_\omega^{X_{T_k^\ell}} \left(e^{\lambda T_m^\ell} \mathbf{1}_{\{D(\ell)=\infty\}} \right) \right) \\ &\geq \mathbb{E}^o(e^{\lambda T_k^\ell} \mathbf{1}_{\{D(\ell)=\infty\}}) \mathbb{E}^o(e^{\lambda T_m^\ell} \mathbf{1}_{\{D(\ell)=\infty\}}), \end{aligned}$$

and hence, by sub-additivity, the following limit exists:

$$\lim_{k \rightarrow \infty} \frac{1}{k} \log \mathbb{E}^o(e^{\lambda T_k^\ell} \mathbf{1}_{\{D(\ell)=\infty\}}) =: g(\ell, \lambda).$$

One can check that if the conclusions of Lemma 3.5.11 hold then also, for $-\lambda > 0$ small enough,

$$\limsup_{k \rightarrow \infty} \frac{1}{k} \log \mathbb{E}^o(e^{\lambda T_k^\ell}) = \lim_{k \rightarrow \infty} \frac{1}{k} \log \mathbb{E}^o(e^{\lambda T_k^\ell} \mathbf{1}_{\{D'=\infty\}}),$$

and hence for such λ ,

$$g(\ell, \lambda) = \limsup_{k \rightarrow \infty} \frac{1}{k} \log \mathbb{E}^o(e^{\lambda T_k^\ell}).$$

An interesting open question is to use this argument, in the nestling setup, to deduce a LDP and to relate the annealed and quenched rate functions.

Bibliographical notes: Large deviations for the position X_n of nestling RWRE in $\mathbb{Z}^d, d > 1$ were first derived in Zerner’s thesis [80]. Zerner uses a martingale differences argument instead of Lemma 3.4.14. With the same technique, he also derives a more general version of Lemma 3.4.10, under the name

“uniform shape theorem”. The large deviations for the hitting times $T_{[nz]}$ are implicit in his approach.

A recent paper of Varadhan [79] develops the quenched large deviations alluded to in remark 2 above, and a corresponding annealed LDP. He also obtains information on the zero set of the annealed and quenched rate functions, and in particular proves in a great generality that they coincide. The techniques are quite different from those presented here.

3.5 Kalikow’s condition

We introduce in this section a condition on the environment, due to Kalikow, which ensures that the RWRE is “ballistic”. Suppose P is elliptic, and let U be a strict subset of \mathbb{Z}^d , with $0 \in U$, and define on $U \cup \partial U$ an auxiliary Markov chain with transition probabilities

$$\hat{P}_U(x, x + e) = \begin{cases} \frac{\mathbb{E}^o[\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}} \omega(x, x+e)]}{\mathbb{E}^o[\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}}]}, & x \in U, |e| = 1 \\ 1 & x \in \partial U, e = 0 \end{cases} \tag{3.5.1}$$

where $\tau_{U^c} = \min\{n \geq 0 : X_n \in \partial U\}$ (note that the expectations in (3.5.1) are finite due the Markov property and ellipticity). The transition kernel \hat{P}_U weights the transitions $x \mapsto x + e$ according to the occupation time of the vertex x before exiting U . We denote by \hat{E}_U expectations with respect to the measure \hat{P}_U .

The following is a basic consequence of the definition of $\hat{P}_U(\cdot, \cdot)$:

Lemma 3.5.2 (Kalikow) *Assume $\hat{P}_U(\tau_{U^c} < \infty) = 1$. Then, $\hat{P}_U(X_{\tau_{U^c}} = v) = \mathbb{P}^o(X_{\tau_{U^c}} = v), v \in \partial U$. In particular, $\mathbb{P}^o(\tau_{U^c} < \infty) = 1$.*

Proof of Lemma 3.5.2:

Set $g_\omega(x) = E_\omega^o(\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}})$. Then

$$\hat{P}_U(x, y) = \frac{E(g_\omega(x)\omega(x, y))}{E(g_\omega(x))}, \quad x \in U, y \in U \cup \partial U. \tag{3.5.3}$$

But, due to the Markov property,

$$g_\omega(x) = \mathbf{1}_{\{x=0\}} + \sum_{z \in U} \omega(z, x)g_\omega(z),$$

and hence, using (3.5.3),

$$\sum_{x \in U} \left(E(g_\omega(x)) \right) \hat{P}_U(x, y) + \mathbf{1}_{\{y=0\}} = E(g_\omega(y)).$$

Set $\hat{\pi}_n(y) = \hat{E}_U(\sum_{j=0}^{\tau_{U^c} \wedge n} \mathbf{1}_{\{X_j=y\}})$. Then, $\hat{\pi}_0(y) = \mathbf{1}_{\{y=0\}}$ and

$$\hat{\pi}_{n+1}(y) = \mathbf{1}_{\{y=0\}} + \sum_{x \in U} \hat{P}_U(x, y) \hat{\pi}_n(x).$$

Then, for $y \in U \cup \partial U$,

$$E(g_\omega(y)) - \hat{\pi}_{n+1}(y) = \sum_{x \in U} \hat{P}_U(x, y) (E(g_\omega(x)) - \hat{\pi}_n(x)).$$

Since $E(g_\omega(y)) - \hat{\pi}_0(y) \geq 0$, it follows by the positivity of $\hat{P}_U(x, y)$ that for $y \in U \cup \partial U$,

$$\hat{E}_U \left(\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=y\}} \right) = \lim_{n \rightarrow \infty} \hat{\pi}_n(y) \leq E(g_\omega(y)).$$

Taking $y \in \partial U$ yields

$$\hat{P}_U(X_{\tau_{U^c}} = y) \leq \mathbb{P}^o(X_{\tau_{U^c}} = y), \quad y \in \partial U.$$

On the other hand, $\sum_{y \in \partial U} \hat{P}_U(X_{\tau_{U^c}} = y) = 1$ because $\hat{P}_U(\tau_{U^c} < \infty) = 1$ by assumption. Hence

$$\mathbb{P}^o(X_{\tau_{U^c}} = y) = \hat{P}_U(X_{\tau_{U^c}} = y), \quad \forall y \in \partial U. \quad \square$$

We are now ready to introduce Kalikow’s condition. Fix a hyperplane by picking a point $\ell \in \mathbb{R}^d \setminus \{0\}$, $|\ell|_1 \leq 1$. Define

$$\varepsilon_\ell := \inf_{U, x \in U} \sum_{|e|=1} (\ell \cdot e) \hat{P}_U(x, x + e)$$

where the infimum is over all connected strict subsets of \mathbb{Z}^d containing 0. We say that *Kalikow’s condition with respect to ℓ* holds if $\varepsilon_\ell > 0$. Note that ε_ℓ acts as a drift in the direction ℓ for the Markov chain \hat{P}_U .

A consequence of Lemma 3.5.2 is the following:

Theorem 3.5.4 *Assume that P satisfies Assumption 3.1.1. If Kalikow’s condition with respect to ℓ holds, then $\mathbb{P}^o(A_\ell) = 1$. If further P is an i.i.d. measure then $v_\ell > 0$.*

Proof. Fix $U_L = \{z \in \mathbb{Z}^d : |z \cdot \ell| \leq L\}$. Let $\hat{X}_{n,L}$ denote the Markov chain with $\hat{X}_{0,L} = 0$ and transition law $\hat{P}_{U_L}(x, x + e)$. Set the local drift at x , $\hat{d}(x) = \sum_{|e|=1} e \hat{P}_{U_L}(x, x + e)$, and recall that $\hat{X}_{n,L} - \sum_{i=0}^{n-1} \hat{d}(\hat{X}_{i,L})$ is a martingale, with bounded increments. It follows that for some constant C ,

$$\hat{P}_{U_L} \left(\sup_{0 \leq n \leq N} |\hat{X}_{n,L} - \sum_{i=0}^{n-1} \hat{d}(\hat{X}_{i,L})| > \delta N \right) \leq C e^{-C \delta^2 N}. \quad (3.5.5)$$

On the other hand, $\sum_{i=0}^{n-1} \hat{d}(\hat{X}_{i,L}) \cdot \ell \geq \varepsilon_\ell (n \wedge \tau_{U_L^c})$ while $|\hat{X}_{n,L} \cdot \ell| \leq L + 1$. We thus conclude from (3.5.5) that

$$\hat{P}_{U_L} \left(\tau_{U_{\ell}} > \frac{L+1}{\varepsilon_{\ell}} + \frac{\delta}{\varepsilon_{\ell}} N \right) \leq C e^{-C\delta^2 N / \varepsilon_{\ell}^2},$$

and hence, for some C_1 independent of L , and all L large,

$$\hat{P}_{U_L} \left(|\hat{X}_{\tau_{U_{\ell}}} \cdot \ell - L| > 1 \right) \leq C_1 e^{-C_1 L}.$$

It follows from Lemma 3.5.2 that

$$\mathbb{P}^o \left(|X_{\tau_{U_{\ell}}} \cdot \ell - L| > 1 \right) \leq C_1 e^{-C_1 L}.$$

A similar argument shows that $\mathbb{P}^o(D = \infty) > 0$: indeed, take now $U_{L,+} = \{z \in \mathbb{Z}^d : 0 \leq z \cdot \ell \leq L\}$. Arguing as above, one finds that for some $C_2 > 0$ independent of L ,

$$\hat{P}_{U_{L,+}} \left(|\hat{X}_{\tau_{U_{L,+}}} - L| \leq 1 \right) > C_2,$$

implying that

$$\mathbb{P}^o \left(|X_{\tau_{U_{L,+}}} - L| \leq 1 \right) > C_2.$$

Thus, $\mathbb{P}^o(D = \infty) > 0$, and then, by an argument as in the proof of Theorem 3.1.2, $\mathbb{P}^o(A_{\ell}) > 0$. By Theorem 3.1.2, it follows that $\mathbb{P}^o(A_{\ell} \cup A_{-\ell}) = 1$. Due to (3.5.5), it holds that $\mathbb{P}^o(\limsup_{n \rightarrow \infty} X_n \cdot \ell = \infty) = 1$. We thus conclude that $\mathbb{P}^o(A_{\ell}) = 1$.

To see that if P is i.i.d. then $v_{\ell} > 0$, recall the regeneration times $\{\tau_i\}_{i \geq 1}$ introduced in Section 3.2. By Lemma 3.2.5, it suffices to prove that $\mathbb{E}^o(\tau_1 | D = \infty) < \infty$. Let $U_{m,k,-} = \{z \in \mathbb{Z}^d : |z| < k, z \cdot \ell < m\}$, and set $T_{m,k} := T_{U_{m,k,-}}$ with $T_m = \lim_{k \rightarrow \infty} T_{m,k} = \min\{n : X_n \cdot \ell \geq m\}$, $m \geq 1$. By Kalikow's condition,

$$\mathbb{E}^o \left(\sum_{n=0}^{T_{m,k}} \mathbf{1}_{\{X_n=x\}} \sum_{e:|e|=1} \omega(x, x+e) \ell \cdot e \right) \geq \varepsilon_{\ell} \mathbb{E}^o \left(\sum_{n=0}^{T_{m,k}} \mathbf{1}_{\{X_n=x\}} \right),$$

and hence, summing over $x \in U_{m,k,-}$ and recalling that $X_{i+1} - X_i - d(\theta^{X_i} \omega)$ is a martingale difference sequence, one gets

$$1 + m \geq \mathbb{E}^o(X_{T_{m,k}} \cdot \ell) \geq \varepsilon_{\ell} \mathbb{E}^o(T_{m,k}),$$

and taking $k \rightarrow \infty$ one concludes that $m + 1 \geq \varepsilon_{\ell} \mathbb{E}^o(T_m)$. In particular,

$$\mathbb{E}^o(\liminf_{m \rightarrow \infty} T_m/m) \leq \liminf_{m \rightarrow \infty} \mathbb{E}^o(T_m/m) \leq 1/\varepsilon_{\ell}. \tag{3.5.6}$$

Since $\tau_i \rightarrow_{i \rightarrow \infty} \infty$, \mathbb{P}^o -a.s., one may find a (random) sequence k_m such that $\tau_{k_m} \leq T_m < \tau_{k_m+1}$. By definition,

$$\ell \cdot X_{\tau_{k_m}} \leq \ell \cdot X_{T_m} \leq \ell \cdot X_{\tau_{k_m+1}},$$

and further,

$$\ell \cdot X_{\tau_k} / k \xrightarrow{k \rightarrow \infty} \mathbb{E}^o(\ell \cdot X_{\tau_1} | \{D = \infty\}) = \frac{1}{\mathbb{P}^o(D = \infty)} < \infty, \mathbb{P}^o - a.s.,$$

due to Lemma 3.2.5 and (3.2.7). Thus, $k_m/m \rightarrow_{m \rightarrow \infty} \mathbb{E}^o(X_{\tau_1} \cdot \ell | \{D = \infty\})^{-1}$, \mathbb{P}^o -a.s. But, since $\tau_{k_m}/k_m \rightarrow_{m \rightarrow \infty} \mathbb{E}^o(\tau_1 | \{D = \infty\}) \in [1, \infty]$, \mathbb{P}^o -a.s., it follows that

$$\begin{aligned} \liminf_{m \rightarrow \infty} \frac{T_m}{m} &\geq \liminf_{m \rightarrow \infty} \frac{\tau_{k_m} k_m}{k_m m} = \lim_{m \rightarrow \infty} \frac{\tau_{k_m} k_m}{k_m m} = \frac{\mathbb{E}^o(\tau_1 | \{D = \infty\})}{\mathbb{E}^o(X_{\tau_1} \cdot \ell | \{D = \infty\})} \\ &= \mathbb{E}^o(\tau_1 | \{D = \infty\}) \mathbb{P}^o(D = \infty). \end{aligned}$$

Since (3.5.6) implies that $\liminf_{m \rightarrow \infty} T_m/m < \infty$, we conclude that $\mathbb{E}^o(\tau_1 | \{D = \infty\}) < \infty$, and hence $v_\ell > 0$. \square

By noting that if Kalikow’s condition holds for some ℓ_0 then it holds for all ℓ in a neighborhood of ℓ_0 , one gets immediately the

Corollary 3.5.7 *Assume that P satisfies assumption 3.2.1. If Kalikow’s condition with respect to some ℓ holds, then there exists a deterministic v such that*

$$\frac{X_n}{n} \xrightarrow{n \rightarrow \infty} v, \quad P - a.s.$$

The following is a sufficient condition for Kalikow’s condition to hold true:

Lemma 3.5.8 *Assume P is i.i.d. and elliptic. Then Kalikow’s condition with respect to ℓ holds if*

$$\inf_{f \in \mathcal{F}} \frac{E \left(\frac{\sum_{e:|e|=1} \omega(0,e)\ell \cdot e}{\sum_{e:|e|=1} \omega(0,e)f(e)} \right)}{E \left(\frac{1}{\sum_{e:|e|=1} \omega(0,e)f(e)} \right)} > 0, \tag{3.5.9}$$

where \mathcal{F} denotes the collection of nonzero functions on $\{e : |e| = 1\}$ taking values in $[0, 1]$.

Proof. Fix U a strict subset of \mathbb{Z}^d , $x \in U$, and let $\tau_x = \min\{n \geq 0 : X_n = x\}$. Define $g(x, y, \omega) := E_\omega^y(\mathbf{1}_{\{\tau_x < \tau_{U^c}\}})$. Note that $g(x, y, \omega)$ is independent of ω_x . Next,

$$\begin{aligned} &\sum_{|e|=1} (\ell \cdot e) \hat{P}_U(x, x + e) \\ &= \frac{E \left(E_\omega^o \left(\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}} \sum_{e:|e|=1} \omega(x, x + e) e \cdot \ell \right) \right)}{E \left(E_\omega^o \left(\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}} \right) \right)} \\ &= \frac{E \left(g(x, 0, \omega) E_\omega^x \left(\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}} \sum_{e:|e|=1} \omega(x, x + e) e \cdot \ell \right) \right)}{E \left(g(x, 0, \omega) E_\omega^x \left(\sum_{n=0}^{\tau_{U^c}} \mathbf{1}_{\{X_n=x\}} \right) \right)}. \end{aligned}$$

Under P_ω^x , the process X_n is a Markov chain with Geometric($\sum_{e:|e|=1} \omega(x, x+e)g(x, x+e, \omega)$) number of visits at x . The last equality and the Markov property then imply

$$\begin{aligned} & \sum_{|e|=1} (\ell \cdot e) \hat{P}_U(x, x+e) \\ &= \frac{E \left(g(x, 0, \omega) \frac{\sum_{e:|e|=1} \omega(x, x+e) e \cdot \ell}{\sum_{e:|e|=1} \omega(x, x+e) g(x, x+e, \omega)} \right)}{E \left(g(x, 0, \omega) \frac{1}{\sum_{e:|e|=1} \omega(x, x+e) g(x, x+e, \omega)} \right)} \\ &= \frac{E \left(\frac{\sum_{e:|e|=1} \omega(x, e) \ell \cdot e}{\sum_{e:|e|=1} \omega(x, e) g(x, x+e, \omega) / g(x, 0, \omega)} \right)}{E \left(\frac{1}{\sum_{e:|e|=1} \omega(x, e) g(x, x+e, \omega) / g(x, 0, \omega)} \right)} \\ &\geq \inf_{f \in \mathcal{F}} \frac{E \left(\frac{\sum_{e:|e|=1} \omega(x, e) \ell \cdot e}{\sum_{e:|e|=1} \omega(x, e) f(e)} \right)}{E \left(\frac{1}{\sum_{e:|e|=1} \omega(x, e) f(e)} \right)} = \inf_{f \in \mathcal{F}} \frac{E \left(\frac{\sum_{e:|e|=1} \omega(0, e) \ell \cdot e}{\sum_{e:|e|=1} \omega(0, e) f(e)} \right)}{E \left(\frac{1}{\sum_{e:|e|=1} \omega(0, e) f(e)} \right)} > 0, \end{aligned}$$

where the first inequality is due to the independence of $g(x, x+e, \omega)/g(x, 0, \omega)$ in ω_x . □

An easy corollary is the following

Corollary 3.5.10 *Assume P satisfies Assumption 3.2.1. If either*
 (a) $\text{supp}(P_d) \subset \{z \in \mathbb{R}^d : \ell \cdot z \geq 0\}$ but $\text{supp}(P_d) \not\subset \{z \in \mathbb{R}^d : \ell \cdot z = 0\}$,
 or
 (b) $E((\sum_{e:|e|=1} \omega(0, e) e \cdot \ell)_+) > \frac{1}{\varepsilon} E((\sum_{e:|e|=1} \omega(0, e) e \cdot \ell)_-)$,
 then Kalikow’s condition with respect to ℓ holds.

In particular, non-nestling walks or walks with drift “either neutral or pointing to the right” satisfy Kalikow’s condition with respect to an appropriate hyperplane ℓ . Further there exist truly nestling walks which do satisfy Kalikow’s condition.

Remark: It is interesting to note that when P is elliptic and i.i.d. and $d = 1$, Kalikow’s condition is equivalent to $v \neq 0$, i.e. to the walk being “ballistic”. For $d > 1$, it is not clear yet whether there exist walks with $\mathbb{P}^o(A_\ell) > 0$ but with zero speed. Such walks, if they exist, necessarily cannot satisfy Kalikow’s condition.

Our next goal is to provide tail estimates on $X_{\tau_1} \cdot \ell$ and on τ_1 . Our emphasis here is in providing a (relatively) simple proof and not the sharpest possible result. For the latter we refer to [71]. In this spirit we throughout assume $\ell = e_1$.

Lemma 3.5.11 *Assume P is i.i.d. and satisfies Kalikow’s condition. Then, there exists a constant c such that*

$$\mathbb{E}^o(\exp c X_{\tau_1} \cdot \ell) < \infty.$$

Proof. Recall the notations of Section 3.2 and write

$$\begin{aligned}
 \mathbb{E}^o(\exp(c X_{\tau_1} \cdot \ell)) &= \sum_{k \geq 1} \mathbb{E}^o\left(\exp(c X_{\tau_1} \cdot \ell) \mathbf{1}_{\{K=k\}}\right) \\
 &= \sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} e^{c x \cdot \ell} E\left(E_\omega^o\left(\mathbf{1}_{\{X_{\overline{S}_k}=x\}} \mathbf{1}_{\{\overline{S}_k < \infty\}}\right) \cdot P_\omega^x(D = \infty)\right) \\
 &= \mathbb{P}^o(D = \infty) E\left(\sum_{k \geq 1} \sum_{x \in \mathbb{Z}^d} e^{c x \cdot \ell} E_\omega^o\left(\mathbf{1}_{\{X_{\overline{S}_k}=x\}} \mathbf{1}_{\{\overline{S}_k < \infty\}}\right)\right) \\
 &= \mathbb{P}^o(D = \infty) \sum_{k \geq 1} \mathbb{E}^o\left(\exp(c X_{\overline{S}_k} \cdot \ell) \mathbf{1}_{\{\overline{S}_k < \infty\}}\right). \tag{3.5.12}
 \end{aligned}$$

But, using the Markov property,

$$\begin{aligned}
 &\mathbb{E}^o\left(\exp(c X_{\overline{S}_k} \cdot \ell) \mathbf{1}_{\{\overline{S}_k < \infty\}}\right) \\
 &\leq \mathbb{E}^o\left(\exp(c X_{\overline{S}_{k-1}} \cdot \ell) \mathbf{1}_{\{\overline{S}_{k-1} < \infty\}}\right) \mathbb{E}^o\left(\exp(c M_0 \cdot \ell) \mathbf{1}_{\{D < \infty\}}\right).
 \end{aligned}$$

Hence,

$$\mathbb{E}^o\left(\exp(c X_{\overline{S}_k} \cdot \ell) \mathbf{1}_{\{\overline{S}_k < \infty\}}\right) \leq \mathbb{P}^o(D = \infty) \sum_{k \geq 0} \left(\mathbb{E}^o(\exp(c M_0 \cdot \ell) \mathbf{1}_{\{D < \infty\}})\right)^k.$$

Note, using $\ell = e_1$, that

$$\begin{aligned}
 &\mathbb{E}^o\left(e^{c M_1 \cdot \ell} \mathbf{1}_{\{D < \infty\}}\right) \\
 &= \sum_{k=1}^\infty e^{ck} \mathbb{E}^o\left(\mathbf{1}_{\{M_0 \cdot \ell = k\}} \mathbf{1}_{\{D < \infty\}}\right) \\
 &= \sum_{k=1}^\infty e^{ck} \sum_{y \in \mathbb{Z}^{d-1}} \mathbb{E}^o\left(E_\omega^o\left(\mathbf{1}_{\{X_{T_k}=(k,y)\}} E_\omega^{(k,y)}(T_0 < T_{k+1})\right)\right). \tag{3.5.13}
 \end{aligned}$$

Using Kalikow’s condition and a computation as in Theorem 3.5.4, one has that for some $c_1 > 0$, $\mathbb{P}^o\left(\sum_{|y| > \frac{2k}{\varepsilon \ell}} \mathbf{1}_{\{X_{T_k}=(k,y)\}}\right) \leq e^{-c_1 k}$, while $\mathbb{P}^o(T_{-k} < T_1) \leq e^{-c_1 k}$. Hence, substituting in (3.5.13), one has

$$\mathbb{E}^o\left(e^{c M_1 \cdot \ell} \mathbf{1}_{\{D < \infty\}}\right) \leq \sum_{k=1}^\infty e^{ck} \left(\left(\frac{4k}{\varepsilon \ell}\right)^{d-1} + 1\right) e^{-c_1 k}.$$

Taking $c < c_1$, the lemma follows. □

A direct consequence of Lemma 3.5.11 is that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}^o(X_{\tau_1} \cdot \ell > vn) \leq -\beta(v) \tag{3.5.14}$$

where $\beta(v) > 0$ for $v > 0$.

With a proof very similar to that of Lemma 3.5.11, we have in fact the

Lemma 3.5.15 *Assume P is i.i.d. and satisfies Kalikow’s condition. Set $X^* = \sup_{0 \leq n \leq \tau_1} |X_n|$. Then, there exists a constant c' such that*

$$\mathbb{E}^o(\exp c' X^*) < \infty.$$

We next turn to obtaining tail estimates on τ_1 . Here, due to the presence of “traps”, one cannot in general expect exponential decay as in (3.5.14). We aim at proving the following result.

Theorem 3.5.16 *Assume P satisfies Assumption 3.2.1. and Kalikow’s condition. Then, with $d \geq 2$, there exists an $\alpha > 1$ such that for all u large,*

$$\mathbb{P}^o(\tau_1 > u) \leq e^{-(\log u)^\alpha}.$$

In particular, τ_1 possesses all moments.

Proof. Recall that we take $\ell = e_1$ and, for $L > 0$, set $c_L = (-L, L) \times \left(-\frac{2L}{\varepsilon_\ell}, \frac{2L}{\varepsilon_\ell}\right)^{d-1}$. Note that, with c as in Lemma 3.5.11,

$$\begin{aligned} \mathbb{P}^o(\tau_1 \geq u) &\leq \mathbb{P}^o\left(\tau_1 \geq u, X_{\tau_1} \cdot \ell \leq L\right) + \mathbb{P}^o\left(X_{\tau_1} \cdot \ell \geq L\right) \\ &\leq e^{-cL/2} + \mathbb{P}^o(\tau_L > \tau_{c_L}) + \mathbb{P}^o(\tau_{c_L} = \tau_L \geq u) \end{aligned}$$

where

$$\begin{aligned} \tau_L &= \inf\{t : X_t \cdot \ell \geq L\} \text{ and} \\ \tau_{c_L} &= \inf\{t : X_t \notin c_L\}. \end{aligned}$$

Hence, by Kalikow’s condition,

$$\mathbb{P}^o(\tau_1 \geq u) \leq e^{-\bar{c}L/2} + \mathbb{P}^o\left(\tau_{c_L} = \tau_L \geq u\right) \tag{3.5.17}$$

for some constant $\bar{c} > 0$.

The heart of the proof of Theorem 3.5.16 lies in the following lemma, whose proof is deferred.

Lemma 3.5.18 *There exist a $\beta < 1$ and $\xi > 1$ such that for any $c > 0$,*

$$\limsup_{L \rightarrow \infty} \frac{1}{L^\xi} \log \mathbb{P}\left(P_\omega^o(X_{\tau_{U_L}} \cdot \ell \geq L) \leq e^{-cL^\beta}\right) < 0,$$

where $U_L = \{z \in \mathbb{Z}^d : |z \cdot \ell| \leq L\}$.

Accepting Lemma 3.5.18, let us complete the proof of Theorem 3.5.16. Toward this end, set $\Delta(u) = a \log u$ with a small enough such that $\varepsilon^{\Delta(u)} > \frac{1}{u^{1/6}}$, and set $L = L(u) = (\log u)^\alpha$, with $N = L/\Delta = \frac{1}{a}(\log u)^{\alpha-1}$ and $2 - \alpha = \beta$ where β is as in Lemma 3.5.18. Observe that

$$\begin{aligned} \mathbb{P}^o(\tau_{c_L} \geq u) &\leq \left(\frac{1}{2}\right)^{(\log u)^{\bar{\alpha}}} + P^o\left(\exists x_1 \in c_L, P_\omega^{x_1}\left(\tau_{c_L} > \frac{u}{(\log u)^{\bar{\alpha}}}\right) \geq \frac{1}{2}\right) \\ &=: \left(\frac{1}{2}\right)^{(\log u)^{\bar{\alpha}}} + P(\mathcal{R}). \end{aligned} \tag{3.5.19}$$

Note that

$$\begin{aligned} &\frac{u}{(\log u)^\alpha} \cdot P_\omega^{x_1}\left(\tau_{c_L} > \frac{u}{(\log u)^{\bar{\alpha}}}\right) \\ &\leq E_\omega^{x_1}(\tau_{c_L}) = E_\omega^{x_1}\left(\sum_{i=1}^{\tau_{c_L}} 1\right) = E_\omega^x\left(\sum_{y \in c_L} \sum_{i=1}^{\tau_{c_L}} \mathbf{1}_{\{X_n=y\}}\right) \\ &= \sum_{y \in c_L} \frac{E_\omega^x(\mathbf{1}_{\{\tau_y < \tau_{c_L}\}})}{E_\omega^y(\mathbf{1}_{\{\tau_y > \tau_{c_L}\}})} \leq |c_L| \frac{1}{\inf_{y \in c_L} P_\omega^y(\mathbf{1}_{\{\tau_y\}} > \tau_{c_L})}, \end{aligned}$$

with $\tau_y = \inf\{t : X_t = y\}$. Hence,

$$\inf_{y \in c_L} P_\omega^y(\tau_y > \tau_{c_L}) \leq \frac{|c_L|(\log u)^{\bar{\alpha}}}{u \inf_{x_1 \in c_L} P_\omega^{x_1}\left(\tau_{c_L} > \frac{u}{(\log u)^\alpha}\right)}.$$

Hence, on \mathcal{R} there exists a y with $\mathbb{P}_\omega^y(\tau_y > \tau_{c_L}) \leq \frac{1}{u^{4/5}}$, for all u large enough.

Set $A_i = \{z \in \mathbb{Z}^d : z \cdot \ell = i\Delta\}$. By ellipticity (recall $\varepsilon^{\Delta(u)} > 1/u^{1/6}$), it follows that on the event \mathcal{R} , there exists an $i_0 \in [-N + 2, N - 1]$ and an $x \in A_{i_0}$ such that

$$P_\omega^{x_0}\left(\tau_{(i_0-1)\Delta} > \tau_{c_L}\right) \leq \frac{1}{\sqrt{u}} \tag{3.5.20}$$

where $\tau_{(i_0-1)\Delta} = \inf\{t : X_t \cdot \ell = (i_0 - 1)\Delta\}$. Set

$$X_i = -\log \min_{z \in A_i \cap c_L} P_\omega^z\left(\tau_{(i-1)\Delta} > \tau_{(i+1)\Delta}\right).$$

Then, the Markov property implies that

$$\begin{aligned} P_\omega^x\left(\tau_{(i-1)\Delta} > \tau_{c_L}\right) &\geq \exp\left(-\sum_{j=i}^{N-1} X_j\right) \varepsilon^{\Delta(u)} \\ &\geq \exp\left(-\sum_{j=i}^{N-1} X_j\right) \frac{1}{u^{1/6}}. \end{aligned}$$

Thus, using (3.5.20)

$$P(\mathcal{R}) \leq P\left(\sum_{i=-N+1}^{N-1} X_i \geq \frac{\log u}{12}\right) \leq 2N \sup_{-N+1 \leq i \leq N-1} P\left(X_i \geq \frac{\log u}{24N}\right). \tag{3.5.21}$$

But, using the shift invariance of P and the definition of $\{X_i\}$,

$$2N \sup_{N+1 \leq i \leq N-1} P \left(X_i \geq \frac{\log u}{24N} \right) \leq |c_L| P \left(P_\omega^o(\tau_{-\Delta} > \tau_\Delta) \leq e^{-(\log u)^{2-\bar{\sigma}b}} \right)$$

for $b = \frac{a}{24}$. The theorem follows by an application of Lemma 3.5.18. \square

Proof of Lemma 3.5.18

The interesting aspect in proving the Lemma is the fact that one constructs lower bounds on $P_\omega^o(X_{\tau_{u_L}} \cdot \ell \geq L)$ for many configurations. Toward this end, fix $1 > \beta > \bar{\beta}$, $\gamma \in (\frac{1}{2}, 1)$, $\chi = \frac{1-\bar{\beta}}{1-\gamma} < \bar{\beta} < 1$ such that $d(\bar{\beta} - \chi) > 1$ (for $\bar{\beta}$ close enough to 1, one may always find a γ close to 1 such that this condition is satisfied, if $d \geq 2$). Set next $L_0 = L^\chi$, $L_1 = L^{\bar{\beta}}$, $N_0 = L_1/L_0$ (for simplicity, assume that L_0, L_1, L_0^γ and N_0 are all integer).

Let R be a rotation of \mathbb{Z}^d such that $\tilde{R}(\ell) = \tilde{R}(e_1) = \frac{v}{|v|}$, and define

$$B_1(z) = \tilde{R} \left(z + [0, L_0]^d \right) \cap \mathbb{Z}^d$$

$$B_2(z) = \tilde{R} \left(z + [-L_0^\gamma, L_0 + L_0^\gamma]^d \right) \cap \mathbb{Z}^d$$

and $\partial_+ B_2(z) = \partial B_2(z) \cap \left\{ x : x \cdot \frac{v}{|v|} \geq L_0 + L_0^\gamma \right\}$.

We say that $z \in L_0 \mathbb{Z}^d$ is *good* if $\sup_{x \in B_1(z)} P_\omega^x(X_{\tau_{B_2(z)}} \notin \partial_+ B_2(z)) \leq \frac{1}{2}$ and say that it is *bad* otherwise. The following estimate is a direct consequence of Kalikow’s condition and Lemma 3.5.15.

Lemma 3.5.22 For $\gamma \in (1/2, 1)$,

$$\limsup_{L_0 \rightarrow \infty} L_0^{1-2\gamma} \log P(0 \text{ is bad}) < 0.$$

Proof of Lemma 3.5.22

Set $u = \max\{\bar{u} : \sup_{x \in B_2(0)} x \cdot \ell \geq \bar{u}\}$ and set $L_u = \sup\{n \geq 0 : X_n \cdot \ell \leq u\}$. Define $\pi(z) = z - \frac{z \cdot v}{|v|^2} v$. Setting $K_n = \sup\{k \geq 0 : \tau_k < n\}$, it holds that $n \leq L_u \Rightarrow K_n \leq u$. Setting $w \in \mathbb{R}^d$ with $w \cdot v = 0$, and $|w|_1 = 1$ one has

$$X_n \cdot w = X_{\tau_{K_n}} \cdot w + (X_n - X_{\tau_{K_n}}) \cdot w$$

$$\leq X_{\tau_{K_n}} \cdot w + X^* \circ \theta_{\tau_{K_n}} \cdot w$$

Hence,

$$\begin{aligned} \mathbb{P}^o\left(\sup_{0 \leq n \leq L_u} X_n \cdot w \geq u^\gamma\right) &\leq \sum_{0 \leq k \leq u} \mathbb{P}^o\left(X_{\tau_k} \cdot w + X^* \circ \theta_{\tau_k} > u^\gamma\right) \\ &\leq \sum_{0 \leq k \leq u} \mathbb{P}^o\left(\left(X_{\tau_k} - X_{\tau_1}\right) \cdot w > \frac{u^\gamma}{3}\right) + u \mathbb{P}^o\left(X_{\tau_1} \cdot w > \frac{u^\gamma}{3}\right) \\ &\quad + u \mathbb{P}^o\left(X^* > \frac{u^\gamma}{3} \mid D = \infty\right) \\ &\leq \sum_{0 \leq k \leq u} \mathbb{P}^o\left(\left(X_{\tau_k} - X_{\tau_1}\right) \cdot w > \frac{u^\gamma}{3}\right) + \frac{2u}{\mathbb{P}^o(D = \infty)} \mathbb{P}^o\left(X^* > \frac{u^\gamma}{3}\right). \end{aligned}$$

Note that by Lemma 3.5.15, $\mathbb{P}^o\left(X^* > \frac{u^\gamma}{3}\right) \leq e^{-c_0 u^\gamma}$ while the random variables $(X_{\tau_{i+1}} - X_{\tau_i}) \cdot w$ are i.i.d., of zero mean and finite exponential moments. In particular,

$$\mathbb{P}^o\left(\left|\frac{(X_{\tau_k} - X_{\tau_1}) \cdot w}{k}\right| > \frac{u^\gamma}{3k}\right) \leq e^{-c_0 k \frac{u^{2\gamma}}{9k^2}} \leq e^{-c_0 u^{2\gamma-1}}$$

by moderate deviations (see e.g., [19, Section 3.7]).

Since $\gamma > 2\gamma - 1$, we conclude that

$$\limsup_{u \rightarrow \infty} \frac{1}{u^{2\gamma-1}} \log \mathbb{P}^o\left(\sup_{0 \leq n \leq L_u} X_n \cdot w \geq u^\gamma\right) < 0$$

and hence

$$\limsup_{u \rightarrow \infty} \frac{1}{u^{2\gamma-1}} \log \mathbb{P}^o\left(\sup_{0 \leq n \leq L_u} |\pi(X_n)| \geq u^\gamma\right) < 0. \tag{3.5.23}$$

Fix now $x \in B_1(0)$. Then, for some $c = c(d)$,

$$\begin{aligned} \mathbb{P}^x\left(X_{\tau_{B_2(0)}} \notin \partial_+ B_2(0)\right) &\leq \mathbb{P}^o\left(\sup_{0 \leq n \leq L_{cu}} \pi(X_n \cdot w) \geq u^\gamma\right) \\ &\quad + \mathbb{P}^o\left(X_{\tau_{V_u}} \cdot \ell < 0\right) \end{aligned}$$

where $V_u = \left\{z : \frac{-L_0^\gamma}{c} \leq z \cdot \ell \leq cL_0^\gamma\right\}$ and the conclusion follows from (3.5.23) and Kalikow's condition. □

Construct now the following subsets of U_L :

$$\text{Set } M = \left\{z \in L_0 \mathbb{Z}^d, z = (0, \bar{z}), \bar{z} \in \left\{-\frac{L_1}{L_0}, \dots, 0, \frac{L_1}{L_0}\right\}^{d-1} L_0\right\}.$$

For $z \in M$,

set $\text{Row}(z) = \cup_{j=N_-(z)}^{N_+(z)} B_1(z + jL_0\ell)$ where

$$N_-(z) = \min\{j : B_1(z + jL_0\ell) \cap \{x : x \cdot \ell \geq 0\} \neq \emptyset\} \geq -c \frac{L}{L_0}$$

$$N_+(z) = \max\{j : B_1(z + jL_0\ell) \cap \{x : x \cdot \ell \geq 0\} \neq \emptyset\} \leq c \frac{L}{L_0}$$

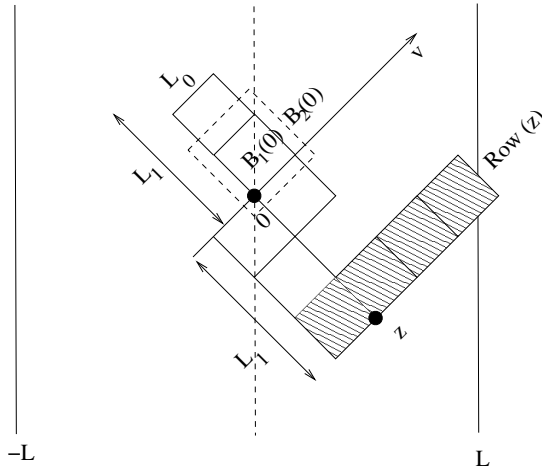


Fig. 3.5.1.

for some constant c depending on v and uniformly bounded, and set $T = \cup_{j \in \{-\frac{L_1}{L_0}, \dots, 0, \dots, \frac{L_1}{L_0}\}^{d-1}} \{jL_0\}$.

The idea behind the proof is that if one of the rows $\{R(z)\}_{z \in M}$, say $R(z_0)$, contains mostly good blocks, a good strategy for the event $(X_{\tau_{u_L}} \cdot \ell) \geq L$ is to force the walker started at x to first move to z , then move to the right successively without leaving $\cup_{z \in \text{Row}(z_0)} B_2(z)$ until exiting from u_L . More precisely, let $N(z_0)$ denote the number of bad blocks in $\cup_{z \in \text{Row}(z_0)} B_1(z)$. Then, for some constants c_i , using ellipticity and the definition of good boxes,

$$\begin{aligned}
 P_\omega^o(X_{\tau_{u_L}} \cdot \ell \geq L) &\geq \varepsilon^{L_1} \left(\frac{1}{2} \varepsilon^{2(d-1)} L_0^\gamma \right)^{cL/L_0} (\varepsilon^{L_0})^{N(z_0)} \\
 &= e^{-c_4(L^{\bar{\beta}} + L^\chi N(z_0))}.
 \end{aligned}$$

Hence, for an arbitrary constant c_6 and all L large enough ($L > g(c_6, c)$ for some fixed function $g(\cdot)$),

$$\begin{aligned}
 P \left(P_\omega^o(X_{\tau_{u_L}} \cdot \ell \geq L) \leq e^{-cL^\beta} \right) &\leq P \left(\{ \bar{A}z_0 \in M : N(z_0) \leq c_6 L^{\bar{\beta}-\chi} \} \right) \\
 &\leq \left[P \left(N(0) \geq c_6 L^{\bar{\beta}-\chi} \right) \right]^{\left(\frac{cL_1}{L_0} \right)^{(d-1)}}
 \end{aligned}$$

using the independence between even rows in Figure 3.5.1. But note that $N(0) = \sum_{i=1}^{L/2L_0} (Y_i + Z_i)$ where $\{Y_i\}$ are i.i.d., $\{Z_i\}$ are i.i.d., $\{0, 1\}$ valued, $P(Y_i = 1) = P(B_1(0))$ is a bad block (the division to (Y_i, Z_i) reflects the division to even and odd blocks, which creates independence). Hence,

$$P\left(N(0) \geq c_6 L^{\bar{\beta}-x}\right) \leq 2P\left(\sum_{i=1}^{L/2L_0} Y_i \geq \frac{c_6}{2} L^{\bar{\beta}-x}\right).$$

Since, by Lemma 3.5.22,

$$\log E(e^{Y_i}) \leq \log\left(1 + e^{-c_7 L_0^{2\gamma-1}}\right) \leq e^{-c_7 L_0^{2\gamma-1}},$$

we conclude from the independence of the Y_i 's that

$$\begin{aligned} P\left(N(0) \geq c_6 L^{\bar{\beta}-x}\right) &\leq 2e^{-\frac{c_6}{2} L^{\bar{\beta}-x}} e^{e^{-c_7 L_0^{2\gamma-1}} \cdot \frac{L}{2L_0}} \\ &\leq e^{-c_8 L^{\bar{\beta}-x}}. \end{aligned}$$

Hence,

$$P\left(P_\omega^\circ\left(X_{\tau_{u_L} \cdot \ell} \geq L\right) \leq e^{-cL^\beta}\right) \leq e^{-c_9 L^{(\bar{\beta}-x)^d}} \leq e^{-L^\xi}$$

for some $\xi > 1$, as claimed. □

Remark The restriction to $d \geq 2$ in Theorem 3.5.16 is essential: as we have seen in the case $d = 1$, one may have ballistic walks (and hence, in $d = 1$, satisfying Kalikow's condition) with moments $m_r := \mathbb{E}^\circ(\tau_1^r)$ of the regeneration time τ_1 being finite only for small enough $r > 1$.

We conclude this section by showing that estimates of the form of Theorem 3.5.16 lead immediately to a CLT. The statement is slightly more general than needed, and does not assume Kalikow's condition but rather some of its consequences.

Theorem 3.5.24 *Assume Assumption 3.2.1, and further assume that $\mathbb{P}^\circ(A_\ell) = 1$ and that the regeneration time τ_1 satisfies $\mathbb{E}^\circ(\tau_1^{2+\delta}) < \infty$ for some $\delta > 0$. Then, under the annealed measure \mathbb{P}° ,*

$$X_n/n \xrightarrow{n \rightarrow \infty} v := \frac{\mathbb{E}^\circ(X_{\tau_2} - X_{\tau_1})}{\mathbb{E}^\circ(\tau_2 - \tau_1)} \neq 0, \quad \mathbb{P}^\circ - a.s., \tag{3.5.25}$$

and $(X_n - nv)/\sqrt{n}$ converges in law to a centered Gaussian vector.

Proof. The LLN (3.5.25) is a consequence of Theorem 3.2.2 and its proof. To see the CLT, set

$$\xi_i = X_{\tau_{i+1}} - X_{\tau_i} - (\tau_{i+1} - \tau_i)v, \quad S_n := \sum_{i=1}^n \xi_i,$$

and $\Xi = \mathbb{E}^\circ(\xi_1 \xi_1^T)$. It is not hard to check that Ξ is non-degenerate, simply because $\mathbb{P}^\circ(|\xi_1| > K) > 0$ for each $K > 0$. Then S_n is under \mathbb{P}° a sum of i.i.d. random variables possessing finite $2 + \delta$ -th moments, and thus $S_{[nt]}/\sqrt{n}$ satisfies the invariance principle, with covariance matrix Ξ . Define

$$\nu_n = \min \left\{ j : \sum_{i=1}^j (\tau_{i+1} - \tau_i) > n \right\}.$$

Note that in \mathbb{P}^o probability, $n/\nu_n \rightarrow \mathbb{E}^o(\tau_2 - \tau_1) < \infty$. Hence, by time changing the invariance principle, see e.g. [2, Theorem 14.4], $S_{\nu_n}/\sqrt{\nu_n}$ converges in \mathbb{P}^o probability to a centered Gaussian variable of covariance Ξ . On the other hand, for any positive η ,

$$\begin{aligned} \mathbb{P}^o(|S_{\nu_n} - (X_n - nv)| > \eta\sqrt{n}) &\leq \mathbb{P}^o(\exists i \leq n : (\tau_{i+1} - \tau_i) > \eta\sqrt{n}/2) \\ &\quad + \mathbb{P}^o(\tau_1 > \eta\sqrt{n}/2) \\ &\leq \frac{(n+1)\mathbb{P}^o(\tau_1 > \eta\sqrt{n}/2)}{\mathbb{P}^o(D')} \rightarrow_{n \rightarrow \infty} 0, \end{aligned}$$

where we used the moment bounds on $\mathbb{E}^o(\tau_2 - \tau_1)^{2+\delta}$ and the fact that $\mathbb{P}^o(D') > 0$ in the last limit. This yields the conclusion. Further, one observes that the limiting covariance of X_n/\sqrt{n} is $\Xi/(\mathbb{E}^o(\tau_2 - \tau_1))$. \square

A direct conclusion of Theorem 3.5.24 is that under Kalikow’s condition, X_n/\sqrt{n} satisfies an annealed CLT.

Bibliographical notes: Lemma 3.5.2, Kalikow’s condition, the fact that it implies $\mathbb{P}^o(A_\ell) = 1$, and Lemma 3.5.8 appeared in [38]. The argument for $v_\ell > 0$ under Kalikow’s condition is due to Sznitman and Zerner [76], who also observed Corollary 3.5.10. [71] proves that in the i.i.d. environment case, $a(0, z) = 0$ if and only if $z = tv$, some $t > 0$. The estimates in Theorem 3.5.16 are a weak form of estimates contained in [71]. Finally, [81] characterizes, under Kalikow’s condition, the speed v as a function of Lyapunov exponents closely related to the functions $a(\lambda, z)$.

In a recent series of papers, Sznitman has shown that many of the conclusions of this section remain valid under a weaker condition, Sznitman’s (T) or (T') conditions, see [74, 73, 75]

Appendix

Markov chains and electrical networks: a quick reminder

With (V, E) as in Section 1.1, let $C_e \geq 0$ be a *conductance* associated to each edge $e \in E$. Assume that we can write

$$\omega_v(w) = \frac{C_{vw}}{\sum_{w \in N_v} C_{vw}} := \frac{C_{vw}}{C_v}.$$

To each such graph we can associate an electrical network: edges are replaced by conductors with conductance C_{vw} . The relation between the electrical network and the random walk on the graph is described in a variety of texts,

see e.g. [25] for an accessible summary or [57] for a crash course. This relation is based on the uniqueness of harmonic functions on the network, and is best described as follows: fix two vertices $v, w \in V$, and apply a unit voltage between v and w . Let $V(z)$ denote the resulting voltage at vertex z . Then,

$$P_w^z(\{X_n\} \text{ hits } v \text{ before hitting } w) = V(z).$$

Recall that for any two vertices v, w , the *effective conductance* $C^{\text{eff}}(v \leftrightarrow w)$ is defined by applying a unit voltage between v and w and measuring the outflow of current at v . In formula, this is equivalent to

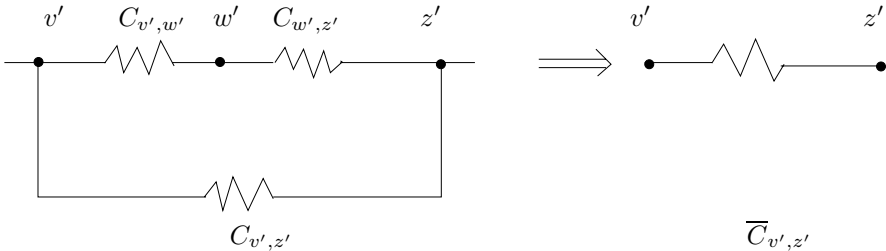
$$C^{\text{eff}}(v \leftrightarrow w) = \sum_{v' \in N_v} [1 - V(v')]C_{vv'} = \sum_{w' \in N_w} V(w')C_{ww'}.$$

For any integer r , the effective conductance $C_{v,r}$ between v and the horocycle of distance r from v is the effective conductance between v and the vertex r' in a modified graph where all vertices in the horocycle have been identified. We set then $C_{v,\infty} := \lim_{r \rightarrow \infty} C_{v,r}$. The effective conductance obeys the following rule:

Combination rule: Edges in parallel can be combined by summing their conductances. Further, the effective conductance between vertices v, w is not affected if, at any vertex $w' \notin \{v, w\}$ with $N_{w'} = \{v', z'\}$, one removes the edges (v', w') and (z', w') and replaces the conductance $C_{v',z'}$ by

$$\bar{C}_{v',z'} = C_{v',z'} + \left(\frac{1}{C_{v',w'}} + \frac{1}{C_{z',w'}} \right)^{-1}.$$

(This formula applies even if an edge $C_{v',w'}$ is not present, by taking $C_{v',z'} = 0$.)



Exercise A.1 Prove formulae (2.1.3) and (2.1.4).

Markov chains of the type discussed here possess an easy criterion for recurrence: a vertex v is recurrent if and only if the effective conductance $C_{v,\infty}$ between v and ∞ is 0. A sufficient condition for recurrence is given by means of the Nash-Williams criterion (see [57, Corollary 9.2]). Recall that an edge-cutset Π separating v from ∞ is a set of edges such that any path starting at v which includes vertices of arbitrarily large distance from v must include some edge in Π .

Lemma A.2 (Nash-Williams) *If Π_n are disjoint edge-cutsets which separate v from ∞ , then*

$$C_{v,\infty} \leq \left(\sum_n \left(\sum_{e \in \Pi_n} C_e \right)^{-1} \right)^{-1}.$$

As an application of the Nash-Williams criterion, we prove that a product of independent Sinai’s walks is recurrent. Recall that a Sinai walk (in dimension 1) is a RWRE satisfying Assumption 2.5.1. For simplicity, we concentrate here on Sinai’s walk without holding times and define a product of Sinai’s walk in dimension d as the RWRE on \mathbb{Z}^d constructed as follows: for each $v \in \mathbb{Z}^d$, set $N_v = \times_{i=1}^d (v_i - 1, v_i + 1)$ and let $\Omega = \times_{i=1}^d (M_1(N_v))^{\mathbb{Z}}$. For $z \in \mathbb{Z}^d$, we set $\omega_{i,z}^+ = \omega_i(z_i, z_i + 1)$, $\omega_{i,z}^- = \omega_i(z_i, z_i - 1)$ and $\rho_i(z) = \omega_{i,z}^- / \omega_{i,z}^+$. We equip Ω with a product of measures $P = \times_{i=1}^d P_i$, such that each P_i is a product measure which also satisfies Assumption 2.5.1. For a fixed $\omega \in \Omega$, define the RWRE in environment ω as the Markov chain (of law P_ω^o) such that $P_\omega^o(X_0 = 0) = 1$ and, for $v \in \{-1, 1\}^d$, $P_\omega^o(X_{n+1} = x + v | X_n = x) = \prod_{i=1}^d \omega_i(x_i, x_i + v_i)$. Define

$$C(x, v) = \prod_{i=1}^d \left(\left(\prod_{j_i=x_i}^{-1} \rho_i(j_i) \right) \left(\prod_{j_i=0}^{x_i-1} \rho_i(j_i)^{-1} \right) (\rho_i(x_i)^{-1})^{(v_i+1)/2} \right),$$

where by definition a product over an empty set of indices equals 1. Then, the resistor network with conductances $C(x, v)$ is a model for the product of Sinai’s RWRE. Define

$$B_i^n(t) = -\frac{1}{\sqrt{n}} \sum_{j=0}^{\lfloor nt \rfloor} \log \rho_i(j) \cdot (\text{sign } t).$$

Then,

$$C(x, v) \leq \varepsilon^{-d} \prod_{i=1}^d e^{\sqrt{n} B_i^n(x_i/n)}.$$

Taking as cutsets Π_n the set of edges $(x, x + v)$ with $|x|_\infty = n$, $v_i \in -1, 1$ and $|x + v|_\infty = n + 1$, we thus conclude that

$$\left(\sum_{e \in \Pi_n} C_e \right) \leq \varepsilon^{-d} \sum_{i=1}^d (e^{\sqrt{n} B_i^n(1)} + e^{\sqrt{n} B_i^n(-1)}) \prod_{j=1, j \neq i}^d \left(\sum_{k=-n}^n e^{\sqrt{n} B_j^n(k/n)} \right) =: D_n.$$

Since P_i are product measures, we have by Kolmogorov’s 0-1 law that $P(\liminf_{n \rightarrow \infty} D_n = 0) \in \{0, 1\}$. On the other hand, for all n large enough, we have by the CLT that

$$P(D_n < e^{-n^{1/4}}) \geq P(B_i^n(1) \leq -1, B_i^n(-1) \leq -1, \sup_{-1 \leq t \leq 1} B_i^n \leq 1/2d, i = 1, \dots, d) \geq c,$$

for some constant $c > 0$ independent of n . Thus, by Fatou's lemma, $P(\liminf_{n \rightarrow \infty} D_n = 0) > 0$, and hence $= 1$ by the above mentioned 0-1 law. We conclude from Nash's criterion (Lemma A.2) that $C_{0,\infty} = 0$, establishing the recurrence as claimed.

Exercise A.3 *Extend the above considerations to Sinai's walk with holding times and non product measures P_i .*

Bibliographical notes: The classical reference for the link between electrical networks and Markov chain is the lovely book [25]. The application to the proof of recurrence of products of Sinai's walks was prompted by a question of N. Gantert and Z. Shi.

References

1. S. Alili, Asymptotic behaviour for random walks in random environments, *J. Appl. Prob.* **36** (1999) pp. 334–349.
2. P. Billingsley, *Convergence of probability measures*, 2-nd. edition, Wiley (1999).
3. E. Bolthausen and I. Goldsheid, Recurrence and transience of random walks in random environments on a strip, *Comm. Math. Phys.* **214** (2000), pp. 429–447.
4. E. Bolthausen, A. S. Sznitman and O. Zeitouni, Cut points and diffusive random walks in random environments, *Ann. Inst. H. Poincaré - Prob. Stat.* **39** (2003), pp. 527–555.
5. R. C. Bradley, A caution on mixing conditions for random fields, *Stat. & Prob. Letters* **1989**, pp. 489–491.
6. M. Bramson and R. Durrett, Random walk in random environment: a counterexample?, *Comm. Math. Phys.* **119** (1988), pp. 199–211.
7. J. Brémont, Marches aléatoire en milieu aléatoire sur \mathbb{Z} ; Dynamique d'applications localement contractantes sur le cercle. *Thesis*, Université de Rennes I, 2002.
8. J. Brémont, Récurrence d'une marche aléatoire symétrique dans \mathbb{Z}^2 en milieu aléatoire, *preprint* (2000).
9. J. Bricmont and A. Kupiainen, Random walks in asymmetric random environments, *Comm. Math. Phys.* **142** (1991), pp. 345–420.
10. W. Bryc and A. Dembo, Large deviations and strong mixing, *Ann. Inst. H. Poincaré - Prob. Stat.* **32** (1996) pp. 549–569.
11. F. Comets. Large deviations estimates for a conditional probability distribution. Applications to random interacting Gibbs measures. *Prob. Th. Rel. Fields* **80** (1989) pp. 407–432.
12. F. Comets, N. Gantert and O. Zeitouni, Quenched, annealed and functional large deviations for one dimensional random walk in random environment, *Prob. Th. Rel. Fields* **118** (2000) pp. 65–114.
13. F. Comets and O. Zeitouni, A law of large numbers for random walks in random mixing environments, to appear, *Annals Probab.* (2003).

14. A. De Masi, P. A. Ferrari, S. Goldstein and W. D. Wick, An invariance principle for reversible Markov processes. Applications to random motions in random environments, *J. Stat. Phys.* **55** (1989), pp. 787–855.
15. A. Dembo, N. Gantert, Y. Peres and O. Zeitouni, Large deviations for random walks on Galton-Watson trees: averaging and uncertainty, *Prob. Th. Rel. Fields* **122** (2002), pp. 241–288.
16. A. Dembo, N. Gantert and O. Zeitouni, Large deviations for random walk in random environment with holding times, to appear, *Annals Probab.* (2003).
17. A. Dembo, A. Guionnet and O. Zeitouni, Aging properties of Sinai's random walk in random environment, *XXX preprint archive*, math.PR/0105215 (2001).
18. A. Dembo, Y. Peres and O. Zeitouni, Tail estimates for one-dimensional random walk in random environment, *Comm. Math. Physics* **181** (1996) pp. 667–684.
19. A. Dembo and O. Zeitouni, *Large deviation techniques and applications*, 2nd edition, Springer, New-York (1998).
20. Y. Derriennic, Quelques application du théorème ergodique sous-additif, *Asterisque* **74** (1980), pp. 183–201.
21. J. D. Deuschel and D. W. Stroock, *Large deviations*, Academic Press, Boston (1989).
22. M. D. Donsker and S. R. S. Varadhan, Asymptotic evaluation of certain Markov process expectations for large time III, *Comm. Pure Appl. Math.* **29** (1976) pp. 389–461.
23. M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain Markov process expectations for large time, IV. *Comm. Pure Appl. Math.* **36** (1983) pp. 183–212.
24. P. Le Doussal, C. Monthus and D. S. Fisher, Random walkers in one-dimensional random environments: exact renormalization group analysis. *Phys. Rev. E* **59** (1999) pp. 4795–4840.
25. P. G. Doyle and L. Snell, *Random walks and electric networks*, Carus Math. Monographs **22**, MAA, Washington (1984).
26. R. Durrett, *Probability: theory and examples*, 2nd ed., Duxbury Press, Belmont (1996).
27. H. Föllmer, *Random fields and diffusion processes*, Lecture Notes in Mathematics **1362** (1988), pp. 101–203.
28. N. Gantert, Subexponential tail asymptotics for a random walk with randomly placed one-way nodes, *Ann. Inst. Henri Poincaré - Prob. Stat.* **38** (2002) pp. 1–16.
29. N. Gantert and O. Zeitouni, Large deviations for one dimensional random walk in a random environment - a survey, *Bolyai Society Math Studies* **9** (1999) pp. 127–165.
30. N. Gantert and O. Zeitouni, Quenched sub-exponential tail estimates for one-dimensional random walk in random environment, *Comm. Math. Physics* **194** (1998) pp. 177–190.
31. A. O. Golosov, Limit distributions for random walks in random environments, *Soviet Math. Dokl.* **28** (1983) pp. 18–22.
32. A. O. Golosov, Localization of random walks in one-dimensional random environments, *Comm. Math. Phys.* **92** (1984) pp. 491–506.
33. A. O. Golosov, On limiting distributions for a random walk in a critical one dimensional random environment, *Comm. Moscow Math. Soc.* **199** (1985) pp. 199–200.

34. A. Greven and F. den Hollander, Large deviations for a random walk in random environment, *Annals Probab.* **22** (1994) pp. 1381–1428.
35. Y. Hu and Z. Shi, The limits of Sinai’s simple random walk in random environment, *Annals Probab.* **26** (1998), pp. 1477–1521.
36. Y. Hu and Z. Shi, The problem of the most visited site in random environments, *Prob. Th. Rel. Fields* **116** (2000), pp. 273–302.
37. B. D. Hughes, *Random walks and random environments*, Oxford University Press (1996).
38. S. A. Kalikow, Generalized random walks in random environment, *Annals Probab.* **9** (1981), pp. 753–768.
39. M. S. Keane and S. W. W. Rolles, Tubular recurrence, *Acta Math. Hung.* **97** (2002), pp. 207–221.
40. H. Kesten, Sums of stationary sequences cannot grow slower than linearly, *Proc. AMS* **49** (1975) pp. 205–211.
41. H. Kesten, The limit distribution of Sinai’s random walk in random environment, *Physica* **138A** (1986) pp. 299–309.
42. H. Kesten, M. V. Kozlov and F. Spitzer, A limit law for random walk in a random environment, *Comp. Math.* **30** (1975) pp. 145–168.
43. E. S. Key, Recurrence and transience criteria for random walk in a random environment, *Annals Probab.* **12** (1984), pp. 529–560.
44. T. Komorowski and G. Krupa, The law of large numbers for ballistic, multi-dimensional random walks on random lattices with correlated sites, *Ann. Inst. H. Poincaré - Prob. Stat.* **39** (2003), pp. 263–285.
45. S. M. Kozlov, The method of averaging and walks in inhomogeneous environments, *Russian Math. Surveys* **40** (1985) pp. 73–145.
46. H.-J. Kuo and N. S. Trudinger, Linear elliptic difference inequalities with random coefficients, *Math. of Computation* **55** (1990) pp. 37–53.
47. G.F. Lawler, Weak convergence of a random walk in a random environment, *Comm. Math. Phys.* **87** (1982) pp. 81–87.
48. G.F. Lawler, A discrete stochastic integral inequality and balanced random walk in a random environment, *Duke Mathematical Journal* **50** (1983) pp. 1261–1274.
49. G.F. Lawler, Estimates for differences and Harnack inequality for difference operators coming from random walks with symmetric, spatially inhomogeneous, increments, *Proc. London Math. Soc.* **63** (1991) pp. 552–568.
50. F. Ledrappier, *Quelques propriétés des exposants caractéristiques*, Lecture Notes in Mathematics **1097**, Springer, New York (1984).
51. R. Lyons, R. Pemantle and Y. Peres, Ergodic theory on Galton–Watson trees: speed of random walk and dimension of harmonic measure, *Ergodic Theory Dyn. Systems* **15** (1995), pp. 593–619.
52. R. Lyons, R. Pemantle and Y. Peres, Biased random walk on Galton–Watson trees, *Probab. Theory Relat. Fields* **106** (1996), pp. 249–264.
53. S. A. Molchanov, *Lectures on random media*, Lecture Notes in Mathematics **1581**, Springer, New York (1994).
54. S. V. Nagaev, Large deviations of sums of independent random variables, *Annals Probab.* **7** (1979) pp. 745–789.
55. S. Olla, Large deviations for Gibbs random fields. *Prob. Th. Rel. Fields* **77** (1988) pp. 343–357.
56. S. Orey and S. Pelikan, Large deviation principles for stationary processes. *Annals Probab.* **16** (1988), pp. 1481–1495.

57. Y. Peres, *Probability on trees: an introductory climb*, Lecture notes in Mathematics **1717** (1999) P. Bernard (Ed.), pp. 195–280.
58. D. Piau, Sur deux propriétés de dualité de la marche au hasard en environnement aléatoire sur \mathbb{Z} , *preprint* (2000).
59. D. Piau, Théorème central limite fonctionnel pour une marche au hasard en environnement aléatoire, *Ann. Probab.* **26** (1998), pp. 1016–1040.
60. A. Pisztor and T. Povel, Large deviation principle for random walk in a quenched random environment in the low speed regime, *Annals Probab.* **27** (1999) pp. 1389–1413.
61. A. Pisztor, T. Povel and O. Zeitouni, Precise large deviation estimates for one-dimensional random walk in random environment, *Prob. Th. Rel. Fields* **113** (1999) pp. 135–170.
62. F. Rassoul-Agha, A law of large numbers for random walks in mixing random environment, to appear, *Annals Probab.* (2003).
63. L. Shen, Asymptotic properties of certain anisotropic walks in random media, *Annals Applied Probab.* **12** (2002), 477–510.
64. M. Sion, On general minimax theorems, *Pacific J. Math* **8** (1958), pp. 171–176.
65. Z. Shi, Sinai’s walk via stochastic calculus, in *Milieux Aléatoires*, F. Comets and E. Pardoux, eds., Panoramas et Synthèses **12**, Société Mathématique de France (2001).
66. Ya. G. Sinai, The limiting behavior of a one-dimensional random walk in random environment, *Theor. Prob. and Appl.* **27** (1982) pp. 256–268.
67. F. Solomon, Random walks in random environments, *Annals Probab.* **3** (1975) pp. 1–31.
68. A. S. Sznitman, *Brownian motion, obstacles and random media*, Springer-Verlag, Berlin (1998).
69. A. S. Sznitman, *Lectures on random motions in random media*, In DMV seminar **32**, Birkhauser, Basel (2002).
70. A. S. Sznitman, Milieux aléatoires et petites valeurs propres, in *Milieux Aléatoires*, F. Comets and E. Pardoux, eds., Panoramas et Synthèses **12**, Société Mathématique de France (2001).
71. A. S. Sznitman, Slowdown estimates and central limit theorem for random walks in random environment, *JEMS* **2** (2000), pp. 93–143.
72. A. S. Sznitman, Slowdown and neutral pockets for a random walk in random environment, *Probab. Th. Rel. Fields* **115** (1999), pp. 287–323.
73. A. S. Sznitman, On a class of transient random walks in random environment, *Annals Probab.* **29** (2001), pp. 724–765.
74. A. S. Sznitman, An effective criterion for ballistic behavior of random walks in random environment, *Probab. Theory Relat. Fields* **122** (2002), pp. 509–544.
75. A. S. Sznitman, On new examples of ballistic random walks in random environment, *Annals Probab.* **31** (2003), pp. 285–322.
76. A. S. Sznitman and M. Zerner, A law of large numbers for random walks in random environment, *Annals Probab.* **27** (1999) pp. 1851–1869.
77. M. Talagrand, A new look at independence, *Annals Probab.* **24** (1996), pp. 1–34.
78. N. S. Trudinger, Local estimates for subsolutions and supersolutions of general second order elliptic quasilinear equations, *Invent. Math.* **61** (1980), pp. 67–79.
79. S. R. S. Varadhan, Large deviations for random walks in a random environment, *preprint* (2002).

80. M. P. W. Zerner, Lyapounov exponents and quenched large deviations for multidimensional random walk in random environment, *Annals Probab.* **26** (1998), pp. 1446–1476.
81. M. P. W. Zerner, Velocity and Lyapounov exponents of some random walks in random environments, *Ann. Inst. Henri Poincaré - Prob. Stat.* **36** (2000), pp. 737–748.
82. M. P. W. Zerner and F. Merkl, A zero-one law for planar random walks in random environment, *Annals Probab.* **29** (2001), pp. 1716–1732.
83. M. P. W. Zerner, A non-ballistic law of large numbers for random walks in i.i.d. random environment, *Elect. Comm. in Probab.* **7** (2002), pp. 181–187.

List of Participants

AMIDI Ali	Beheshti University, Tehran, Iran
ARNAUDON Marc	Univ. de Poitiers, F
ASCI Claudio	Universita degli Studi di L'Aquila, Italy
BAHADORAN Christophe	Univ. Blaise Pascal, Clermont-Ferrand, F
BALDI Paolo	Universita Roma Tor Vergata, Italy
BARDET Jean-Baptiste	Ecole Polytechnique Fédér. de Lausanne, CH
BEN AROUS Gérard	Ecole Polytechnique Fédér. de Lausanne, CH
BERARD Jean	Univ. Claude Bernard, Lyon, F
BERNARD Pierre	Univ. Blaise Pascal, Clermont-Ferrand, F
BOLTHAUSEN Erwin	Univ. de Zurich, CH
BOUGEROL Philippe	Univ. Pierre et Marie Curie, Paris, F
BOURRACHOT Ludovic	Univ. Blaise Pascal, Clermont-Ferrand, F
BRETON Jean-Christophe	Univ. Lille 1, F
CAMPILLO Fabien	INRIA, Marseille, F
CERNY Jiri	Ecole Polytechnique Fédér. de Lausanne, CH
CHAMPAGNAT Nicolas	Ecole Normale Supérieure de Paris, F
CLIMESCU-HAULICA Adriana	Comm. Research Center, Ottawa, Canada
DA SILVA Soares Ana	Univ. Libre de Bruxelles, Belgique
DARWICH Abdul	Univ. d'Angers, F
DEMBO Amir	Stanford University, USA
DJELLOUT Hacene	Univ. Blaise Pascal, Clermont-Ferrand, F
DUDOIGNON Lorie	INRIA, Marseille, F
FEDRIGO Mattia	Scuola Normale Superiore di Pisa, Italy
FERRIERE Régis	Ecole Normale Supérieure de Paris, F
GIACOMIN Giambattista	Univ. Denis Diderot, Paris, F
GILLET Florent	Univ. Henri Poincaré, Nancy, F
GROSS Thierry	Univ. Denis Diderot, Paris, F
GUILLIN Arnaud	Univ. Blaise Pascal, Clermont-Ferrand, F
GUIONNET Alice	Ecole Normale Supérieure de Lyon, F
HOFFMANN Marc	Univ. Denis Diderot, Paris, F
KERKYACHARIAN Gérard	Univ. Denis Diderot, Paris, F
KISTLER Nicolas	Univ. de Zurich, CH
LAREDO Catherine	INRA, Jouy-en-Josas, F

LOECHERBACH Eva	Univ. de Paris Val de Marne, F
LOPEZ-MIMBELA Jose Alfredo	CIMAT, Guanajuato, Mexico
LORANG Gerard	Centre Universitaire de Luxembourg
MILLET Annie	Univ. de Paris 10, F
MOUSSET Sylvain	Univ. Pierre et Marie Curie, Paris, F
NICAISE Florent	Univ. Blaise Pascal, Clermont-Ferrand, F
NUALART Eulalia	Ecole Polytechnique Fédér. de Lausanne, CH
OCONE Daniel	Rutgers University, Piscataway, NJ, USA
PARDOUX Etienne	Univ. de Provence, Marseille, F
PAROISSIN Christian	Univ. René Descartes, Paris, F
PIAU Didier	Univ. Claude Bernard, Lyon, F
PICARD Dominique	Univ. Denis Diderot, Paris, F
PICARD Jean	Univ. Blaise Pascal, Clermont-Ferrand, F
PLAGNOL Vincent	Ecole Normale Supérieure de Paris, F
RASSOUL-AGHA Firas	Courant Institute, New York, USA
ROUAULT Alain	Univ. de Versailles-St Quentin, F
ROUX Daniel	Univ. Blaise Pascal, Clermont-Ferrand, F
ROZENHOLC Yves	INRA, Paris, F
SKORA Dariusz	University of Wroclaw, Poland
SOOS Anna	Babes-Bolyai Univ., Cluj-Napoca, Romania
SORTAIS Michel	Ecole Polytechnique Fédér. de Lausanne, CH
WU Liming	Univ. Blaise Pascal, Clermont-Ferrand, F

List of Short Lectures

Claudio ASCI, Generating uniform random vectors.

Christophe BAHADORAN, Boundary conditions for driven conservative particle systems.

Jean-Baptiste BARDET, Limit theorems for coupled analytic maps.

Jean BÉRARD, Genetic algorithms in random environments.

Erwin BOLTHAUSEN, A fixed point approach to weakly self-avoiding random walks.

Philippe BOUGEROL, A path representation of the eigenvalues of the GUE random matrices.

Jean-Christophe BRETON, Multiple stable stochastic integrals: representation, absolute continuity of the law.

Jiří ČERNÝ, Critical path analysis for continuum percolation.

Adriana CLIMESCU-HAULICA, Cramér decomposition and noise modelling: applications from/to communications theory.

Ana DA SILVA SOARES, Files d'attente fluides.

Amir DEMBO, Random polynomials having few or no real zeros.

Mattia FEDRIGO, A multifractal model for network data traffic.

Florent GILLET, Algorithmes de tri: analyse du coût, stabilité face aux erreurs.

Alice GUIONNET, Enumerating graphs, matrix models and spherical integrals.

Eva LÖCHERBACH, On the invariant density of branching diffusions.

José Alfredo LÓPEZ-MIMBELA, A proof of non-explosion of a semilinear PDE system.

Florent NICAISE, Infinite volume spin systems: an application to Girsanov results on the Poisson space.

Eulalia NUALART, Potential theory for hyperbolic SPDE's.

Dan Ocone, Finite-fuel singular control with discretionary stopping.

Didier PIAU, Mutation-replication statistics of polymerase chain reactions.

Yves ROZENHOLC, Classification trees and colza diversity.

Michel SORTAIS, Large deviations in the Langevin dynamics of short range disordered systems.