

Mélykúti Bence:

**The Mixing Rate of
Markov Chain Monte Carlo Methods
and some Applications of
MCMC Simulation in Bioinformatics**

Szakdolgozat

Eötvös Loránd Tudományegyetem
Természettudományi Kar
matematikus szak

Témavezető:
Miklós István (Növényrendszertani és Ökológiai Tanszék)
Márkus László (Valószínűségelméleti és Statisztika Tanszék)

2006.

MÉLYKÚTI Bence

**The Mixing Rate of
Markov Chain Monte Carlo Methods
and some Applications of
MCMC Simulation in Bioinformatics**

Thesis
for the degree
MSc in mathematics

Eötvös Loránd University
Faculty of Science
Budapest, Hungary

Advisors:

MIKLÓS István (Department of Plant Taxonomy and Ecology)

MÁRKUS László (Department of Probability Theory and Statistics)

2006

Contents

1	Summary	7
2	Introduction	10
2.1	Sorting signed permutations by reversals	10
2.2	Multiple sequence alignment	13
2.3	On the notations to be used	15
3	Markov chain Monte Carlo methods	17
3.1	Markov chains	17
3.2	An overview of MCMC simulation	21
3.3	The Metropolis Algorithm	23
3.4	The Metropolis–Hastings Algorithm	24
3.5	Metropolised Independent Sampler	25
3.6	Partial Independent Sampler	26
3.7	Gibbs Sampler	31
4	Theoretical results on the mixing rate	33
4.1	Measuring the distance between distributions	33
4.2	Convergence to steady state	35
4.3	Gershgorin’s bound	42
4.4	Liu’s result for MIS	44
4.5	Conductance	44
4.6	Canonical (distinguished) paths, Poincaré coefficient	46
4.7	Similar bounds with different coefficients	49
4.8	Multicommodity flow	52
4.9	Probabilistic inequalities	54
4.10	Dobrushin’s inequality	55
4.11	Convergence rate for the Gibbs Sampler	56
4.12	The coupling method	57

5 Applications	59
5.1 Ladder	59
5.2 MIS is slow at sorting by reversals	66
5.3 Dragon’s wing	69
5.4 Another example	72
5.5 The ParIS needs to choose the whole sorting sequence with positive probability	75
5.6 An application of the coupling method	81
5.7 Conclusions and future work	86
Acknowledgements	87
References	88
Abbreviations	91

1 Summary

In the last decades, Markov chain simulation has become an important and widely popular computational paradigm. Its extreme flexibility and power made it an indispensable tool in many application areas, including statistics, computer science, physics, material science, chemistry, biology, engineering, economics and finance.

The chief applications of Markov chain Monte Carlo (MCMC) methods are random sampling from specified probability distributions and estimating the expectation of certain functions.

A common property of all MCMC methods is that one has to iterate a typically simple, stochastic step a large number of times (i.e. make transitions in a Markov chain), and after that, one has to stop the chain and read out the result. If one is interested in random sampling, then the result is the last state of the chain: the probability distribution of this state is an approximation of the target distribution, so the last state is a random sample from a distribution that approximates the target distribution.

However, it is a difficult question, how long the iteration should be run to get sufficiently close to the target distribution, because a Markov chain cannot give any direct indications about this.

The scope of this thesis is twofold. Firstly, we study theoretical results of Markov chain theory in this crucial area, and secondly, we apply some of these results to two problems of biocomputing.

In Section 2 we introduce our motivating problems. We explain two problems of bioinformatics we wish to solve using MCMC techniques. The first problem is to construct a *fast* algorithm that samples from the set of all optimal sorting sequences of signed permutations, uniformly at random, if reversals are the only allowed sorting transformations. The second problem is sampling from the set of all multiple sequence alignments, uniformly at random. Again, we wish to find a fast, approximative stochastic algorithm.

In the first part of Section 3 we present some notions and theorems of Markov chain theory that we will use in our investigations. After this short

introduction, an overview of the most classical MCMC methods is given, and we describe a less known method, the Partial Independent Sampler (ParIS), which becomes fundamental when we investigate our two biologically motivated problems.

In Section 4 we define the *variation distance* of probability distributions, which makes it possible for us to talk about the closeness of distributions. We give a sufficient condition that ensures convergence to a target distribution. Nevertheless, the largest part of this section is a collection of various approaches to bounding the running time a Markov chain (or specifically, an MCMC algorithm) needs to get close to its target distribution. This running time is the so-called *mixing time*.

We present our own results in Section 5.

We prove that there is a signed permutation with the property that the Metropolised Independent Sampler (MIS) is slow (i.e. it has exponential mixing time) at sampling from its optimal sorting sequences (Section 5.2). We suggest that the ParIS is superior to the MIS (Section 5.1).

We conducted investigations what window size the ParIS should use to have an acceptable (i.e. polynomial) mixing time. This question turned out to be a difficult one. We hope, but we could not prove yet, that there exists a window size distribution that results in a fast (polynomially mixing) ParIS algorithm for sampling from the set of all minimum sequences of sorting reversals, uniformly at random.

Section 5.3 shows a graph to which the uniform window size distribution applied gives a fast ParIS method. In contrast, Section 5.4 is about an example which satisfies the converse: even the uniform window size distribution results in a slow ParIS method.

Section 5.5 describes a signed permutation with the property that when sampling from its set of all minimum sequences of sorting reversals the ParIS needs to choose the whole sorting path as a window with positive probability to guarantee the irreducibility of the Markov chain.

As to our other problem, we found a sophisticated argument to prove

that in an oversimplified model of sequence alignments, the ParIS with window size 2, applied to sample from the set of all sequence alignments, has polynomial mixing time (Section 5.6). Admittedly, this result has almost no biological relevance, but from a mathematician's viewpoint its proof is interesting indeed, and the proof might be developed to yield similar bounds for models of more relevance.

2 Introduction

First we outline the motivation of our investigations, more precisely, we introduce the problems we would like to solve by using MCMC techniques. To keep it short, these parts are intended to be easily understood, but they are not claimed to be precise in terms of notation techniques.

2.1 Sorting signed permutations by reversals

First, we introduce basics of a mathematical model of genome rearrangements.

A **signed permutation** is a permutation $\sigma = (\sigma_1, \dots, \sigma_n)$ on the set $\{1, \dots, n\}$, where each number is also assigned a sign of plus or minus, e.g. $(+2, -4, -1, +5, -3)$. A **reversal** $\rho(i, j)$ ($i \leq j$) on σ transforms σ to

$$\sigma' = (\sigma_1, \dots, \sigma_{i-1}, -\sigma_j, -\sigma_{j-1}, \dots, -\sigma_i, \sigma_{j+1}, \dots, \sigma_n).$$

Biologists often use the word *inversion* for reversal. Since inversion has a different meaning in mathematics, we shall always use the word *reversal*.

The minimum number of reversals needed to transform one signed permutation into another one is called the **reversal distance** between them. We can assume that one of the signed permutations is the identity permutation $id = (+1, +2, \dots, +n)$ and the other one is defined relative to this one. The problem of **sorting signed permutations by reversals** is to find, for a given signed permutation σ , a sequence of minimum length of reversals that transforms σ into the identity permutation id .

A **transposition** is a transformation where a connected part (a block) is cut out of the signed permutation and it is relocated somewhere. In other words, two consecutive blocks are swapped:

$$(+2, \underline{-4, -1}, \underline{+5, -3}) \longrightarrow (+2, +5, -3, -4, -1).$$

An **inverted transposition** is a transposition, in which the moving part is inverted before relocation:

$$(+2, \underline{-4, -1}, +5, -3) \longrightarrow (+2, +5, -3, +1, +4).$$

Genetic and DNA data on many organisms is accumulating rapidly, and consequently the ability to compare genomes of different species is growing dramatically. Signed permutations and their transformations are useful tools in the comparative study of genomes. A genome of a species can be thought of as a set of ordered sequences of genes, each gene having an orientation. Different species often have similar genes that were inherited from common ancestors. However, these genes have been shuffled by mutations that modified the order and/or the directionality of genes. In this case the gene orders themselves make two species differ. One can find a more detailed explanation of (mitochondrial) genome rearrangements in [14].

A possible mathematical model for comparing two genomes represents genomes by signed permutations. For the sake of simplicity we assume that both genomes consist of exactly one chromosome. We also assume that both genomes share the same genes and each gene has a copy number one. The order of genes are represented by permutations and the orientation of genes in the genome is given by signs of plus and minus. Mutations are the transformations of signed permutations that were defined above.

One usually investigates the problem of sorting by reversals only. There are two reasons for this. First, sorting by reversals is of real biological relevance. It is widely accepted that reversal distance between two genomes (between the two corresponding signed permutations) provides a good estimate of the evolutionary distance between two species.

Secondly, in 1995, Hannenhalli and Pevzner developed a theory which yields a polynomial time algorithm for computing the reversal distance between two signed permutations and finding one of possibly many optimal sorting sequences. The presentation of the Hannenhalli–Pevzner theory is out of the scope of this thesis. The reader is referred to the original paper [10], or the book on computational molecular biology by Pevzner [19]. Much work, including considerable simplifications, has been done on the theory since its birth [13, 4, 5, 23].

The model would be more accurate in estimating the evolutionary dis-

tance, if other types of transformations were also allowed. There are approximation algorithms, but, unfortunately, today there is no exact polynomial time algorithm for computing the distance between two signed permutations using reversals plus other transformations. Therefore we restrict our attention to sorting by reversals only.

We aim to find a polynomial mixing time algorithm that samples from the set of all optimal sorting sequences of signed permutations by reversals, uniformly at random.

Suppose that we investigate signed permutations on $\{1, \dots, n\}$. The Hannenhalli–Pevzner theory tells us that the optimal sorting sequence starting from σ consists of not more than $n + 1$ reversals.

Let G be a directed graph defined as follows. Let the set V of nodes contain σ , id , and every other signed permutation that can be reached from σ by a sequence of reversals of which each reversal decreases the reversal distance between the current signed permutation and id . There is an oriented edge in the graph from π to π' if and only if the reversal distance between π' and id is one less than the distance between π and id . This definition gives a directed acyclic graph G . Our problem is equivalent to sampling from the set of directed paths in G leading from σ to id .

As we have mentioned, the paths are not longer than $n + 1$. The difficulty is that there may be loads of them: their number may be exponential in n . This means that we cannot get a polynomial time algorithm by enumerating all paths. Today there is no known polynomial time algorithm either to calculate the number of paths or to sample from their set by any other means, uniformly at random. We shall use MCMC methods instead, and our main concern is what can be said about the running time of these algorithms.

In Section 3.6 it is explained what strategies one can think of to construct an MCMC algorithm to solve this problem. In later sections we shall compare our two algorithms, the MIS and the Paris.

2.2 Multiple sequence alignment

A good introduction to this topic can be found in [8].

The problem of sequence alignment is the following: transform one sequence of characters over an alphabet Σ (typically $\Sigma = \{\text{A, C, G, T}\}$) into another sequence, or transform d sequences one into another, by the use of three elementary transformations: *insertion*, *deletion*, *substitution*. Let $-$ denote a gap. An insertion is a $- \rightarrow \text{X}$, a deletion is a $\text{Y} \rightarrow -$ transformation. See the example with $d = 3$. Lines without the gap signs are the sequences.

A	-	G	G	T	C	T	A
-	-	G	A	G	C	T	G
A	C	T	-	C	C	-	G

Each alignment has a score. The score represents a distance between sequences: the more they differ, the bigger the distance is. The score is usually defined by a sum of scores over all columns:

$$\sum_{i \text{ is a column}} \text{score}(i),$$

where the score of a single column is given by a sum-of-pairs: each column has $\binom{d}{2}$ pairs of letters, and each type of pairs is given a score. The problem of finding an extremal score alignment is provably NP-complete.

Each sequence alignment is assigned a probability, that is determined by the scores. The problem under our investigations is to sample random elements from this distribution.

The fastest known algorithm to sample from this distribution takes $O(2^d \prod_j L_j)$ steps, where L_1, L_2, \dots, L_d are the lengths of sequences. We wish to find a faster approximative algorithm by using an MCMC method.

A popular approach is the use of parallel tempering. (For a description of the algorithm, see [17].) Parallel tempering runs more different MCMC algorithms simultaneously. Estimating its running time seems to be beyond our reach. However, one can get a lower bound for the running time, if one gives a lower bound for the running time of one of those subalgorithms. For

that reason, we are going to examine a very special case: sampling from the set of sequence alignments, uniformly at random.

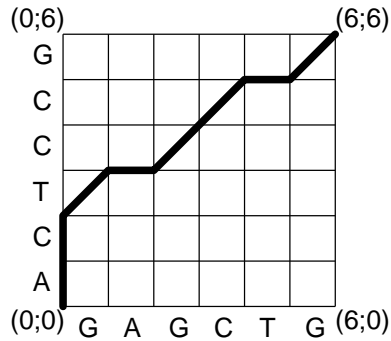


Figure 1: The alignment of two sequences

A sequence alignment can be given graphically. Suppose that $d = 2$. Pick the rectangle in \mathbb{Z}^2 with edge lengths L_1 and L_2 , and vertices $(0;0)$, $(L_1;0)$, $(L_1;L_2)$ and $(0;L_2)$. We associate the edge at the bottom with the first sequence and the edge on the left with the second one, with orientation given by the following rule: the first character of the first sequence is associated with the left end of the edge, the last one with the right. The second sequence is positioned with its first character at the bottom, the last one at the top. An alignment is represented by a path that connects $(0;0)$ and $(L_1;L_2)$. The path starts at $(0;0)$. It consists of steps from a point in \mathbb{Z}^2 to another one. The points that it crosses represent the gradual alignment of the sequences. Each step is of the following three kinds: \rightarrow , \uparrow , \nearrow (or equivalently, $(+1;0)$, $(0;+1)$, $(+1;+1)$). From the viewpoint of the first sequence, \uparrow is an insertion, \rightarrow is a deletion and \nearrow is a substitution or no transformation (depending on what the corresponding letters are: they may be different or identical). From the viewpoint of the second sequence, \uparrow is a deletion, \rightarrow is an insertion.

Figure 1 is the graphical representation of the following sequence alignment:

```

- - G A G C T G
A C T - C C - G

```

For $d \geq 3$, one can follow a similar method in a d dimensional space.

We found the problem to sample from the set of sequence alignments, uniformly at random, too difficult to handle. Hence, we made significant simplifications. We hope, but cannot prove that our simplified model is of real relevance in the environment described above.

Our model is the one that is represented by a square in \mathbb{Z}^2 , or a d dimensional hypercube in \mathbb{Z}^d . The other condition is far more restrictive: we use two types of steps only, \rightarrow and \uparrow . It is clear that this approach is inspired by the geometrical interpretation, and it does not have much sense in the world of DNA sequence alignments.

In fact, there is a perfect sampling method for this greatly simplified model. One can always easily calculate how many paths lead through a given point in \mathbb{Z}^2 , so one can compute the ratio between these quantities of any two points. Choosing each step between the two possible directions with probability distribution proportional to these quantities, one can build up a random path, with uniform distribution on the set of all paths, as one desired.

We are more concerned in this thesis with local decisions that arise in random walks (see the *random walk procedure* of Section 3.6) than with global ones of methods like the perfect sampling. The reason is that we do not hope that we can develop a global description of optimal sorting sequences in the problem of sorting signed permutations by reversals, or that we can perfectly understand the problem of sampling from the set of sequence alignments.

2.3 On the notations to be used

Giving value to a variable will be expressed by the combination of a colon and an equality sign, with the positioning of the colon indicating the new expression. For instance, $b =: a/5$ is to mean ‘whatever b is, choose a such that it satisfies the equality $b = a/5$ ’.

In certain cases we will use the notation *wedge* for the minimum, *vee* for the maximum of two numbers. Symbolically, $a \wedge b := \min\{a, b\}$ and $a \vee b := \max\{a, b\}$.

In some cases the sign \cdot (*dot*) will be used to express multiplication.
 χ_S denotes the indicator function of set S :

$$\chi_S(x) = \begin{cases} 1, & \text{if } x \in S, \\ 0, & \text{if } x \notin S. \end{cases}$$

3 Markov chain Monte Carlo methods

3.1 Markov chains

Throughout this work we will assume familiarity with the elementary theory of Markov chains. For those who have difficulties with this subject we advise studying the appropriate parts of Brémaud's textbook [6], which is a useful book and starts from the very basics of probability theory. First, let us recall some well-known definitions.

Definition 3.1.1 Let (Ω, \mathcal{A}, P) be a probability space, I a nonempty countable set (later to be called the **state space**) and $X_n : \Omega \rightarrow I$ a random variable for all $n \in \mathbb{N}$. If for all choices of nonnegative integers $n_1 < n_2 < \dots < n_k < m$, and for all $B \subseteq I$ and $i_1, i_2, \dots, i_k \in I$

$$P(X_m \in B \mid X_{n_1} = i_1, X_{n_2} = i_2, \dots, X_{n_k} = i_k) = P(X_m \in B \mid X_{n_k} = i_k),$$

(in other words, (X_n) has the **Markov property**), then the sequence (X_n) is called a **Markov chain**.

If for all $n \in \mathbb{N}$, $i, j \in I$

$$P(X_{n+1} = j \mid X_n = i) = p_{ij}$$

(that is, $P(X_{n+1} = j \mid X_n = i)$ does not depend on n), then it is a **homogeneous** Markov chain. The matrix $P = ((p_{ij}))_{i,j \in I}$ is the **transition matrix** of the homogeneous Markov chain. (One should be aware that P will denote both the probability measure and the transition matrix.)

In all cases we will examine homogeneous Markov chains, and for this reason, we will omit the word homogeneous.

Definition 3.1.2 The square matrix P indexed by I is a **stochastic matrix** if for all $i, j \in I$

$$p_{ij} \geq 0, \quad \sum_{k \in I} p_{ik} = 1.$$

Theorem 3.1.3 *Every transition matrix is a stochastic matrix. And conversely, for every stochastic matrix P , there exists a Markov chain with transition matrix P .*

Definition 3.1.4 The **transition graph** of a Markov chain defined by transition matrix P is a graph whose nodes are the states of the chain. There is an oriented edge in the graph from i to j (labelled by p_{ij}) if and only if $p_{ij} > 0$.

Definition 3.1.5 State $j \in I$ is **accessible** from state $i \in I$ if there exists some $n \in \mathbb{N}$ such that

$$p_{ij}^{(n)} := P(X_{m+n} = j \mid X_m = i) > 0.$$

States i and j are said to **communicate**, if j is accessible from i and i is accessible from j .

Definition 3.1.6 A Markov chain is **irreducible**, if all its states communicate.

Definition 3.1.7 The **period** of state $i \in I$ is

$$d(i) := \gcd \{n \geq 1 \mid p_{ii}^{(n)} > 0\}.$$

If $d(i) = 1$, then i is said to be **aperiodic**. Since communicating states have the same period, it makes sense to call an irreducible Markov chain **aperiodic**, if all its states are aperiodic.

Remark 3.1.8 Given $A \in \mathcal{A}$, $i \in I$, we will sometimes abbreviate $P(A \mid X_0 = i)$ to $P_i(A)$. If μ is a probability distribution on I , then $P_\mu(A) := \sum_{i \in I} \mu(i) P_i(A)$, which is the probability of A starting from an initial state chosen according to the distribution μ .

Definition 3.1.9 State $i \in I$ is called **recurrent** if

$$\sum_{n=1}^{\infty} P_i(X_1 \neq i, X_2 \neq i, \dots, X_{n-1} \neq i, X_n = i) = 1,$$

and otherwise it is called **transient**. We call a recurrent state i **positive recurrent** if

$$\sum_{n=1}^{\infty} nP_i(X_1 \neq i, X_2 \neq i, \dots, X_{n-1} \neq i, X_n = i) < \infty,$$

and otherwise **null recurrent**.

Proposition 3.1.10 *An irreducible Markov chain with finite state space is positive recurrent (i.e. all its states are positive recurrent).*

Definition 3.1.11 A Markov chain is **ergodic** if it is irreducible, positive recurrent and aperiodic.

There are a number of non-equivalent definitions of ergodicity in the literature. This definition binds us to using Brémaud's [6] terminology.

Definition 3.1.12 A probability distribution π satisfying

$$\pi^T = \pi^T P$$

is called a **stationary distribution** of the transition matrix P , or of the corresponding Markov chain.

Remark 3.1.13 If a Markov chain is started with a stationary distribution, then it keeps this probability distribution during all forthcoming steps. In this case we say that the chain is **stationary**, or equivalently, the chain is in a **stationary regime**, in **equilibrium** or in **steady state**.

There is a useful technical extension of the notion of stationary distribution.

Definition 3.1.14 A nonnegative, nonnull vector $x = (x_i)_{i \in I}$ is called an **invariant measure** of the stochastic matrix P , if

$$x^T = x^T P.$$

Theorem 3.1.15 *Let P be the transition matrix of an irreducible, recurrent Markov chain. Then there exists an invariant measure of P , whose each entry is positive.*

Theorem 3.1.16 *The invariant measure of the transition matrix of an irreducible, recurrent Markov chain is unique up to a multiplicative factor.*

Corollary 3.1.17 *Every irreducible, recurrent Markov chain has a stationary distribution, and this distribution is unique.*

Theorem 3.1.18 (Ergodic theorem) *Let (X_n) be an irreducible, positive recurrent Markov chain with the initial distribution μ and stationary distribution π , and let $f : I \rightarrow \mathbb{R}$ be such that*

$$\sum_{i \in I} |f(i)|\pi(i) < \infty.$$

Then, P_μ -almost-surely,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} f(X_k) = \sum_{i \in I} f(i)\pi(i).$$

Corollary 3.1.19 *Let (X_n) be an irreducible, positive recurrent Markov chain with the initial distribution μ and stationary distribution π , and let $f = \chi_{\{i\}}$ for a fixed $i \in I$. Then, P_μ -almost-surely,*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^{N-1} \chi_{\{X_k=i\}} = \pi(i).$$

Corollary 3.1.20 *By Proposition 3.1.10, in an irreducible Markov chain with finite state space the relative frequency of visits to certain states converges to the distribution π . More generally, this proposition is true for any irreducible positive recurrent Markov chain. Obviously, this includes ergodic chains.*

Definition 3.1.21 Suppose that π is a positive stationary distribution of a Markov chain. If for all $i, j \in I$ the **detailed balance equation** is satisfied:

$$\pi(i)p_{ij} = \pi(j)p_{ji},$$

then the Markov chain is said to be **reversible** with respect to π .

Proposition 3.1.22 *If a Markov chain is reversible with respect to π , then π is a stationary distribution of the Markov chain.*

Definition 3.1.23 Let π be a positive probability distribution on I . One can define $L_2(\pi)$ as a real vector space \mathbb{R}^I endowed with the scalar product

$$\langle x, y \rangle_\pi := \sum_{i \in I} x(i)y(i)\pi(i),$$

for $x = (x(i))_{i \in I}$, $y = (y(i))_{i \in I}$.

3.2 An overview of MCMC simulation

Markov chain simulation is a powerful algorithmic tool for random sampling (especially of combinatorial structures) from a specified probability distribution. The main idea is as follows.

Suppose we would like to generate samples from a large but finite set I of structures from a distribution π (later to be called the **target distribution**), or we wish to estimate the expectation of a scalar valued function f on I .

As to the first problem, construct a Markov chain which converges asymptotically to the stationary distribution π . Start the chain from an arbitrary state, and simulate it *long enough*, until it gets close to steady state. If one stops the chain, the distribution of the final state will be close to the desired distribution π .

For the second problem, calculate the empirical average of values of f in the visited states weighted by the number of visits to these states. By *Ergodic theorem 3.1.18*, this quantity (the empirical average) converges to the expectation value almost surely. Faster convergence can be realized by summing only after an initialization period (or *burn-in* period) of m steps:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=m+1}^{m+N} f(X_k) = \sum_{i \in I} f(i)\pi(i).$$

Obviously, the simulation needs a lot of computational work, and therefore it is always done using computers.

One might ask, what does *long enough* mean. This turns out to be a very delicate question. There is no direct indication of the chain being close to equilibrium.

A trivial approach in applications of MCMC simulation to the second problem is to run the simulation more times starting from different states, for the same number of steps. If one gets similar estimates for the expectation value in these different simulations, one may feel that this number of steps are enough. If one gets estimates significantly different, then it is obvious that results are worthless and one needs considerably longer run.

In the first problem, there are statistical methods to monitor the samples. However, they can only prove that one needs to run the simulation longer, but cannot prove that the simulation can be terminated.

Indications of such monitoring can be misleading. Suppose that the transition graph of the Markov chain consists of two almost separated halves, with only a very small number of edges with tiny transition probabilities between the two halves. Picturesquely speaking, the chain has a *bottleneck*. In this case, it is likely that the chain started from a certain state remains in one half of the state space for the whole simulation. This means that although our monitoring methods indicate that we can already use the chain to take samples, we should not, because we will not get samples from the other half of states.

The main concern of this thesis are to find out what theoretical results are in the literature on the running time needed, and how they can be applied to our problems.

Before getting down to this subject, let us study some of the most widely known MCMC methods. We will rely mostly on Liu's book [17] in this topic.

In most cases we shall omit the proof of convergence to the target distribution. With Metropolis-type algorithms these proofs are easy: one has to prove reversibility with respect to the target distribution, and apply Proposition 3.1.22.

3.3 The Metropolis Algorithm

Metropolis et al. [18] introduced the fundamental idea of evolving a Markov process to achieve the sampling of a probability distribution. Note that in Markov chain theory one usually knows the transition rule of a certain chain and investigates its properties, e.g. what its stationary distribution is. In MCMC simulation one is given a distribution and wants to find an efficient transition rule of a Markov chain, of which the stationary distribution is at hand. An immanent property of most MCMC sampling methods is that they can only provide statistically *dependent* samples, often highly correlated ones. Various attempts have been made in specific fields of applications to overcome this problem.

The Metropolis algorithm can be used to generate random samples from distribution π , that is known only up to a normalizing constant, or only the proportions of probabilities $\pi(i)$ ($i \in I = \{1, 2, \dots, r\}$) are known. Assume that $b(1), \dots, b(r) > 0$, and the target distribution has the form $\pi(i) = b(i)/B$. Evaluating $B = \sum_{i \in I} b(i)$ is trivial in theory, but in practice it is often more difficult, than the original problem of sampling from the distribution π was: either because r is too big, or because only the proportions of values $b(i)$ are known. That is the reason why an indirect sampling method is needed, for instance, the Metropolis algorithm.

Metropolis Algorithm

Starting with an initial state i_0 , the Metropolis algorithm iterates the following two steps.

1. Propose a small, random ‘unbiased perturbation’ of the current state i_t so as to generate a new one j . More accurately, generate j from a symmetric probability transition function (i.e. for all i, j $p_{ij} \geq 0$, $\sum_k p_{ik} = 1$ and $p_{ij} = p_{ji}$).
2. Generate a random variable $U \sim \text{Uniform}[0, 1]$ independently from

earlier variables.

$$i_{t+1} := \begin{cases} j, & \text{if } U \leq \pi(j)/\pi(i_t) = b(j)/b(i_t), \\ i_t, & \text{otherwise.} \end{cases}$$

In other words, accept the new state with probability $\min\{1, \pi(j)/\pi(i_t)\}$.

If it is not accepted, then reject it, and remain in the same state.

The symmetric probability transition function of Step 1 is often called the **proposal function**. Intuitively, the symmetry requirement means that there is no ‘trend bias’ at the proposal step.

3.4 The Metropolis–Hastings Algorithm

Hastings’ algorithm [11] is a generalization of the Metropolis algorithm. It omits the requirement of symmetry of the transition rule. There are new restrictions, though. These will be formulated right after the definition of the algorithm.

Suppose that the target distribution π is given on a set of r elements. First, define an irreducible Markov chain on $I = \{1, 2, \dots, r\}$ with transition matrix P . Using P one can define a new transition matrix Q on the same state space. For this, let us introduce for all $i, j \in I$, $i \neq j$ quantities

$$\alpha_{ij} := \min \left\{ 1, \frac{\pi(j)p_{ji}}{\pi(i)p_{ij}} \right\} = \min \left\{ 1, \frac{b(j)p_{ji}}{b(i)p_{ij}} \right\},$$

namely the **acceptance probabilities**.

Metropolis–Hastings Algorithm

Start with an initial state i_0 , and iterate the following two steps.

1. After t steps, visiting state i_t , generate j from the proposal distribution $p_{i_t j}$.
2. Draw a random variable $U \sim \text{Uniform}[0, 1]$.

$$i_{t+1} := \begin{cases} j, & \text{if } U \leq \alpha_{i_t j}, \\ i_t, & \text{otherwise.} \end{cases}$$

In other words, accept the new state with the acceptance probability α_{ij} . If it is not accepted, then reject it.

One can easily check that the transition matrix Q defined by this algorithm is the following:

$$q_{ij} = p_{ij}\alpha_{ij} = \min \left\{ p_{ij}, \frac{\pi(j)}{\pi(i)}p_{ji} \right\}, \quad \text{if } i \neq j,$$

$$q_{ii} = p_{ii} + \sum_{k \neq i} p_{ik}(1 - \alpha_{ik}) = 1 - \sum_{k \neq i} p_{ik}\alpha_{ik}.$$

This new Markov chain is to be used for the sampling. If P is chosen to meet the following two requirements:

$$i, j \in I \quad p_{ij} > 0 \Rightarrow p_{ji} > 0,$$

$$\exists i \in I : p_{ii} > 0,$$

then Q is irreducible and aperiodic, and it converges to the equilibrium distribution π . This is a sufficient, but not necessary condition.

3.5 Metropolised Independent Sampler

This specialized version of the Metropolis–Hastings algorithm was also proposed by Hastings, in [11].

Assume that we have a strictly positive *trial distribution* p on I at hand, that is thought to be ‘similar’ to the strictly positive *target distribution* π , from which we can draw independent samples. The Metropolised Independent Sampler (or **MIS**, for short) generates its proposed moves *independently* of the current state of the chain (i.e. p_{ij} does not depend on i). In comparison to Metropolis–Hastings algorithms, which usually make dependent local moves, MIS makes independent global jumps. We define the algorithm by giving the transition matrix Q of the corresponding Markov chain.

Metropolised Independent Sampler

Start with an initial state i_0 , and choose forthcoming states using the following transition rule:

$$q_{ij} = p(j) \min \left\{ 1, \frac{\pi(j)p(i)}{\pi(i)p(j)} \right\}, \quad \text{if } i \neq j,$$
$$q_{ii} = 1 - \sum_{k \neq i} p(k) \min \left\{ 1, \frac{\pi(k)p(i)}{\pi(i)p(k)} \right\}.$$

A frequently used notion related to MIS is the **importance weight**, or **importance ratio**: for $i \in I$, $w(i) := \pi(i)/p(i)$.

The intuition behind this algorithm is that a transition from i to j is accomplished by drawing an independent sample from p , and ‘thinning it down’ based on a comparison of the corresponding importance weights $w(j)$ and $w(i)$.

3.6 Partial Independent Sampler

Partial Independent Sampler (or **ParIS**, for short) is a modification of MIS for solving problems of a special family of combinatorial problems. This method will be explained using the problem *sorting signed permutations by reversals*.

Let us investigate signed permutations of the set $\{1, \dots, n\}$. Our objective is to sample elements from all optimal sorting sequences, uniformly at random. In Section 2.1 we explained how these sequences can be represented by paths in an oriented graph G : nodes of the graph are signed permutations, its edges connect two signed permutations if and only if one can be obtained from the other by exactly one reversal that leads closer to the target signed permutation (in most cases, id).

The main difficulty is that in general, this graph is not known entirely. We only know that from a certain node (which represents a signed permutation), which other nodes lead one step closer to the target, more strictly speaking, we can calculate this in polynomial time. It is clear that there are $\binom{n+1}{2}$ possible reversals that act on a given signed permutation. This gives an $O(n^2)$

upper bound for the out-degree of each node: there are $\binom{n+1}{2} = O(n^2)$ possible neighbours. The reversal distance between a signed permutation and id can be calculated in $O(n)$ time [4]. Consequently, there is an algorithm to find (or simply, count) all neighbours of a given node in $O(n^3)$ steps.

One can hardly find any better idea to draw samples from these paths uniformly at random, than to start from the initial signed permutation σ , and choose next state from its neighbours that are closer to the target signed permutation, uniformly at random. And iterate this process, always take a step leading closer to the target. Let us call this method the **random walk procedure**.

It is obvious that this algorithm can do very poorly. Imagine that we have two neighbours of the initial signed permutation to start with, but the structure of the underlying graph is such that there is only one path starting with one neighbour and there are many paths starting with the other. In this case the method would yield samples from a distribution that is far from being uniform: the first path would have probability $1/2$, the others would have much smaller probabilities.

The idea is to use MCMC simulation for sampling. The states of the chain will be the optimal sorting sequences. This definition often causes confusion, so we emphasize once again that the states are specific sequences of signed permutations and not the signed permutations themselves. Target distribution π is usually the uniform distribution. We need an effective transition rule to define our Markov chain.

MIS seems to be an applicable one. Trial distribution p is defined by the probabilities of paths in the graph of signed permutations under the random walk procedure. Starting from the initial signed permutation σ , $p(i)$ for a path i can be computed by multiplying the probabilities of each step one after another. Assume that i leads through nodes $i^0(= \sigma), i^1, i^2, \dots, i^n(= id)$. Then

$$p(i) = \prod_{r=0}^{n-1} \frac{1}{d(i^r)},$$

where $d(i^r)$ is the number of neighbours of i^r that are one step closer to i^n ,

than i^r . We will present three examples in Sections 5.1, 5.2, and 5.3, in which the MIS is too slow in terms of computational complexity.

ParIS is a sophisticated version of the MIS, which is intended to be faster than MIS. The concept is that we do not go back to start from scratch in every step of the Markov chain and draw independent samples, as with MIS, we only modify the path that corresponds to the current state of the Markov chain by cutting a part out of it (a ‘window’) and by replacing it with a new (preferably different) subsequence. The new subsequence has to be of the same length to get an optimal sorting sequence again.

We introduce the following notation: for states i, j and window w ,

$$p_{ij}^w := P^{\text{proposal}}(X_{t+1} = j, w \mid X_t = i),$$

that is, the proposal probability of transition to state j using window w , conditioned that the chain is in i .

Partial Independent Sampler

Choose an initial state (path) i_0 with the random walk procedure, and iterate the following two steps.

1. Pick a connected part (a window w) of $i := i_t$ at random, say,

$$w = \langle i^m, i^{m+1}, \dots, i^{m+k} \rangle \subseteq \langle i^0, i^1, \dots, i^n \rangle.$$

Using the random walk procedure, pick an optimal sorting sequence from i^m to i^{m+k} , for example

$$\langle i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k} \rangle.$$

2. Accept the new state (path)

$$j = \langle i^0, i^1, \dots, i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k}, \dots, i^n \rangle$$

with acceptance probability

$$\alpha_{ij}^w := \min \left\{ 1, \frac{\pi(j)p_{ji}^w}{\pi(i)p_{ij}^w} \right\}.$$

Equivalently, the transition probabilities of the ParIS are

$$q_{ij} = \sum_w p_{ij}^w \alpha_{ij}^w, \quad \text{if } i \neq j,$$

$$q_{ii} = 1 - \sum_{k \neq i} q_{ik}.$$

Proposition 3.6.1 *This Markov chain is reversible with respect to π . Consequently, if it is irreducible and recurrent, then its unique stationary distribution is π .*

PROOF We can use Proposition 3.1.22. Indeed,

$$\begin{aligned} \pi(i)q_{ij} &= \pi(i) \sum_w p_{ij}^w \alpha_{ij}^w = \pi(i) \sum_w p_{ij}^w \min \left\{ 1, \frac{\pi(j)p_{ji}^w}{\pi(i)p_{ij}^w} \right\} = \\ &= \sum_w \min \left\{ \pi(i)p_{ij}^w, \pi(j)p_{ji}^w \right\}, \end{aligned}$$

what equals to $\pi(j)q_{ji}$ because of the symmetry of the last term. ■

Step 1 can be done using different strategies. One can bind oneself to pick windows with fixed length, or one can draw a window size at random from a fixed distribution. Any strategy works, as long as p_{ij}^w is a fixed, well-defined quantity.

It has turned out that in the problem of sorting signed permutations by reversals, the probability of cutting the whole path out must be positive, otherwise it may happen that the Markov chain is not irreducible. (See Section 5.5). This means that applying ParIS to this problem cannot be done with a fixed window size, only if the size is the entire path, but this is not ParIS any more, but the MIS itself.

There are $n - k + 1$ possibilities to pick a window of length k from a path of length n ($k \leq n$). We have not carried out any investigations on which the optimal distribution is for possible locations of the window, we will always use uniform distribution.

If one is interested in sampling from the set of optimal sequences of sorting reversals, uniformly at random; if it is assumed that a certain window size distribution is fixed; the window is positioned in a place chosen from the possibilities uniformly at random; and the new subsequence is drawn by the random walk procedure, then the description of the algorithm can be simplified. In this case, if i can be transformed into j by cutting window w out, and the random window size W of w is k , then

$$p_{ij}^w = P(W = k) \frac{1}{n - k + 1} p(\langle i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k} \rangle),$$

consequently,

$$\begin{aligned} \frac{\pi(j)p_{ji}^w}{\pi(i)p_{ij}^w} &= \frac{p(\langle i^m, i^{m+1}, \dots, i^{m+k-1}, i^{m+k} \rangle)}{p(\langle i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k} \rangle)} = \\ &= \frac{\prod_{s=m+1}^{m+k-1} d(i^s)^{-1}}{\prod_{s=m+1}^{m+k-1} d(j^s)^{-1}} = \frac{\prod_{s=0}^{n-1} d(i^s)^{-1}}{\prod_{s=0}^{n-1} d(j^s)^{-1}} = \frac{p(i)}{p(j)}. \end{aligned}$$

Note that using the central term, $p(i)/p(j)$ can be evaluated in $O(kn^3)$ steps.

Partial Independent Sampler for this problem

Choose an initial state (path) i_0 with the random walk procedure, and iterate the following two steps.

1. Pick a connected part (a window w) of $i := i_t$ at random, say,

$$w = \langle i^m, i^{m+1}, \dots, i^{m+k} \rangle \subseteq \langle i^0, i^1, \dots, i^n \rangle.$$

Using the random walk procedure, pick an optimal sorting sequence from i^m to i^{m+k} , for example

$$\langle i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k} \rangle.$$

2. Accept the new state (path)

$$j = \langle i^0, i^1, \dots, i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k}, \dots, i^n \rangle$$

with acceptance probability

$$\alpha_{ij}^w := \min \left\{ 1, \frac{p(i)}{p(j)} \right\}.$$

The transition probabilities of this specific ParIS method are:

$$\begin{aligned}
q_{ij} &= \sum_w p_{ij}^w \min \left\{ 1, \frac{p_{ji}^w}{p_{ij}^w} \right\} = \sum_w \min \{ p_{ij}^w, p_{ji}^w \} = \\
&= \sum_w P(W = k) \frac{1}{n - k + 1} \cdot \\
&\quad \cdot \min \left\{ p(\langle i^m, j^{m+1}, \dots, j^{m+k-1}, i^{m+k} \rangle), p(\langle i^m, i^{m+1}, \dots, i^{m+k} \rangle) \right\}, \quad \text{if } i \neq j, \\
q_{ii} &= 1 - \sum_{k \neq i} q_{ik}.
\end{aligned}$$

3.7 Gibbs Sampler

Gibbs sampling can be applied to draw samples from a multidimensional probability distribution π , if the conditional distributions of specific coordinates given the rest of coordinates are known.

Assume that we can decompose the random variable which we want to simulate into d components: $X = (X_1, \dots, X_d)$. A d dimensional random variable is a special case of this. In Gibbs sampling one chooses a coordinate index by a certain strategy (either randomly or systematically in a well-determined order), say k , and then updates the corresponding coordinate with a sample drawn from the conditional distribution π given $X_{[-k]}$, where $X_{[-k]}$ means $(X_1, \dots, X_{k-1}, \cdot, X_{k+1}, \dots, X_d)$.

Random-Scan Gibbs Sampler

Start with an initial state i_0 , and after t iterations, being in state $i_t = (i_{t,1}, \dots, i_{t,d})$, conduct the following two steps:

1. Pick a coordinate k randomly from $\{1, \dots, d\}$ according to a fixed, strictly positive probability vector $(\alpha_1, \dots, \alpha_d)$ (e.g. $(1/d, \dots, 1/d)$).
2. Draw $i_{t+1,k}$ from the conditional distribution $\pi(\cdot \mid i_{t,[-k]})$ and leave the remaining components unchanged, that is

$$i_{t+1} = (i_{t+1,1}, \dots, i_{t+1,d}) := (i_{t,1}, \dots, i_{t,k-1}, i_{t+1,k}, i_{t,k+1}, \dots, i_{t,d}).$$

Systematic-Scan or Periodic Gibbs Sampler

Fix a permutation σ on the set $\{1, \dots, d\}$. Let \tilde{t} denote $t \bmod d$. Start with an initial state i_0 , and after t iterations, being in state $i_t = (i_{t,1}, \dots, i_{t,d})$:

- Draw $i_{t+1, \sigma(\tilde{t})}$ from the conditional distribution $\pi(\cdot \mid i_{t, [-\sigma(\tilde{t})]})$ and leave the remaining components unchanged, that is

$$i_{t+1} := (i_{t,1}, \dots, i_{t, \sigma(\tilde{t})-1}, i_{t+1, \sigma(\tilde{t})}, i_{t, \sigma(\tilde{t})+1}, \dots, i_{t,d}).$$

For details, the reader is referred to [6] or [17].

4 Theoretical results on the mixing rate

We have mentioned earlier that our main concern is the mixing rate of Markov chains, especially those that are used in Monte Carlo simulation methods. **Mixing rate** means the number of simulation steps a Markov chain needs to get sufficiently close to its equilibrium distribution. n denoting the size of a problem under investigation, we shall call a chain informally **rapidly mixing**, if this number is bounded by a polynomial in n . It is worth mentioning that the Markov chain corresponding to this problem typically has a number of states exponential in n . So mixing rapidly usually means the need for a number of steps being dramatically less, than the size of the state space itself. To have an efficient algorithm it is essential that the chain is rapidly mixing. However, one has to keep in mind that polynomial bounds with large exponents are only theoretically satisfying, and they are worthless for practical purposes.

We shall see that the mixing rate of a Markov chain on a finite state space is closely related to the second largest eigenvalue modulus of its transition matrix.

In this section we will use various sources of information, and for that reason, references will be given more precisely than before.

4.1 Measuring the distance between distributions

Since we are interested in the rate of convergence to steady state in Markov chains, we have to define what sort of convergence we are investigating. Therefore we define a distance between two distributions. There are more approaches to this.

Let I be a countable space and let α and β be two probability distributions on I .

Definition 4.1.1 The **variation distance** $d_V(\alpha, \beta)$ between α and β is

defined by

$$d_V(\alpha, \beta) := \frac{1}{2} \|\alpha - \beta\|_{\ell_1} = \frac{1}{2} \sum_{i \in I} |\alpha(i) - \beta(i)|.$$

Proposition 4.1.2 *The variation distance is a distance indeed, and it has another form:*

$$d_V(\alpha, \beta) = \sup_{S \subseteq I} |\alpha(S) - \beta(S)|.$$

PROOF This proposition is basically Lemma 1.1, in [6], Chapter 4. ■

Definition 4.1.3 The **relative pointwise distance** of α with respect to β is

$$\Delta(\alpha, \beta) := \sup_{i \in I} \frac{|\alpha(i) - \beta(i)|}{\beta(i)}.$$

For example, this definition is used in [12], with $\alpha(i) = p_{hi}^{(n)}$, maximizing also over all $h \in I$, and $\beta = \pi$, the stationary distribution.

Definition 4.1.4 The χ^2 -**distance** (or χ^2 -**contrast**) $\chi^2(\alpha, \beta)$ of α with respect to β is given by

$$\chi^2(\alpha, \beta) := \sum_{i \in I} \frac{(\alpha(i) - \beta(i))^2}{\beta(i)}.$$

χ^2 -distance is ‘stronger’ than variation distance; that means, if a Markov chain converges to the target distribution geometrically in χ^2 -distance, then it also converges geometrically in variation distance. The reason for this is that for any two probability distributions α and β on I

$$4d_V(\alpha, \beta)^2 \leq \chi^2(\alpha, \beta).$$

This can easily be proved using the Cauchy-Schwarz inequality. See [6], Chapter 6, proof of Theorem 3.2 for the short calculation.

Much more information can be found on the rich theory of probability metrics in [20]. However, in the future we shall use *variation distance* only.

If X is a random variable with values in I , let $\mathcal{L}(X)$ denote the distribution of X . If μ is a probability distribution on I , let $d_V(X, \mu)$ denote $d_V(\mathcal{L}(X), \mu)$ for short.

If the Markov chain (X_n) **converges in variation** to π , which means

$$\lim_{n \rightarrow \infty} d_V(X_n, \pi) = 0,$$

or equivalently,

$$\lim_{n \rightarrow \infty} \sum_{i \in I} |P(X_n = i) - \pi(i)| = 0,$$

then

$$\lim_{n \rightarrow \infty} E(f(X_n)) = \sum_{i \in I} f(i)\pi(i)$$

for all bounded functions $f : I \rightarrow \mathbb{R}$. The proof is straightforward. Suppose that M is an upper bound of $|f|$. Then

$$\begin{aligned} \left| E(f(X_n)) - \sum_{i \in I} f(i)\pi(i) \right| &= \left| \sum_{i \in I} f(i)(P(X_n = i) - \pi(i)) \right| \leq \\ &\leq M \sum_{i \in I} |P(X_n = i) - \pi(i)| \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

4.2 Convergence to steady state

The main qualitative result concerns ergodic Markov chains. The reader is advised to recall Corollary 3.1.20, and compare it with the following theorem.

Theorem 4.2.1 *Let P be the transition matrix of an ergodic Markov chain on the countable state space I . Let μ and ν be probability distributions on I . Then,*

$$\lim_{n \rightarrow \infty} d_V(\mu^T P^n, \nu^T P^n) = 0.$$

In particular, if ν is the stationary distribution π , then

$$\lim_{n \rightarrow \infty} d_V(\mu^T P^n, \pi^T) = 0.$$

If $\mu = \delta_j$ (δ_j is the Dirac-delta, that puts all its mass onto state j), then

$$\lim_{n \rightarrow \infty} \sum_{k \in I} |p_{jk}^{(n)} - \pi(k)| = 0.$$

PROOF The proof, which uses a coupling argument, can be found in [6], Chapter 4, proof of Theorem 2.1. ■

When the state space is finite, much more can be said about the asymptotic behaviour of Markov chains, because one can use the full machinery of linear algebra. In fact, the asymptotic behaviour of the distribution at time n depends only on the asymptotic behaviour of the n -step transition matrix P^n , and consequently, on the eigenstructure of P . Building on this, we will be able to carry out a quantitative analysis of the convergence rate.

The Perron–Frobenius theorem plays a central role in these investigations. This theorem is of linear algebra, not of Markov chain theory. Although we want to formulate it with conditions as general as possible, we have to use some notions of Markov chain theory. Before formulating the theorem, some definitions have to be introduced. We continue using Brémaud’s book [6] for this part, as well.

Definition 4.2.2 The matrix M with real coefficients is called **nonnegative** (respectively, **positive**) if all its entries are nonnegative (resp., positive). A nonnegative square matrix M is called **primitive** if there exists a positive integer k such that M^k is positive.

The definition of the communication graph is very similar to that of the transition graph (Definition 3.1.4).

Definition 4.2.3 The **communication graph** of a nonnegative square matrix $M = ((m_{ij}))_{i,j=1}^r$ is the oriented graph whose nodes are $\{1, 2, \dots, r\}$. There is an oriented edge from node i to node j if and only if $m_{ij} > 0$.

Definition 4.2.4 A nonnegative square matrix M is called **irreducible** (resp., **aperiodic**) if it has the same communication graph as the transition matrix of an irreducible (resp., aperiodic) Markov chain.

Proposition 4.2.5 *A nonnegative square matrix is primitive if and only if it is irreducible and aperiodic.*

Now we are ready to formulate the main theorem.

Theorem 4.2.6 (Perron–Frobenius theorem) *Let M be a nonnegative primitive $r \times r$ matrix. There exists a real eigenvalue μ_1 with algebraic as well as geometric multiplicity one, such that $\mu_1 > 0$ and $\mu_1 > |\mu_j|$ for any other eigenvalue μ_j . Moreover, the left eigenvector u_1 and the right eigenvector v_1 associated with μ_1 can be chosen positive and such that $u_1^T v_1 = 1$.*

Let $\mu_2, \mu_3, \dots, \mu_r$ be the eigenvalues of M other than μ_1 ordered in such a way that

$$\mu_1 > |\mu_2| \geq \dots \geq |\mu_r|,$$

and if $|\mu_2| = |\mu_j|$ for some $j \geq 3$, then $m_2 \geq m_j$, where m_j is the algebraic multiplicity of μ_j . Then

$$M^n = \mu_1^n v_1 u_1^T + \Theta(n^{m_2-1} |\mu_2|^n), \quad (1)$$

elementwise, where for $f : \mathbb{N} \rightarrow \mathbb{R}_+$, $\Theta(f(n))$ represents a function of n such that there exist $c_1, c_2 \in \mathbb{R}$, $0 < c_1 \leq c_2$, such that $c_1 f(n) \leq \Theta(f(n)) \leq c_2 f(n)$ for all n sufficiently large.

If in addition, M is stochastic, then $\mu_1 = 1$.

If M is stochastic but not irreducible, then the algebraic and geometric multiplicities of the eigenvalue 1 are equal to the number of communication classes.

If M is stochastic and irreducible with period $d > 1$, then there are exactly d distinct eigenvalues of modulus 1, namely the d th roots of unity, and all other eigenvalues have modulus strictly less than 1.

PROOF The proof can be found in Seneta's book [22]. ■

Some easy consequences for transition matrices are immediate.

Let P be an irreducible aperiodic $r \times r$ transition matrix, and let $\lambda_1, \lambda_2, \dots, \lambda_r$ denote its eigenvalues. We already know that $\lambda_1 = 1$ is among them. (The corresponding right eigenvector is the one whose all entries are 1. Let $\mathbf{1}$ denote this vector.) It is also known that with the exception of $\lambda_1 = 1$, all

eigenvalues are in the open unit disk of \mathbb{C} , and in the reversible case, they are all real¹. Therefore, with proper ordering,

$$1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_r > -1.$$

We also know from the theorem that $\lambda_r = -1$ if and only if the chain is periodic of period 2.

We will abbreviate *second largest eigenvalue* (λ_2) to **SLE**. **SLEM** will stand for *second largest eigenvalue modulus*, that is

$$\rho := \max \{ \lambda_2, |\lambda_r| \}.$$

These quantities are crucial when investigating the mixing rate of Markov chains.

Applying (1) of the Perron-Frobenius theorem for the transition matrix of an irreducible aperiodic Markov chain on a finite state space, one gets the following corollary.

Corollary 4.2.7 *If P is a transition matrix on $I = \{1, 2, \dots, r\}$ that is irreducible and aperiodic, then*

$$v_1 = \mathbf{1}, \quad u_1 = \pi,$$

where π is the unique stationary distribution. Therefore

$$P^n = \mathbf{1}\pi^T + \Theta(n^{m_2-1}\rho^n).$$

¹Suppose that P is reversible with respect to π . The operator P given by

$$(P\varphi)(i) := \sum_{j=1}^r p_{ij}\varphi(j), \quad i \in \{1, \dots, r\}$$

is a self-adjoint operator on $L_2(\pi)$. Indeed,

$$\langle P\varphi, \psi \rangle_\pi = \sum_i (P\varphi)(i)\psi(i)\pi(i) = \sum_i \sum_j p_{ij}\varphi(j)\psi(i)\pi(i),$$

by swapping i for j , and applying the *detailed balance equation*,

$$= \sum_{i,j} \varphi(i)p_{ji}\psi(j)\pi(j) = \sum_i \varphi(i) \sum_j p_{ij}\psi(j)\pi(i) = \langle \varphi, P\psi \rangle_\pi.$$

It can already be seen that convergence to equilibrium of an irreducible aperiodic finite state space Markov chain is geometric, with relative speed equal to the SLEM. The SLEM can be interpreted as the asymptotic rate of convergence to the stationary distribution. We make this statement more precise.

Theorem 4.2.8 *If P is a transition matrix on $I = \{1, 2, \dots, r\}$ that is irreducible and reversible with the stationary distribution π , then for all $n \geq 1$ and all $i \in I$,*

$$d_V(\delta_i^T P^n, \pi^T) \leq c(P, \pi) \rho^n,$$

where

$$c(P, \pi) = \min \left\{ \frac{1}{\rho} \left(\frac{p_{ii}^{(2)}}{\pi(i)} \right)^{\frac{1}{2}}, \frac{1}{2} \left(\frac{1 - \pi(i)}{\pi(i)} \right)^{\frac{1}{2}} \right\},$$

and $c(P, \pi)$ does not depend on n . Recall that

$$\delta_i^T P^n = (p_{i1}^{(n)}, p_{i2}^{(n)}, \dots, p_{ir}^{(n)}),$$

which is the distribution of the n th state starting the chain from initial state i .

PROOF This theorem is basically a combination of Theorem 3.1 and Theorem 3.3 of [6], Chapter 6. The proofs can be found there. ■

There are different quantities in the literature to characterize the mixing rate or the convergence rate of a Markov chain. We give only a few examples. A good description of these can be found in [2], Chapter 4, where they are compared to other parameters of a Markov chain (*maximal mean commute time, average hitting time* and a ‘*flow*’ parameter). Aldous and Fill prove a number of inequalities that describe the relationships between these parameters, they give illustrations of properties of chains which are closely connected to the parameters, and they present methods of bounding the parameters.

Aldous and Fill define the **variation threshold time** by

$$\tau_1 := \min \{n \in \mathbb{N} : \max_{i,j \in I} d_V(\delta_i^T P^n, \delta_j^T P^n) \leq e^{-1}\},$$

and the **relaxation time** by

$$\tau_2 := \frac{1}{1 - \lambda_2}.$$

Sinclair in [24] uses the function $\tau_i(\varepsilon)$ ($i \in I$) defined by

$$\tau_i(\varepsilon) := \min \{n_0 \in \mathbb{N} : d_V(\delta_i^T P^n, \pi^T) \leq \varepsilon \text{ for all } n \geq n_0\}$$

for $\varepsilon > 0$ and stationary distribution π .

Theorem 4.2.9 *The quantity $\tau_i(\varepsilon)$ satisfies*

1. $\tau_i(\varepsilon) \leq \frac{1}{1-\rho}(\ln \pi(i)^{-1} + \ln \varepsilon^{-1})$,
2. $\max_{i \in I} \tau_i(\varepsilon) \geq \frac{\rho}{2(1-\rho)} \ln(2\varepsilon)^{-1}$.

PROOF Part 1 follows from Proposition 3 of [7]. Part 2 is a discrete-time version of Proposition 8 of [3]. ■

Part 1 gives an upper bound on the time to reach steady state from initial state i . The converse, part 2 says that convergence cannot be rapid unless ρ is bounded away from 1.

It is interesting to note that in the latter bound there is a maximization over initial states. It is possible that a chain converges to equilibrium fast from certain states even when ρ is close to 1. However, even if such a state exists, finding it requires more detailed information about the chain than is usually available in more complex applications.

In our applications we shall use $\tau_i(\varepsilon)$ in Section 5.6 and the SLE (consequently, τ_2) in every other example.

There are two more issues we discuss in this section. First, we explain how theorems like Theorem 4.2.8 can inspire the definition of relaxation time τ_2 . The second is that we investigate under what conditions the SLEM can be exchanged for SLE in formulae, for example, why we defined relaxation time by λ_2 instead of ρ .

Definition 4.2.10 The value $1 - \rho$ is called the **spectral gap**.

Suppose that we have an infinite sequence of problems that are very similar but have different, strictly increasing sizes. As we have already defined, n denoting the size of a problem in this sequence, we shall call the sequence of Markov chains corresponding to the sequence of these problems *rapidly mixing*, if the number of simulation steps the appropriate Markov chain needs to get sufficiently close to its equilibrium distribution is bounded by a polynomial in n , with one fixed polynomial for the entire sequence.

We present a typical setting in which the sequence of Markov chains is rapidly mixing.

Let $\rho(n)$ denote the SLEM of the Markov chain corresponding to the unique problem of the sequence of size n . If the spectral gap corresponding to $\rho(n)$ is $1/q(n)$, where $q(x) = a_k x^k + a_{k-1} x^{k-1} + \dots + a_0$ ($a_k > 0$) is a fixed polynomial, then

$$\lim_{n \rightarrow \infty} \rho(n)^{q(n)} = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{q(n)} \right)^{q(n)} = e^{-1} > 0,$$

but for any $\varepsilon > 0$, defining $q'(x) := x^{k+\varepsilon}$,

$$\lim_{n \rightarrow \infty} \rho(n)^{q'(n)} = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{q(n)} \right)^{q'(n)} = 0.$$

This means that if spectral gaps of SLEMs $\rho(n)$ have lower bounds of the form $1/q(n)$ for a polynomial $q(x) = a_k x^k + \dots + a_0$ ($a_k > 0$), then Theorem 4.2.8 yields that the sequence is rapidly mixing.

It is a little nuisance that one must consider the SLEM when investigating the mixing rate. It would be easier if one could forget about λ_r , and focus attention on the SLE only. Moreover, there are better tools to estimate the SLE than λ_r , as we shall see. We present a crude approach to overcome this problem.

Suppose that one adds a holding probability of $1/2$ to each state, in other words, one replaces the transition matrix P with $P' := \frac{1}{2}(P + I)$, where I is

the $r \times r$ identity matrix. P' is irreducible and aperiodic, just as P . It can easily be verified that π remains the unique stationary distribution of P' . Eigenvalues of P' are

$$\lambda'_i = \frac{\lambda_i + 1}{2}, \text{ for all } i \in \{1, 2, \dots, r\}.$$

This ensures that all eigenvalues are positive, while the spectral gap decreases only by a factor of 2. Consequently, one can use the new Markov chain for MCMC simulation and the mixing rate remains not much worse than the original was.

In Section 4.6 we shall prove that in the reversible case, if for all $i \in \{1, 2, \dots, r\}$ $p_{ii} \geq 1/2$, then all eigenvalues are nonnegative. In our applications when using variants of the Metropolis–Hastings algorithm, we will have symmetric transition matrices with uniform distribution as the target distribution. This guarantees reversibility. The latter condition will also be met in most of our examples and hence we will have $\rho = \lambda_2$.

In applications it is typically very hard, practically impossible, to calculate the SLEM explicitly. Therefore other methods are needed to investigate the mixing rate. In the coming sections we are to collect the most important representatives of these methods.

4.3 Gershgorin's bound

Gershgorin's bound is a surprising theorem of linear algebra that is sometimes applicable in Markov chain theory. The proof can be found in books on numerical analysis or in [6].

Theorem 4.3.1 *Let A be an $r \times r$ matrix with complex entries. Let us introduce the following notations: for $i, j \in \{1, 2, \dots, r\}$*

$$r_i := \sum_{j \neq i} |a_{ij}|, \quad s_j := \sum_{i \neq j} |a_{ij}|.$$

Then for any eigenvalue λ , there exists a $k \in \{1, 2, \dots, r\}$ such that

$$|\lambda - a_{kk}| \leq r_k,$$

and there exists a $k' \in \{1, 2, \dots, r\}$ such that

$$|\lambda - a_{k'k'}| \leq s_{k'}.$$

In other words, all eigenvalues can be located within the union of r discs on the complex plane with centres a_{kk} .

Corollary 4.3.2 *Let P be an $r \times r$ stochastic matrix. Then, for all eigenvalues λ , there exists a $k \in \{1, 2, \dots, r\}$ such that*

$$|\lambda - p_{kk}| \leq 1 - p_{kk}. \quad (2)$$

When the eigenvalue λ is real, then

$$-1 + 2 \min_i p_{ii} \leq \lambda. \quad (3)$$

PROOF (2) follows from the fact that for a stochastic matrix $r_k = 1 - p_{kk}$.

If λ is real, then either $\lambda - p_{ii} < 0$ for all i or there exists an i such that $\lambda - p_{ii} \geq 0$.

In the first case, $|\lambda - p_{ii}| = p_{ii} - \lambda$ for all i . By (2)

$$p_{kk} - \lambda \leq 1 - p_{kk}$$

for some k , which proves (3):

$$-1 + 2 \min_i p_{ii} \leq -1 + 2p_{kk} \leq \lambda.$$

If $\lambda - p_{ii} \geq 0$ for some i , then $\lambda - \min_i p_{ii} \geq 0$. Using this inequality and $-1 + \min_i p_{ii} \leq 0$ we get (3). ■

Unfortunately, Theorem 4.3.1 cannot be used to prove that a Markov chain Monte Carlo is rapidly mixing because it cannot bound the SLE away from 1. The largest eigenvalue $\lambda_1 = 1$ is also contained in a disc, whose centre (p_{ii} for some i) is a real number, moreover $p_{ii} \in [0, 1]$. The radius of this disc must be positive, otherwise the chain would not be irreducible: the chain could not leave or could not enter state i (if, indirectly, we used the bound

with r_k or $s_{k'}$, respectively). If this radius is positive, then this disc contains real numbers arbitrarily close to 1, and each of them is a potential SLE.

Corollary 4.3.2 can be used to refute rapid mixing. Let $\varepsilon := 1/2^r$, $p_{ii} := 1 - \varepsilon$ for all $i \in \{1, \dots, r\}$ and $p_{ij} := \varepsilon/(r - 1)$, if $i \neq j$. It is obvious that this chain converges in variation to the uniform distribution, and it can be seen that it does this very slowly. Indeed, Corollary 4.3.2 proves $\lambda_2 \geq 1 - 2\varepsilon = 1 - 2/2^r$.

4.4 Liu's result for MIS

Liu [16] managed to identify the eigenvalues and corresponding right eigenvectors of the transition matrix in the very special case of Metropolised Independent Sampling. Recall from Section 3.5 that for $i \in I$, the importance weight is defined by $w(i) := \pi(i)/p(i)$. Without loss of generality, we assume that $I = \{1, 2, \dots, r\}$ and the states are labelled according to the magnitudes of their importance weights:

$$w(1) \geq w(2) \geq \dots \geq w(r).$$

In this case, the eigenvalues are $1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_r \geq 0$, where for $k \geq 2$

$$\lambda_k = \sum_{i=k}^r \left(p(i) - \frac{\pi(i)}{w(k-1)} \right) = \sum_{i=k}^r \pi(i) \left(\frac{1}{w(i)} - \frac{1}{w(k-1)} \right).$$

It is easy to check that the SLEM is

$$\rho = \lambda_2 = \sum_{i=2}^r \left(p(i) - \frac{\pi(i)}{w(1)} \right) = \sum_{i=1}^r \left(p(i) - \frac{\pi(i)}{w(1)} \right) = 1 - \frac{1}{w(1)}.$$

4.5 Conductance

The first analyses of the efficiency of more complex Markov chains arising in combinatorial applications were carried out using a quantity called *conductance*. The conductance Φ is essentially the *edge expansion* of the transition graph. Φ may be viewed as the probability that the chain in steady state escapes from a subset $S \subseteq I$ in one step, minimized over all subsets S . It

is intuitively reasonable that Φ is related to the mixing rate: if this escape probability is small for some S then the cut edges separating S from $I \setminus S$ form a bottleneck, which prevents rapid convergence to steady state. Conversely, if Φ has a large value, then the chain cannot be trapped by any small region of the state space, and hence it should mix rapidly. After we gave some explanation on what was about to happen, we now introduce the definition of conductance.

For a nonempty set $S \subseteq I = \{1, 2, \dots, r\}$ the **capacity** of S is

$$\pi(S) := \sum_{i \in S} \pi(i),$$

and the **ergodic flow** out of S is

$$F(S) := \sum_{i \in S, j \in I \setminus S} \pi(i)p_{ij}.$$

Note that $0 \leq F(S) \leq \pi(S) \leq 1$.

Definition 4.5.1 One defines the **conductance** of a pair (P, π) by

$$\Phi := \min \left\{ \frac{F(S)}{\pi(S)} \mid S \subset I, 0 < \pi(S) \leq \frac{1}{2} \right\}.$$

If the SLEM is the SLE, then Φ characterizes the rapid mixing property by giving an upper and also a lower bound for the SLE.

Theorem 4.5.2 (Cheeger's inequality) *The SLE λ_2 satisfies*

$$1 - 2\Phi \leq \lambda_2 \leq 1 - \frac{\Phi^2}{2}.$$

PROOF The proof can be found in [6], Chapter 6, Section 4.2, or in [7]. The latter gives historical remarks and references where refinements can be found. ■

We see that a Markov chain is rapidly mixing if and only if $\Phi \geq 1/p(n)$ for some polynomial p with a positive leading coefficient.

4.6 Canonical (distinguished) paths, Poincaré coefficient

If one wants to prove positive propositions (i.e. a certain chain is rapidly mixing), one has to give good lower bounds on Φ . In some cases such bounds can be obtained directly, using elementary arguments or geometric ideas. (For such investigations, see the references given by Sinclair in [24].) However, this is not typical.

A strong technique was introduced by Jerrum and Sinclair in [12] (and in Sinclair's PhD thesis). The idea is to choose a *canonical path* γ_{ij} in the transition graph for each ordered pair $(i, j) \in I \times I$ ($i \neq j$). If the paths can be chosen in such a way that no edge is overloaded by paths, then the chain cannot contain a bottleneck, so Φ is not too small. (The existence of a constriction would imply that for any choice of canonical paths, the edges in the bottleneck must be overloaded.) Some applications of this technique can be found in their paper.

They use an interesting approach to bounding the number of canonical paths that use a given edge. Instead of counting these directly, they set up an injective mapping for each oriented edge, from the set of canonical paths containing this edge into the state space. They call it an *encoding* technique, since the paths (or pairs of endpoints of paths) are encoded by single states.

Although being powerful, using this technique is not easy, since one always has to start from scratch. That means one has to prove that a good choice of canonical paths means a good lower bound for the conductance.

Diaconis and Stroock [7] found a 'user-friendly interface' to use a similar argument. They observed that canonical path arguments can lead directly to bounds on the relaxation time, independently of the conductance Φ . First, we present a theorem they used to prove their result.

The **variance** of x with respect to π is given by

$$\text{Var}_\pi(x) := \langle x, x \rangle_\pi - \langle x, \mathbf{1} \rangle_\pi^2.$$

Recall that the inner product $\langle \cdot, \cdot \rangle_\pi$ was defined in Section 3.1, Definition 3.1.23.

Let I denote the $r \times r$ identity matrix. (Admittedly, we create some confusion by denoting both the state space and the identity matrix by I . However, we are sure that the reader can always tell which is the one under consideration.) The **Dirichlet form** associated with a reversible pair (P, π) is defined by

$$\mathcal{E}_\pi(x, x) := \langle (I - P)x, x \rangle_\pi,$$

and it equals

$$\mathcal{E}_\pi(x, x) = \frac{1}{2} \sum_{i,j \in I} (x(i) - x(j))^2 \pi(i) p_{ij}.$$

Rayleigh's theorem gives a characterization of the eigenvalues of P using the notions of variance and Dirichlet form, given the right eigenvectors v_i ($i \in I = \{1, 2, \dots, r\}$) associated with eigenvalues λ_i , respectively.

Theorem 4.6.1 (Rayleigh's theorem) *Let P be the transition matrix of an irreducible Markov chain on a finite state space, and let π be its stationary distribution. Assume that P is reversible with respect to π . Then, for $i \geq 2$*

$$1 - \lambda_i = \inf \left\{ \frac{\mathcal{E}_\pi(x, x)}{\text{Var}_\pi(x)} \mid x \neq 0, \text{ for all } j \in \{1, \dots, i-1\} \langle x, v_j \rangle_\pi = 0 \right\}.$$

Any vector achieving the infimum above is an eigenvector of P corresponding to the eigenvalue λ_i .

Specifically, since $v_1 = \mathbf{1}$,

$$1 - \lambda_2 = \inf \left\{ \frac{\mathcal{E}_\pi(x, x)}{\text{Var}_\pi(x)} \mid x \neq 0, \langle x, \mathbf{1} \rangle_\pi = 0 \right\},$$

and the fact that $\mathcal{E}_\pi(x, x) = \mathcal{E}_\pi(x - c\mathbf{1}, x - c\mathbf{1})$ for any $c \in \mathbb{R}$, implies that this is equivalent to

$$1 - \lambda_2 = \inf \left\{ \frac{\mathcal{E}_\pi(x, x)}{\text{Var}_\pi(x)} \mid x \text{ is nonconstant} \right\}.$$

PROOF The proof of this theorem can be found in [6], Chapter 6, Section 2.2. ■

We are now ready to proceed to formulate Diaconis and Stroock's results, usually referred as *geometric bounds*.

In the transition graph associated with transition matrix P on the state space $I = \{1, 2, \dots, r\}$, we shall denote an oriented edge $i \rightarrow j$ by e . Define for any such oriented edge e ,

$$Q(e) := \pi(i)p_{ij}.$$

For each ordered pair of distinct states $(i, j) \in I \times I$, select arbitrarily one and only one path from i to j , that is a sequence i, i_1, \dots, i_m, j such that $p_{ii_1}, p_{i_1i_2}, \dots, p_{i_mj} > 0$, which does not use the same edge twice. Let Γ be the set of paths so selected. For a path $\gamma_{ij} \in \Gamma$, let

$$|\gamma_{ij}|_Q := \sum_{e \in \gamma_{ij}} \frac{1}{Q(e)} = \frac{1}{\pi(i)p_{ii_1}} + \frac{1}{\pi(i_1)p_{i_1i_2}} + \dots + \frac{1}{\pi(i_m)p_{i_mj}}.$$

Definition 4.6.2 The **Poincaré coefficient** is defined by

$$\kappa := \kappa(\Gamma) := \max_e \sum_{\gamma_{ij} \ni e} |\gamma_{ij}|_Q \pi(i) \pi(j).$$

(The sum goes for all paths that contain a fixed edge e .)

Theorem 4.6.3 (Diaconis, Stroock) *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Denoting by λ_2 its SLE,*

$$\lambda_2 \leq 1 - \frac{1}{\kappa}.$$

We can also give a lower bound of the smallest eigenvalue. For each state $i \in I$, select exactly one closed path σ_i from i to i that does not pass twice through the same edge, and it consists of an *odd* number of edges (for this to be possible, we assume that P is aperiodic). Let Σ be the collection of paths so selected. For a path $\sigma_i \in \Sigma$, let

$$|\sigma_i|_Q := \sum_{e \in \sigma_i} \frac{1}{Q(e)}.$$

Define

$$\alpha := \alpha(\Sigma) := \max_e \sum_{\sigma_i \ni e} |\sigma_i|_Q \pi(i).$$

Theorem 4.6.4 (Diaconis, Stroock) *Let P be an irreducible and aperiodic transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Then*

$$\lambda_r \geq -1 + \frac{2}{\alpha}.$$

At the end of Section 4.2 we announced proving that in a reversible Markov chain if for all $i \in \{1, 2, \dots, r\}$ $p_{ii} \geq 1/2$, then all eigenvalues are non-negative, hence the SLEM and the SLE are identical. We can present this implication as an easy corollary of the last theorem.

Proposition 4.6.5 *If an irreducible Markov chain with state space $I = \{1, 2, \dots, r\}$ is reversible with respect to π and has the property that for all $i \in I$ $p_{ii} \geq 1/2$, then all eigenvalues of its transition matrix are nonnegative.*

PROOF We shall use the last theorem of which the conditions are satisfied. For all $i \in I$ we choose the single $i \rightarrow i$ edge as path σ_i . Then

$$\alpha = \max_e \sum_{\sigma_i \ni e} \frac{1}{\pi(i)p_{ii}} \pi(i) = \max_{i \in I} \frac{1}{p_{ii}} = \frac{1}{\min_{i \in I} p_{ii}} \leq 2,$$

and consequently

$$\lambda_r \geq -1 + \frac{2}{\alpha} \geq -1 + 1 = 0. \quad \blacksquare$$

4.7 Similar bounds with different coefficients

We shall investigate a modification of Theorem 4.6.3 that Sinclair developed and presented in [24].

He also gave a short formula how the canonical path approach can give a bound of the conductance. We use the notation Γ of the previous section. The *maximum loading* of any edge is measured by the quantity

$$\vartheta := \vartheta(\Gamma) := \max_e \frac{1}{Q(e)} \sum_{\gamma_{ij} \ni e} \pi(i)\pi(j).$$

Sinclair says that we may view the Markov chain as a flow network, in which $\pi(i)\pi(j)$ units of flow travel from i to j along γ_{ij} , and $Q(e)$ plays the role of the capacity of e . The quantity ϑ measures the maximum flow along any edge as a fraction of its capacity. Note that this picture is somewhat misleading: it may happen that $\vartheta > 1$, that is, the flow along the edge exceeds the capacity.

The following result shows that the existence of a good choice of paths implies a large value for the conductance.

Proposition 4.7.1 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . With any choice of canonical paths,*

$$\Phi \geq \frac{1}{2\vartheta}.$$

This can be combined with Cheeger's inequality, Theorem 4.5.2.

Corollary 4.7.2 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . With any choice of canonical paths, the SLE satisfies*

$$\lambda_2 \leq 1 - \frac{1}{8\vartheta^2}.$$

This bound may be rather weak because of the square in the denominator. Obviously, this is not a problem for someone trying to prove polynomial mixing time, only for those who try to get sharp estimates.

We turn to Sinclair's modified geometric bound. We replace κ with

$$K := K(\Gamma) := \max_e \frac{1}{Q(e)} \sum_{\gamma_{ij} \ni e} |\gamma_{ij}| \pi(i)\pi(j),$$

where $|\gamma_{ij}|$ is the length (i.e. the number of edges) of the path γ_{ij} .

Theorem 4.7.3 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Then*

$$\lambda_2 \leq 1 - \frac{1}{K}.$$

The following simplified form is an easy consequence that is often useful.

Corollary 4.7.4 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Let*

$$\ell := \ell(\Gamma) := \max_{\gamma \in \Gamma} |\gamma|,$$

that is the length of a longest path in Γ . With any choice of canonical paths, the SLE satisfies

$$\lambda_2 \leq 1 - \frac{1}{\ell\vartheta}.$$

This corollary can be applied in the same situations as Corollary 4.7.2. However, the maximum path length ℓ is often less than the estimate obtained for ϑ . In such cases, Corollary 4.7.4 will give a sharper bound than Corollary 4.7.2. Sinclair presents examples that confirm this observation.

Sinclair also argues that K and $\ell\vartheta$ are usually more useful quantities to work with in practice than Diaconis and Stroock's κ .

When one tries to apply these theorems to a certain Markov chain with a more complicated transition graph, one might encounter the problem, that any definition of paths will result in Γ having a very symmetrical structure. This might mean that the majority of edges in the transition graph are hardly used in any paths, but some edges are really overloaded. Because of the maximization over edges in the formulae, any overload results in a large $\kappa(\Gamma)$ or $K(\Gamma)$, and a useless upper bound of the SLE. Although there may exist a good choice of paths, if we cannot find it, we only get weak upper bounds.

One may have the fuzzy idea why not pick paths at random to avoid choosing a Γ with a symmetrical structure. In fact, there are results in this spirit. We continue the investigation of bounding the relaxation time of Markov chains with these results.

4.8 Multicommodity flow

The following theorems about multicommodity flows are due to Sinclair, and they can be found in the same paper [24]. These can be seen as natural generalizations of the path-counting ideas of the last sections.

Let us view the transition graph G as a flow network by assigning the capacity $Q(e) = \pi(i)p_{ij}$ to each oriented edge e whose initial and end vertices are i and j . Suppose that for each ordered pair of distinct vertices i and j , a quantity $\pi(i)\pi(j)$ of some commodity denoted by (i, j) is to be transported from i to j along the edges of the network. The task is to construct a flow which minimizes the total throughput through any oriented edge e in G as a fraction of its capacity $Q(e)$. Recall that the flow through an edge may exceed its capacity. This viewpoint is already familiar from the previous section. But now we are allowing multiple rather than canonical paths between vertices. Fortunately we can get very similar bounds and proofs can be carried over quite easily.

Definition 4.8.1 Let Π_{ij} be the set of all directed paths from i to j in the transition graph G . Let $\Pi := \bigcup_{i \neq j} \Pi_{ij}$. A **flow** in G is a function $f : \Pi \rightarrow [0, \infty[$ which satisfies

$$\sum_{\gamma \in \Pi_{ij}} f(\gamma) = \pi(i)\pi(j)$$

for all $i, j \in I$, $i \neq j$.

f can be extended to a function on all oriented edges by setting

$$f(e) := \sum_{\gamma \ni e} f(\gamma),$$

that is, $f(e)$ is the total flow routed by f through e .

We introduce the quantity $\vartheta(f)$ analogous to ϑ :

$$\vartheta(f) := \max_e \frac{f(e)}{Q(e)}.$$

As we have suggested, Proposition 4.7.1 and Corollary 4.7.2 carry over immediately to this more general setting.

Proposition 4.8.2 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . With any flow f ,*

$$\Phi \geq \frac{1}{2\vartheta(f)}.$$

Corollary 4.8.3 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . With any flow f , the SLE satisfies*

$$\lambda_2 \leq 1 - \frac{1}{8\vartheta(f)^2}.$$

We can also generalize the quantity K to a flow f . Let a function \bar{f} be defined on all oriented edges by

$$\bar{f}(e) := \sum_{\gamma \ni e} f(\gamma)|\gamma|,$$

where $|\gamma|$ is the length of the path γ . (Sinclair uses the expression *elongated flow* through e for $\bar{f}(e)$.) Let

$$K(f) := \max_e \frac{\bar{f}(e)}{Q(e)}.$$

Again, Theorem 4.7.3 and Corollary 4.7.4 carry over, and one can get the following.

Theorem 4.8.4 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Then, with any flow f ,*

$$\lambda_2 \leq 1 - \frac{1}{K(f)}.$$

Corollary 4.8.5 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Let*

$$\ell(f) := \max_{\gamma \in \Pi, f(\gamma) > 0} |\gamma|.$$

With any flow f , the SLE satisfies

$$\lambda_2 \leq 1 - \frac{1}{\ell(f)\vartheta(f)}.$$

There are examples in [24] for the use of multicommodity flows and they show that the flexibility provided by flows can yield better bounds than canonical path approaches.

4.9 Probabilistic inequalities

Aldous and Fill ([2], Chapter 4) proved a probabilistic version of inequalities that we discussed in previous sections.

For each ordered pair of distinct states $(i, j) \in I \times I$, let γ_{ij} be a random path from i to j of the form $i = I_0, I_1, \dots, I_M = j$ of random length $M = |\gamma_{ij}|$, such that no edge is traversed more than once.

Theorem 4.9.1 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Then*

$$\lambda_2 \leq 1 - \left(\max_e \frac{1}{Q(e)} \sum_{i \in I} \sum_{j \in I} \pi(i)\pi(j) E\left(|\gamma_{ij}| \chi_{\{e \in \gamma_{ij}\}}\right) \right)^{-1}.$$

Schweinsberg [21] proved a corollary of this theorem, which can be useful if one can easily define a short path γ_{ij} only when j is in some subset $B \subseteq I$. He used this result to study a Markov chain on the set of n -leaf cladograms.

Corollary 4.9.2 *Let P be an irreducible transition matrix on the finite state space I , with stationary distribution π , and assume that P is reversible with respect to π . Let B be a subset of I . Suppose, for all $i \in I$ and $j \in B$, that γ_{ij} is a path from i to j , possibly random, which has at most L edges. Then*

$$\lambda_2 \leq 1 - \left(\frac{4L}{\pi(B)} \max_e \frac{1}{Q(e)} \sum_{i \in I} \sum_{j \in B} \pi(i)\pi(j) P(e \in \gamma_{ij}) \right)^{-1}.$$

4.10 Dobrushin's inequality

We already know that the SLEM determines the rate of convergence to steady state. Various methods were presented in previous sections for obtaining upper bounds of the SLEM. For theoretical purposes, there is another upper bound for the SLEM due to R. L. Dobrushin, that is often very tractable. If the state space is finite and the chain is ergodic, it guarantees a computable geometric rate of convergence to equilibrium.

In this section we follow Brémaud's book [6]. The reader is directed to this book for references.

Let G, H, I denote countable sets.

Definition 4.10.1 Let P be a stochastic matrix indexed by $H \times G$. Its **Dobrushin's ergodic coefficient** $\delta(P)$ is defined by

$$\delta(P) := \frac{1}{2} \sup_{i,j \in H} \sum_{k \in G} |p_{ik} - p_{jk}| = \sup_{i,j \in H} d_V(p_i, p_j).$$

First, note that $0 \leq \delta(P) \leq 1$. If the stochastic matrix P has two orthogonal rows (e.g. $i, j \in H$, such that for all $k \in G$ $p_{ik}p_{jk} = 0$), then $\delta(P) = 1$, which is totally useless as we shall see shortly. This happens quite often: when H is infinite, it is typical that P has orthogonal rows. However, for finite state spaces this notion is usually very powerful.

Theorem 4.10.2 (Dobrushin's inequality) Let $P_1 = ((a_{ij}))$ and $P_2 = ((b_{ij}))$ denote two stochastic matrices indexed by $H \times G$ and $I \times H$, respectively. Let $P_2P_1 = ((c_{ij}))$ be their product indexed by $I \times G$ just as one would expect:

$$c_{ij} := \sum_{k \in H} b_{ik}a_{kj}.$$

Then

$$\delta(P_2P_1) \leq \delta(P_2)\delta(P_1).$$

One can get a geometric convergence rate in terms of Dobrushin's coefficient.

Theorem 4.10.3 *Let P be a stochastic matrix indexed by I . Let μ and ν be probability distributions on I . Then*

$$d_V(\mu^T P^n, \nu^T P^n) \leq d_V(\mu^T, \nu^T) \delta(P)^n.$$

Proofs of both theorems can be found in [6], Chapter 6, Section 7.1.

4.11 Convergence rate for the Gibbs Sampler

One can find in [6] how Theorem 4.10.3 can prove a geometric rate of convergence of the Periodic Gibbs Sampler (for the latter, see Section 3.7). Recall that the random variable was demanded to be decomposable, namely, into d components.

Theorem 4.11.1 *Let P be the transition matrix of a Periodic Gibbs Sampler. Let π denote its stationary distribution and let μ be a probability distribution on its state space. Then*

$$d_V(\mu^T P^n, \pi^T) \leq \frac{1}{2} d_V(\mu^T, \pi^T) (1 - e^{-d\Delta})^n,$$

where Δ depends on π and can be given explicitly.

It is more surprising that the distance in variation to the target distribution decreases in every step.

Theorem 4.11.2 *Let μ be a probability distribution on the state space and let ν be the probability distribution obtained by applying a Gibbs step at an arbitrary component. Then $d_V(\nu^T, \pi^T) \leq d_V(\mu^T, \pi^T)$.*

Brémaud's textbook claims that it is an experimental fact that the Gibbs Sampler is not the best MCMC simulation algorithm. Brémaud suggests that the monotonicity property is not an advantage but a possible indicator of short-sighted strategies.

Jun S. Liu's thesis may have more information on this topic (*Correlation Structure and Convergence Rate of the Gibbs Sampler*, Ph.D. thesis, University of Chicago, 1991).

4.12 The coupling method

Coupling is an old idea of W. Doeblin from 1938 that was revived in Markov chain theory in the 1970s.

Recall from Section 4.1 that if X is a random variable with values in I , $\mathcal{L}(X)$ denotes the distribution of X . $d_V(X_1, X_2)$ denotes $d_V(\mathcal{L}(X_1), \mathcal{L}(X_2))$ for short.

In Section 4.1 one can also find that the definition of *convergence in variation* concerns only the marginal distributions of the Markov chain, not the process itself. Therefore, if there exist two Markov chains (X_n) and (X'_n) with $\mathcal{L}(X_n) = \mathcal{L}(X'_n)$ for all $n \in \mathbb{N}$, and there exists a third one (X''_n) , such that $\mathcal{L}(X''_n) = \pi$ for all $n \in \mathbb{N}$, then

$$\lim_{n \rightarrow \infty} d_V(X'_n, X''_n) = 0$$

is sufficient to yield

$$\lim_{n \rightarrow \infty} \sum_{i \in I} |P(X_n = i) - \pi(i)| = 0,$$

or in other words, to provide the convergence of (X_n) in variation to π .

This trivial observation is useful because of the freedom in the choice of (X'_n) and (X''_n) . The most interesting case is when one uses dependent versions and there exists a finite random time τ such that $X'_n = X''_n$ for all $n \geq \tau$.

τ is finite if and only if $\lim_{n \rightarrow \infty} P(\tau > n) = 0$. In this case, the *coupling inequality* implies that (X_n) converges in variation to π .

Definition 4.12.1 Let (X'_n) and (X''_n) be two stochastic processes that take their values in the same countable state space. They are said to **couple**, if there exists an almost surely finite random time τ such that if $n \geq \tau$, then $X'_n = X''_n$. The random variable τ is called a **coupling time** of the two processes.

Theorem 4.12.2 (Coupling inequality) *If τ is a coupling time for (X'_n) and (X''_n) , then*

$$d_V(X'_n, X''_n) \leq P(\tau > n).$$

PROOF The proof can be found in [6], Chapter 4, Section 1. ■

Combining previous observations, one can prove the convergence of (X_n) in variation to π by giving appropriate Markov chains (X'_n) and (X''_n) and proving that they couple. This method can be carried out in many situations and there are many applications of it.

For historical comments, typical applications and additional information, see Lindvall's textbook [15].

In the following section (Section 5, Applications), a number of results described in Section 4 are applied to our examples. The Applications section can be seen as a collection of examples of applications of these ideas.

Liu's theorem of Section 4.4 is used in Sections 5.1, 5.2 and 5.3. *Canonical paths* arguments of Section 4.6 are used in Section 5.1 and 5.3. An argument based on the notion *conductance* of Section 4.5 is used in Section 5.4. We shall use a *coupling argument* (Section 4.12) to prove the bound of Section 5.6.

5 Applications

The purpose of this section is twofold. We present results of our investigations into applying MCMC methods in two fields of bioinformatics. Secondly, these examples demonstrate the use of techniques of previous sections.

5.1 Ladder

We present an example in which the ParIS is superior to MIS in terms of rate of convergence to the target distribution: the ParIS is rapidly mixing, the MIS is not.

As we have already made clear, our main subject is about random walks on directed acyclic graphs. We introduce this example as a Markov chain whose states correspond to paths. Let us take a graph of a form of a ladder (see Figure 2). We investigate paths from the top left to the bottom right corner that consist of steps *right* and *down*, therefore once crossed from the left to the right, cannot go back to the left side. The states of the Markov chain will be these paths. Obviously, there is a one-to-one correspondence between these paths and rungs of the ladder: each path contains exactly one rung. Suppose that the ladder has $n \geq 2$ rungs, and the state space is $I = \{1, 2, \dots, n\}$.

We choose the uniform distribution on I to be the target distribution.

First, the MIS. We define the primal Markov chain, whose transition matrix is P , with the *random walk procedure* defined in Section 3.6. Let us recall this method. Imagine that we are at the top left corner. We can choose between two directions to start with: *right* or *down*. We pick one of them with probability $1/2 - 1/2$. If we picked *right*, then the path is determined: we cross using the first rung. If we picked *down*, then move to the next node and choose between the two directions again with the same probability distribution. Go on with the same procedure. If we arrive at the bottom left corner, we use the last rung. It is obvious that using this distribution the

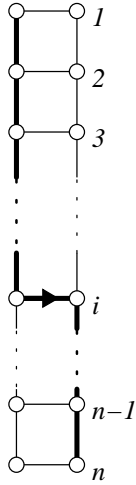


Figure 2: Ladder

probability of path i is

$$\left(\frac{1}{2}\right)^{i \wedge (n-1)}.$$

This is a natural approach from the viewpoint of somebody staying at the top left corner to the problem of defining a method to choose from the paths uniformly at random, without knowing anything about either the special structure or the number of paths. Of course we know everything about the paths and we also know that this distribution is far from being uniform. Remember that the states are not the nodes of the ladder, but the paths. As with every MIS method, this Markov chain also has the property that the probability distribution of its state at time $t + 1$ is independent from its state at time t : for each $i, j \in I$

$$p_{ij} = \left(\frac{1}{2}\right)^{j \wedge (n-1)}.$$

Proposition 5.1.1 *For the SLE of the transition matrix of this MIS method*

$$\lambda_2 \geq 1 - \frac{n}{2^{n-1}},$$

hence this MIS method has exponential relaxation time.

PROOF We start with calculating the acceptance probabilities α_{ij} ($i, j \in I$, $i \neq j$):

$$\alpha_{ij} = \min \left\{ 1, \frac{\frac{1}{n} \left(\frac{1}{2}\right)^{i \wedge (n-1)}}{\frac{1}{n} \left(\frac{1}{2}\right)^{j \wedge (n-1)}} \right\} = \begin{cases} 1, & \text{if } i < j, \\ 1/2^{i-j}, & \text{if } j < i \leq n-1, \\ 1/2^{n-1-j}, & \text{if } j < i = n. \end{cases}$$

We get the following transition matrix Q ($i \neq j$):

$$q_{ij} = p_{ij} \alpha_{ij} = \begin{cases} 1/2^j, & \text{if } i < j \leq n-1, \\ 1/2^i, & \text{if } j < i \leq n-1, \\ 1/2^{n-1}, & \text{if } j < i = n, \\ 1/2^{n-1}, & \text{if } i < j = n, \end{cases}$$

and with more calculation

$$q_{ii} = p_{ii} + \sum_{j \neq i} p_{ij} (1 - \alpha_{ij}) = 1 - \frac{i \wedge (n-1)}{2^{i \wedge (n-1)}}.$$

$$Q = \begin{pmatrix} 1 - \frac{1}{2} & 1/4 & 1/8 & \dots & 1/2^i & \dots & 1/2^{n-1} & 1/2^{n-1} \\ 1/4 & 1 - \frac{2}{4} & 1/8 & & 1/2^i & & 1/2^{n-1} & 1/2^{n-1} \\ 1/8 & 1/8 & 1 - \frac{3}{8} & & 1/2^i & & 1/2^{n-1} & 1/2^{n-1} \\ \vdots & & & \ddots & & & \vdots & \vdots \\ 1/2^i & 1/2^i & 1/2^i & & 1 - \frac{i}{2^i} & & 1/2^{n-1} & 1/2^{n-1} \\ \vdots & & & & & \ddots & & \vdots \\ 1/2^{n-1} & 1/2^{n-1} & 1/2^{n-1} & \dots & 1/2^{n-1} & & 1 - \frac{n-1}{2^{n-1}} & 1/2^{n-1} \\ 1/2^{n-1} & 1/2^{n-1} & 1/2^{n-1} & \dots & 1/2^{n-1} & \dots & 1/2^{n-1} & 1 - \frac{n-1}{2^{n-1}} \end{pmatrix}$$

Since this is our first example, we shall present calculations in detail. As we shall see, the structure of the transition matrix makes it possible to estimate the SLEM directly: we will not use theoretical results on Markov chains, only elementary linear algebra.

We give a lower bound of the second largest eigenvalue of Q , λ_2 . Let I denote the identity matrix. Eigenvalues of Q are exactly the roots of $\det(Q - xI)$. If we subtract the last line of $Q - xI$ from the preceding line componentwise,

we get

$$\begin{pmatrix} 1 - \frac{1}{2} - x & 1/4 & \dots & 1/2^i & \dots & 1/2^{n-1} & 1/2^{n-1} \\ 1/4 & 1 - \frac{2}{4} - x & & 1/2^i & & 1/2^{n-1} & 1/2^{n-1} \\ \vdots & & \ddots & & & \vdots & \vdots \\ 1/2^i & 1/2^i & & 1 - \frac{i}{2^i} - x & & 1/2^{n-1} & 1/2^{n-1} \\ \vdots & & & & \ddots & & \vdots \\ 0 & 0 & \dots & 0 & & 1 - \frac{n}{2^{n-1}} - x & -1 + \frac{n}{2^{n-1}} + x \\ 1/2^{n-1} & 1/2^{n-1} & \dots & 1/2^{n-1} & \dots & 1/2^{n-1} & 1 - \frac{n-1}{2^{n-1}} - x \end{pmatrix}$$

and the new determinant equals $\det(Q - xI)$. One can see that choosing $x = 1 - n/2^{n-1}$ each entry of the line before the last one is zero. It means that $1 - n/2^{n-1}$ is a root of the characteristic polynomial, it is an eigenvalue, and this quantity is a lower bound of λ_2 .

In fact, using Liu's result from Section 4.4, we can prove that $1 - n/2^{n-1}$ is the SLEM itself:

$$\rho = \lambda_2 = 1 - \left(\frac{1/n}{1/2^{n-1}} \right)^{-1} = 1 - \frac{n}{2^{n-1}}.$$

This needs little and trivial calculation only, and it uses the trial distribution p_{ij} directly, without any need for modified transitions q_{ij} .

We will prove shortly that there is no polynomial p with the property that

$$\lim_{n \rightarrow \infty} \left(1 - \frac{n}{2^{n-1}} \right)^{p(n)} = 0.$$

This means that there exists no polynomial upper bound of the time that this Markov chain needs to suitably approximate the target distribution.

To prove the former proposition note that for some n_0 sufficiently large,

$$1 - \frac{n}{2^{n-1}} \geq 1 - \frac{1}{1.9^{n-1}},$$

if $n_0 \leq n$. Again, with any polynomial p given, of which the leading coefficient is positive, for a sufficiently large n_1 (which is dependent on p and we choose

it to satisfy $n_0 \leq n_1$), if $n_1 \leq n$, then $0 \leq p(n) \leq 1.9^{n-1}$ and

$$\left(1 - \frac{n}{2^{n-1}}\right)^{p(n)} \geq \left(1 - \frac{1}{1.9^{n-1}}\right)^{p(n)} \geq \left(1 - \frac{1}{1.9^{n-1}}\right)^{1.9^{n-1}} \xrightarrow{n \rightarrow \infty} e^{-1} > 0.$$

■

Now we turn our attention to the ParIS. The task is the same: sampling random elements from uniform distribution on the set of paths leading from the top left to the bottom right corner. Instead of choosing whole new paths in each step, we pick a (connected) part (a ‘window’) of the current path randomly and alter only this part with a certain probability. To keep calculations as easy as possible, we choose solely two-edge-long parts.

A path consists of n edges. If this path is state i of the Markov chain, in other words, it crosses from the left to the right on the i th rung, then the rung is its i th edge. Let us use the following algorithm. We pick two successive edges of the current path with uniform distribution on the set of these pairs: $[j, j + 1]$ ($1 \leq j \leq n - 1$). If both j and $j + 1$ are on the same side of the ladder, the path cannot be altered with replacing only these two edges. If one of them is a rung, then we alter the path with probability $1/2$, or leave it unchanged with probability $1/2$. If we alter the path, there is only one way to do this: if $i = j$, then the new path will have one more edge on the left side, and its $i + 1$ st edge will cross from the left to the right, which means a transition from state i to state $i + 1$; if $i = j + 1$, then the new path will have one edge less on the left side, and its $i - 1$ st edge will cross from the left to the right, so we step from state i to state $i - 1$.

Proposition 5.1.2 *This ParIS method (with window size 2) has SLEM*

$$\rho \leq 1 - \frac{1}{n^3},$$

consequently, this chain has polynomial relaxation time.

PROOF The entries of the primal transition matrix are:

$$\begin{aligned}
p_{i,i+1} &= \frac{1}{n-1} \frac{1}{2}, \quad \text{if } 1 \leq i \leq n-1, \\
p_{i,i-1} &= \frac{1}{n-1} \frac{1}{2}, \quad \text{if } 2 \leq i \leq n, \\
p_{11} &= p_{nn} = 1 - \frac{1}{n-1} \frac{1}{2}, \\
p_{ii} &= 1 - \frac{1}{n-1}, \quad \text{if } 2 \leq i \leq n-1, \\
p_{ij} &= 0, \quad \text{otherwise.}
\end{aligned}$$

The acceptance probabilities α_{ij} ($i, j \in I$, $i \neq j$) are:

$$\begin{aligned}
\alpha_{i,i+1} &= \min \left\{ 1, \frac{\frac{1}{n} p_{i+1,i}}{\frac{1}{n} p_{i,i+1}} \right\} = \min \left\{ 1, \frac{\frac{1}{n} \frac{1}{2(n-1)}}{\frac{1}{n} \frac{1}{2(n-1)}} \right\} = 1, \quad \text{if } 1 \leq i \leq n-1, \\
\alpha_{i,i-1} &= \min \left\{ 1, \frac{\frac{1}{n} p_{i-1,i}}{\frac{1}{n} p_{i,i-1}} \right\} = \min \left\{ 1, \frac{\frac{1}{n} \frac{1}{2(n-1)}}{\frac{1}{n} \frac{1}{2(n-1)}} \right\} = 1, \quad \text{if } 2 \leq i \leq n.
\end{aligned}$$

We will not need to know any other acceptance probabilities.

The transition probabilities are ($i \neq j$):

$$q_{ij} = p_{ij} \alpha_{ij} = p_{ij},$$

because either $\alpha_{ij} = 1$ or $p_{ij} = 0$, and

$$q_{ii} = p_{ii} + \sum_{j \neq i} p_{ij} (1 - \alpha_{ij}) = \begin{cases} (1 - \frac{1}{2(n-1)}) + 0 = 1 - \frac{1}{2(n-1)}, & \text{if } i \in \{1, n\}, \\ (1 - \frac{1}{n-1}) + 0 = 1 - \frac{1}{n-1}, & \text{if } 2 \leq i \leq n-1. \end{cases}$$

Finally, one can observe that $Q = P$.

In order to prove fast convergence, we give a sufficiently good upper bound of the SLEM by the use of Theorem 4.6.3 and Theorem 4.6.4.

For $i, j \in I$, $i < j$, let us choose the path $i, i+1, \dots, j$, and for $i > j$, $i, i-1, \dots, j$. Since the transition matrix Q is symmetric, $\pi(1) = \dots = \pi(n) = 1/n$ is a stationary distribution. One can get

$$|\gamma_{ij}|_Q = |j - i| n 2(n-1),$$

$$\kappa = \max_e \sum_{\gamma_{ij} \ni e} |\gamma_{ij}|_Q \frac{1}{n^2}.$$

For an oriented edge $e = i \rightarrow i + 1$,

$$\sum_{\gamma_{ij} \ni e} |\gamma_{ij}|_Q \frac{1}{n^2} = \left(\sum_{k=1}^i |\gamma_{k,i+1}|_Q + \sum_{k=1}^i |\gamma_{k,i+2}|_Q + \cdots + \sum_{k=1}^i |\gamma_{k,n}|_Q \right) \frac{1}{n^2},$$

and in the following steps by summing arithmetical sequences,

$$\begin{aligned} &= \left(\frac{(i+1)i}{2} + \frac{(i+3)i}{2} + \cdots + \frac{(2n-i-1)i}{2} \right) n \frac{1}{n^2} = \\ &= \frac{((i+1) + (2n-i-1))(n-i)}{2} \frac{i}{2} 2(n-1) \frac{1}{n} = (n-i)i(n-1). \end{aligned}$$

If $e = j \rightarrow j - 1$, then the quantity under consideration is equal to the quantity associated with the oriented edge $1 + n - j \rightarrow 2 + n - j$. Therefore, it is $(j-1)(n-(j-1))(n-1)$.

The function $i \mapsto (n-i)i$ has its maximum among integers in $n/2$, if n is even, and in both $(n-1)/2$ and $(n+1)/2$, if n is odd. One can get

$$\kappa = \begin{cases} \frac{n^2}{4}(n-1), & \text{if } n \text{ is even} \\ \frac{n-1}{2} \frac{n+1}{2}(n-1), & \text{if } n \text{ is odd} \end{cases} \leq \frac{n^2}{4}(n-1).$$

By using Theorem 4.6.3,

$$\lambda_2 \leq 1 - \frac{4}{n^2(n-1)} \leq 1 - \frac{1}{n^3}.$$

Proceeding to eigenvalue λ_n we could use Proposition 4.6.5. As an alternative, to see another example, we use Theorem 4.6.4 directly. We choose the oriented edge $i \rightarrow i$ by itself as a path σ_i .

$$\sum_{\sigma_i \ni e} |\sigma_i|_Q \pi(i) = \begin{cases} 0, & \text{if } e \neq j \rightarrow j \text{ for any } j \in I, \\ |\sigma_j|_Q \pi(j) = \frac{1}{\pi(j)p_{jj}} \pi(j) = \frac{1}{p_{jj}}, & \text{if } e = j \rightarrow j \text{ for some } j \in I. \end{cases}$$

$1/p_{jj}$ is maximal if $2 \leq j \leq n-1$, so

$$\alpha = \frac{1}{1 - \frac{1}{n-1}},$$

and by Theorem 4.6.4,

$$\lambda_n \geq -1 + 2 \left(1 - \frac{1}{n-1} \right) = 1 - \frac{2}{n-1}.$$

These two bounds guarantee that the SLEM is λ_2 . By defining the function $p(n) = n^{3+\varepsilon}$ (for any $\varepsilon > 0$),

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n^3}\right)^{p(n)} = 0. \quad \blacksquare$$

5.2 MIS is slow at sorting by reversals

In this section we wish to introduce a signed permutation that has the property that the MIS has exponential relaxation time when sampling from its minimal sequences of sorting reversals. From a certain viewpoint this improves the result with the Ladder, because it is not based on an artificially constructed graph, instead it shows explicitly that there are troublesome signed permutations which the MIS algorithm might encounter in everyday practice.

Proposition 5.2.1 *For every positive, even integer n there exists a signed permutation of size $5n - 1$, such that sampling from its minimal sequences of sorting reversals the MIS has a relaxation time exponential in input size $5n - 1$.*

PROOF Let us consider the following signed permutation:

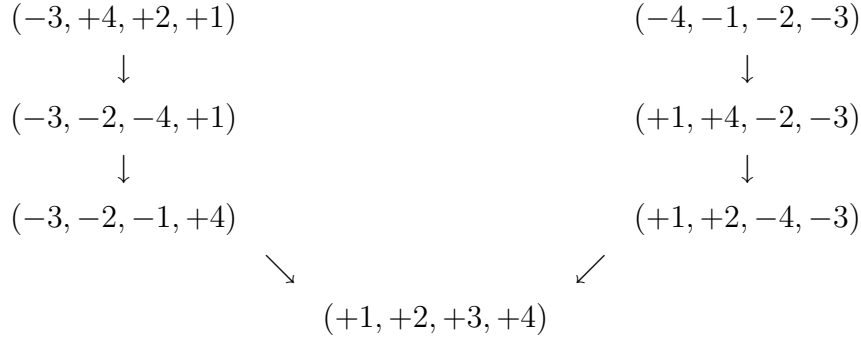
$$(-4, +3, +2, +1).$$

By the Hannenhalli–Pevzner theory there are exactly 26 different optimal sorting sequences of this signed permutation, that is, sequences of reversals to turn it into *id*. Each optimal sequence uses 4 reversals.

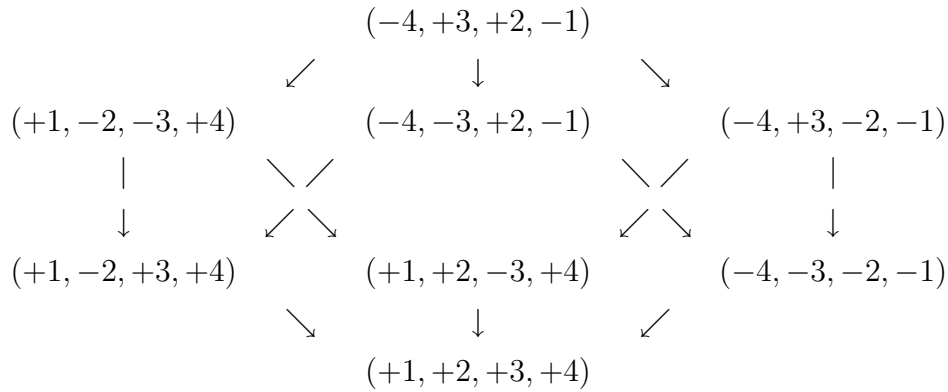
The first optimal reversal transforms $(-4, +3, +2, +1)$ into one of the following six signed permutations:

$$\begin{aligned} &(-3, +4, +2, +1); \quad (-4, -1, -2, -3); \quad (-4, +3, +2, -1); \\ &(-1, -2, -3, +4); \quad (-4, +3, -2, +1); \quad (-4, -3, +2, +1). \end{aligned}$$

A common property of the first two $((-3, +4, +2, +1), (-4, -1, -2, -3))$ is that both have unique optimal sorting sequences:



Each of the other four has six optimal sorting sequences. Each gives a diagram similar to this one:



So $(-4, +3, +2, +1)$ has 26 optimal sorting sequences. Now, the real example is ‘ n instances of $(-4, +3, +2, +1)$ placed one next to another with spacing’:

$$\left(\underline{-4, +3, +2, +1, +5, -9, +8, +7, +6, \dots}, \dots, +5(n-1), \underline{-(5n-1), +(5n-2), +(5n-3), +(5n-4)} \right).$$

To prove exponential relaxation time we shall use Liu’s result from Section 4.4. Our target distribution π is the uniform distribution on all optimal sorting sequences. Therefore the largest importance ratio $w(1) = \pi(1)/p(1)$ is given by the smallest proposal probability $p(1)$.

An optimal sorting sequence of reversals comprises independent optimal sorting sequences of the small, four-character-long blocks. One has to pick one of 26 sequences for every block and then carry out the reversals given by

these n sequences in arbitrary order (but obeying the chosen and thereafter fixed order in every four-reversal-long sequence). As a result of this, there are

$$\pi(1)^{-1} = \frac{(4n)!}{(4!)^n} 26^n$$

optimal sorting sequences.

The least probable path of the random walk procedure yields the smallest proposal probability in the MIS. This probability is realized, for instance, when one chooses to sort small blocks one after another, sorting each by a sorting sequence of the ‘branching type’, that is, the case of $(-4, +3, +2, -1)$, where there were 3, then 2 possible choices for the next step. The probability of one such sorting sequence is

$$\begin{aligned} p(1) &= \frac{1}{6n} \frac{1}{6(n-1)+3} \frac{1}{6(n-1)+2} \frac{1}{6(n-1)+1} \cdot \\ &\cdot \frac{1}{6(n-1)} \frac{1}{6(n-2)+3} \frac{1}{6(n-2)+2} \frac{1}{6(n-2)+1} \cdots \frac{1}{6} \frac{1}{3} \frac{1}{2} \frac{1}{1} = \\ &= \frac{(6n-1)(6n-2)(6n-7)(6n-8)\dots 5 \cdot 4}{(6n)!}, \end{aligned}$$

and consequently, assuming that n is even,

$$\begin{aligned} (4n)!p(1) &= \frac{(6n-1)(6n-2)\dots(3n+5)(3n+4)}{6n(6n-1)\dots(5n+2)(5n+1)} \cdot \\ &\cdot \frac{(3n-1)(3n-2)\dots 5 \cdot 4}{5n(5n-1)\dots(4n+2)(4n+1)} < \\ &< \frac{1}{1} \left(\frac{3}{5}\right)^n. \end{aligned}$$

The spectral gap is

$$\frac{1}{w(1)} = \frac{p(1)}{\pi(1)} < \left(\frac{3}{5}\right)^n \frac{26^n}{(4!)^n} = \frac{1}{c^n} < \frac{1}{(c^{1/5})^{5n-1}},$$

where

$$c = \frac{5}{3} \frac{4!}{26} = 1.538 \cdots > 1,$$

and this means exponential relaxation time with respect to input size $5n-1$. ■

5.3 Dragon's wing

Similarly to the example with the *Ladder*, the name was inspired by the shape of the defining graph, and secondly, by the attractiveness of this name. There are n paths and each path consists of n edges which lead from the top (the tip of the wing) to the bottom (the dragon's body). (See Figure 3.) The task is the same: sampling from the set of all paths, uniformly at random.

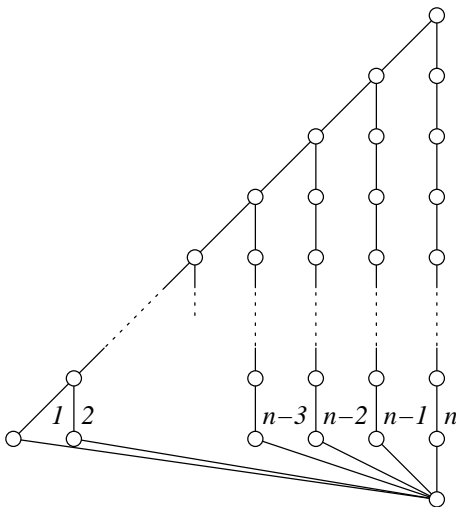


Figure 3: Dragon's wing

Proposition 5.3.1 *The MIS applied to sample from the set of paths of the Dragon's wing has the same SLEM and relaxation time as it had with the Ladder (Section 5.1), therefore it has exponential relaxation time.*

PROOF This is caused by the fact that the *random walk procedure* yields the same trial distribution. This does not seem symbolically at first glance because we change the order of the paths to satisfy Liu's condition for importance weights from Section 4.4: $w(1) \geq w(2) \geq \dots \geq w(n)$. The trial distribution is as follows: for $i, j \in \{1, \dots, n\}$

$$p_{ij} = p(j) = \left(\frac{1}{2}\right)^{(n+1-j) \wedge (n-1)}.$$

Liu's theorem tells us that

$$\rho = \lambda_2 = 1 - \frac{1}{w(1)} = 1 - \frac{n}{2^{n-1}}. \quad \blacksquare$$

However, we have a problem with the ParIS. If we fix a window size k ($2 \leq k \leq n-1$), then the Markov chain will not be irreducible: there will always be paths out of reach. If we fix the window size to be n , then we get the MIS.

To overcome this problem, we pick a random window size $k \in \{2, 3, \dots, n\}$ in every step independently with uniform distribution, and we take the ParIS step with this window size. That means, we pick a connected part of length k of the current path randomly and replace it with another one, using the random walk procedure starting from the starting node of the window, if the random walk ends up in the ending node of the window. If it ends up somewhere else, we do not change the path, so we stay in the same state of the Markov chain. (An obvious necessary condition to be able to leave state i is to draw $k \geq i \vee 2$ and to have the last node of the path in the window.) This MCMC method turns out to be rapidly mixing.

Proposition 5.3.2 *The ParIS with uniform window size distribution on the set $\{2, 3, \dots, n\}$ applied to the Dragon's wing has polynomial relaxation time.*

PROOF Relying on the last formulae of Section 3.6, transition probabilities q_{ij} are

$$q_{ij} = \sum_{k=j}^n \frac{1}{n-1} \frac{1}{n-k+1} \min \left\{ \left(\frac{1}{2} \right)^{k-j+1}, \left(\frac{1}{2} \right)^{(k-i+1) \wedge (k-1)} \right\}, \quad \text{if } i < j,$$

$$q_{ij} = \sum_{k=i}^n \frac{1}{n-1} \frac{1}{n-k+1} \min \left\{ \left(\frac{1}{2} \right)^{(k-j+1) \wedge (k-1)}, \left(\frac{1}{2} \right)^{k-i+1} \right\}, \quad \text{if } i > j,$$

that is,

$$q_{ij} = \frac{1}{n-1} \sum_{k=i \vee j}^n \frac{1}{n-k+1} \left(\frac{1}{2} \right)^{(k-(i \wedge j)+1) \wedge (k-1)}, \quad \text{if } i \neq j.$$

We will prove shortly that $q_{ii} \geq 1/2$, if $n \geq 6$. Hence Proposition 4.6.5 gives that all eigenvalues of transition matrix Q are nonnegative.

Indeed, if $i \neq j$, then

$$q_{ij} = \sum_w \min \{p_{ij}^w, p_{ji}^w\} \leq \min \left\{ \sum_w p_{ij}^w, \sum_w p_{ji}^w \right\} \leq \sum_w p_{ij}^w.$$

It is obvious by the definition of p_{ij}^w that

$$\sum_{j \neq i} \sum_w p_{ij}^w \leq P^{\text{proposal}}(\text{the last node of the path is in the window}).$$

These two observations yield

$$\begin{aligned} q_{ii} &= 1 - \sum_{j \neq i} q_{ij} \geq 1 - \sum_{j \neq i} \sum_w p_{ij}^w \geq \\ &\geq P^{\text{proposal}}(\text{the last node of the path is not in the window}). \end{aligned}$$

With straightforward calculation one can get

$$\begin{aligned} q_{ii} &\geq P^{\text{proposal}}(\text{the last node of the path is not in the window}) = \\ &= \frac{1}{n-1} \left(\frac{1}{2} + \frac{2}{3} + \frac{3}{4} + \dots + \frac{n-2}{n-1} \right) = \\ &= \frac{1}{n-1} \left(\frac{1}{2} + \left(\frac{1}{2} + \frac{1}{6} \right) + \left(\frac{1}{2} + \frac{1}{6} + \frac{1}{12} \right) + \dots \right) \geq \\ &\geq \frac{1}{n-1} \left(\frac{1}{2}(n-2) + \frac{1}{6}(n-3) \right), \end{aligned}$$

and the last term is greater than or equal to $1/2$, if $n \geq 6$.

Just as with the Ladder, we use Theorem 4.6.3 to prove polynomial relaxation time. We choose the same paths in the Markov chain: for $i < j$, we choose the path $i, i+1, \dots, j$, and for $i > j$, $i, i-1, \dots, j$.

To keep calculations easier, each sum in transition probabilities q_{ij} is replaced by its first summand. This is a lower bound of q_{ij} and therefore it yields an upper bound of κ .

$$q_{ij} \geq \frac{1}{n-1} \frac{1}{n - (i \vee j) + 1} \left(\frac{1}{2} \right)^{(i \vee j) - (i \wedge j) + 1}, \quad \text{if } i \neq j.$$

With some calculation similar to that of the Ladder one can get

$$\kappa \leq \frac{1}{2}n^3(n-1),$$

which means polynomial relaxation time. ■

5.4 Another example

We describe a quite complicated graph. It can be seen as a far descendant of the *Dragon's wing* (Section 5.3) with the property that the Paris method applied to sample from the set of its paths is not rapidly mixing even with uniform window size distribution. Figure 4 shows this graph.

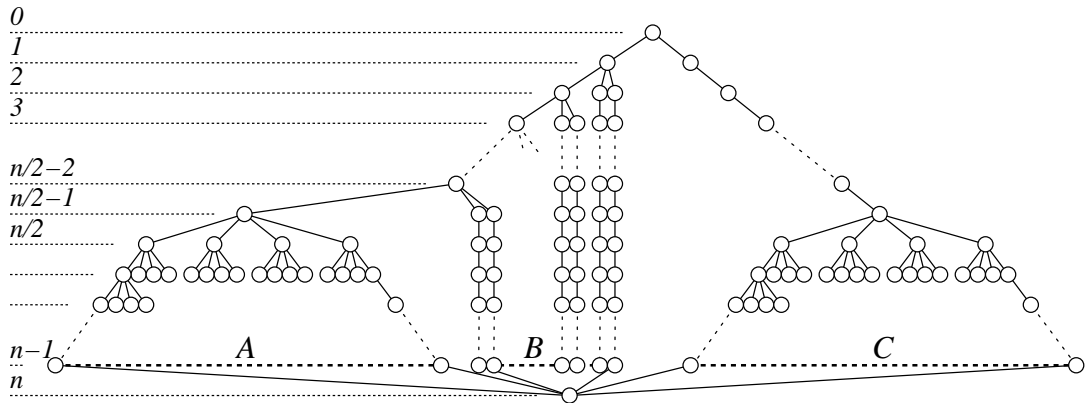


Figure 4: The graph that defines this example. Here $a = 3$, $b = 4$.

Suppose that n is even. Pick integers $a, b \geq 2$, with $a \leq b$. Every path starts at the top node and leads to the one at the bottom. We construct the set of all paths I as the disjoint union of three sets: $I = A \cup B \cup C$.

The top node has two neighbours. If we go and keep to the left, then there is a split in every step from the first node until the $(n/2 - 2)$ nd one. In every step the current path splits into a paths. These new paths do not branch and they lead straight to the last node. They form set B . From the $(n/2 - 1)$ st node until the $(n - 2)$ nd one there is a split in every step into b branches. These new branches also continue branching into b branches in every step and they form set A .

If we start from the top node to the right, there is no branching until the $(n/2 - 2)$ nd node. From the $(n/2 - 1)$ st one, we create an exponentially branching part, identical to set A . We call this set of paths C .

Obviously, $|A| = |C| = b^{n/2}$, $|B| = (a - 1)(n/2 - 2)$.

Assume that we would like to sample random elements from the set of all directed paths leading from the top node to the bottom one, uniformly at random. Relying on what we learnt from the Dragon's wing, we use ParIS with uniform window size distribution on the set $\{2, 3, \dots, n\}$. Unfortunately, we will find that the method has exponential relaxation time. In fact, one should not be surprised, since the underlying idea of the construction is the creation of a bottleneck between $A \cup B$ and C .

Proposition 5.4.1 *The ParIS method with uniform window size distribution on the set $\{2, 3, \dots, n\}$ applied to sample random elements from the set of all directed paths of this graph leading from the top node to the bottom one has exponential relaxation time.*

PROOF We use the notion of conductance and Cheeger's inequality (Theorem 4.5.2) to prove slow mixing. The reader is advised to recall the notations of Section 4.5.

For the subset C of paths (or equivalently, subset of the state space of the Markov chain) $\pi(C) \leq 1/2$, therefore, by definition, $\Phi \leq F(C)/\pi(C)$. If one can show that

$$\frac{F(C)}{\pi(C)} \leq \frac{1}{c^n},$$

for some constant $c > 1$, then, by Cheeger's inequality,

$$1 - 2\frac{1}{c^n} \leq 1 - 2\frac{F(C)}{\pi(C)} \leq 1 - 2\Phi \leq \lambda_2,$$

which means exponential relaxation time.

In the ParIS method, the transition probabilities are

$$q_{ij} = q_{ji} = \sum_w \min \{p_{ij}^w, p_{ji}^w\}.$$

With this specific graph, all transitions we need imply a unique window size and window location (that is, the whole path), hence these sums consist of one summand only. Values p_{ij}^w which will be needed are as follow:

$$\begin{aligned} p_{ij}^w &= \frac{1}{n-1} \cdot 1 \cdot \frac{1}{2b^{n/2}}, \quad \text{if } i \in A, j \in C, \\ p_{ij}^w &= \frac{1}{n-1} \cdot 1 \cdot \frac{1}{2a^{n/2-2}b^{n/2}}, \quad \text{if } i \in C, j \in A, \\ p_{ij}^w &= \frac{1}{n-1} \cdot 1 \cdot \frac{1}{2b^{n/2}}, \quad \text{if } i \in B, j \in C, \\ p_{ij}^w &\geq \frac{1}{n-1} \cdot 1 \cdot \frac{1}{2a^{n/2-2}}, \quad \text{if } i \in C, j \in B. \end{aligned}$$

As a result,

$$q_{ij} = q_{ji} = \frac{1}{n-1} \cdot 1 \cdot \frac{1}{2a^{n/2-2}b^{n/2}}, \quad \text{if } i \in A, j \in C,$$

and $a \leq b$ implies

$$q_{ij} = q_{ji} = \frac{1}{n-1} \cdot 1 \cdot \frac{1}{2b^{n/2}}, \quad \text{if } i \in B, j \in C.$$

In our case, π is the uniform distribution. Consequently, for any $S \subseteq I$, $k \in I$,

$$\frac{F(S)}{\pi(S)} = \frac{\sum_{i \in S, j \in I \setminus S} \pi(i) q_{ij}}{|S| \pi(k)} = \frac{\sum_{i \in S, j \in I \setminus S} q_{ij}}{|S|}.$$

Specifically, for $S = C$,

$$\frac{F(C)}{\pi(C)} = \left(\sum_{i \in C, j \in A} q_{ij} + \sum_{i \in C, j \in B} q_{ij} \right) / |C|.$$

The first summand is

$$\frac{\sum_{i \in C, j \in A} q_{ij}}{|C|} = \frac{|C||A| q_{ij}}{|C|} = \frac{|A|}{(n-1) 2 a^{n/2-2} b^{n/2}} = \frac{1}{(n-1) 2 a^{n/2-2}},$$

while the second one is

$$\frac{\sum_{i \in C, j \in B} q_{ij}}{|C|} = \frac{|C||B| q_{ij}}{|C|} = \frac{|B|}{(n-1) 2 b^{n/2}} = \frac{(a-1)(n/2-2)}{(n-1) 2 b^{n/2}}.$$

Hence $F(C)/\pi(C) \leq 1/c^n$ for some constant $c > 1$, if n is big enough. This proves our statement. ■

The role of set B can be seen clearly now. It makes transitions from C to A very improbable, what is used in the first summand, still it has relatively few elements, thus letting the second summand be small.

Note that the number of paths (i.e. the number of states of the Markov chain) is exponential in n and this fact has a great effect on the mixing rate. Exponentially many paths can indeed emerge in the problem of sorting by reversals. Therefore we need to cope with this situation if we try to find a fast algorithm to sample random elements from all optimal sorting sequences.

5.5 The ParIS needs to choose the whole sorting sequence with positive probability

In this section we show a signed permutation which proves that in the problem of sampling from optimal sequences of sorting reversals the ParIS needs to choose the whole sorting path as a window with positive probability to guarantee the irreducibility of the Markov chain.

We assume familiarity with the Hannenhalli–Pevzner theory. One can find introductions to the topic in [10, 19, 13, 5, 23].

Proposition 5.5.1 *Let n be a positive integer. The set of optimal sorting sequences of reversals of any signed permutation σ whose breakpoint graph consists of two hurdles and an additional cycle as seen in Figure 5, can be partitioned into two disjoint subsets, such that signed permutations through which sorting sequences of one subset lead are disjoint from signed permutations of sorting sequences of the other subset, except for the two endpoints, σ and id .*

Corollary 5.5.2 *Any ParIS method to sample from the set of all optimal sequences of sorting reversals of σ has to cut the whole sequence out with positive probability to ensure that the Markov chain is irreducible.*

To prove the proposition we will use an easy lemma, which is proved in [23].

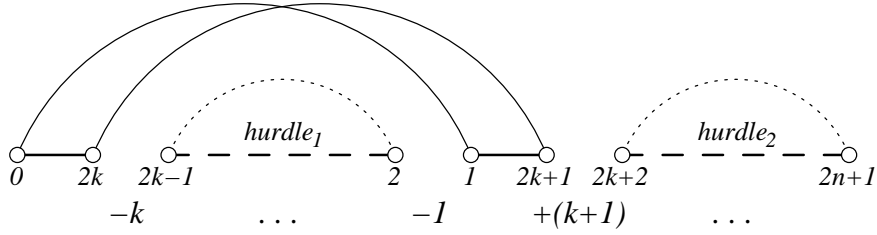


Figure 5: The breakpoint graph of the example, σ .

Lemma 5.5.3 *If there is no fortress in the signed permutation, then a reversal is a sorting reversal if and only if, first, it does not introduce a new fortress, and secondly, $\Delta c = -1$ and $\Delta h = -2$; or $\Delta c = 0$ and $\Delta h = -1$; or $\Delta c = 1$ and $\Delta h = 0$. Consequently, under the fortress-free assumption, a reversal that increases the number of hurdles cannot be a sorting reversal.*

PROOF of the Proposition. Lemma 7 of [23] guarantees that all our permutations will be fortress-free. The structure of σ with the two hurdles implies that σ consists of integers with negative signs between the position of $-k$ and -1 , and integers with positive signs to the right from -1 :

$$\sigma = \left(-k, \{ -2, -3, \dots, -(k-1) \}, -1, \{ +(k+1), +(k+2), \dots, +n \} \right).$$

There are three types of possible starting sorting reversals of permutation σ :

1. a hurdle merging,
2. the reversal that sorts the additional cycle, that is, that reverses the interval from $-k$ to -1 ,
3. a hurdle cutting: the cutting of either $hurdle_1$ or $hurdle_2$.

We prove the Proposition by showing that optimal sorting sequences of Case 1 do not contain any permutations that any sequence of Case 2 or 3 contains, except for σ and id . The sorting sequences of Case 1 shall be the elements of the first set of the partition we are looking for, the second set

shall consist of sequences of Case 2 and 3.

Case 1 We claim that if we start with a hurdle merging, the permutation will retain the following structure until the last reversal:

$$(-, -, \dots, -, \pm 1, +, +, \dots, +)$$

(Claim 1). We also claim that the cycle formed by $0, 2k, \dots, 1, 2k+1, \dots$ of the unsigned permutation will be in existence until the last reversal: either in this form (see Figure 6) or in the form $0, 2k, \dots, 2k+1, 1, \dots$ (see Figure 7)

(Claim 2). This implies that nothing but the last reversal will take ± 1 to the first position (Claim 3).

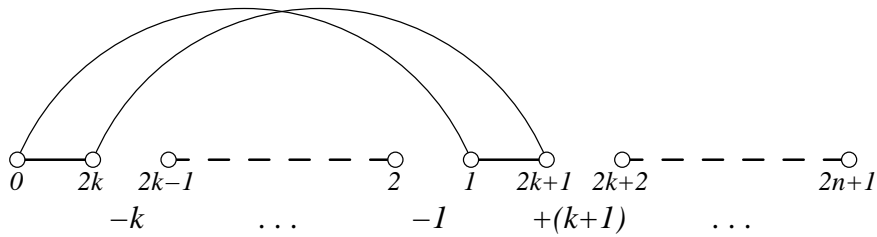


Figure 6: The additional cycle is oriented.

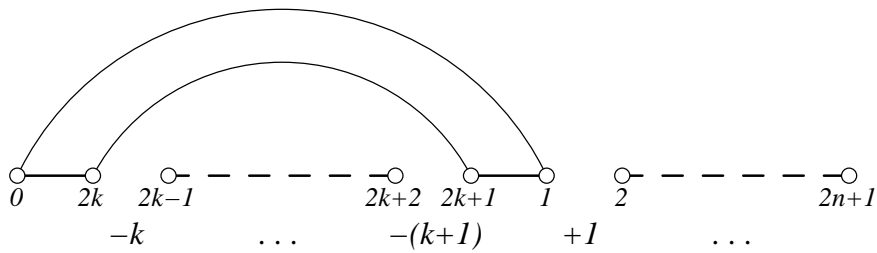


Figure 7: The additional cycle is unoriented.

Two reality edges of the same cycle are said to be *convergent*, if in a traversal of the cycle, they induce the same circular ordering of the vertices of the breakpoint graph, otherwise the edges are *divergent*. It is well-known

that any reversal that acts on divergent reality edges will split the cycle to which the edges belong, and any reversal that acts on convergent edges will not split the cycle to which they belong.

Assume that desire edge $2j, 2j + 1$ ($0 \leq j \leq n$) of the unsigned permutation is oriented. Then the interval on which the corresponding sorting reversal acts is one of the following four:

$$(\dots, \underline{-j}, \dots, +(j+1), \dots), \quad (4)$$

$$(\dots, \underline{+(j+1)}, \dots, -j, \dots),$$

$$(\dots, +j, \dots, \underline{-(j+1)}, \dots),$$

$$(\dots, -(j+1), \underline{\dots}, +j, \dots). \quad (5)$$

It can be seen that one endpoint of the interval and the neighbour of the other endpoint that does not belong to the interval have different signs (Claim 4).

Since we started with a hurdle merging, there are no hurdles left in the permutation. Lemma 5.5.3 guarantees that no hurdle will arise later and the permutation remains hurdle-free. The Lemma also implies that each reversal must increase the number of cycles, therefore must act on divergent edges of the same cycle. Claim 4 ensures that Claim 1 holds.

Indeed, a reversal reverses all signs on which it has an effect. If there are both positive and negative numbers in the reverted interval, then it is clear that the reversal transforms the permutation $(-, \dots, -, \pm 1, +, \dots, +)$ to some permutation with the same structure. If there are not, then we can use Claim 4: there are both positive and negative numbers in or right next to the reverted interval. Nothing but these two cases can cause a problem:

$$(-, \dots, -, \underline{-}, \dots, \underline{-}, +1, +, \dots, +) \rightarrow (-, \dots, -, +, \dots, +, +1, +, \dots, +),$$

$$(-, \dots, -, -1, \underline{+}, \dots, \underline{+}, +, \dots, +) \rightarrow (-, \dots, -, -1, -, \dots, -, +, \dots, +).$$

But none of the four cases of Claim 4 fits any of these. The first one looks like (4), with $j + 1 = 1$, but then the interval should contain a virtual -0 which

is impossible. The second one looks like (5), with $j + 1 = 1$ again, but then 0 should be in the interval, which is also impossible. This proves Claim 1.

However, such a reversal might occur:

$$(-, -, \dots, -, \underline{-1, +, +, \dots, +}) \longrightarrow (-, -, \dots, -, +1).$$

This is like (5), with $+j$ in the last position.

In contrast, sorting σ starting with a hurdle merging (that is, in Case 1), reversals like

$$(\underline{-, -, \dots, -, +1, +, +, \dots, +}) \longrightarrow (-1, +, +, \dots, +) \quad (6)$$

cannot occur (Claim 5), no matter that this looks like (4), with $-j$ in the first position.

Such a reversal can only be a sorting one if it acts on divergent edges of the same cycle. It is clear that since we are in Case 1 the first reversal is a hurdle merging and it leads to a configuration shown in Figure 7. The second reversal cannot be like (6), since if this was the case, then it would act on two different cycles, it would not be a sorting reversal. Therefore it must act on divergent reality edges that earlier belonged to the hurdles. As long as this cycle is unoriented, (6) is not a sorting reversal, and as a result, $-k$ remains fixed in the first position.

After a number of reversals the additional cycle might become oriented again (Figure 6). It cannot happen that a sorting reversal sorts the additional cycle to get some permutation

$$(+1, +, +, \dots, +),$$

because this permutation contains a hurdle, unless it is *id*. In the latter case we finished proving Claim 2. When this is not the case, then the additional cycle remains. As long as it is oriented, it can only be turned into two trivial cycles, if this sorting reversal is the last reversal. Note that if the cycle is oriented, then (6) cannot happen, because the ± 1 has negative sign.

If the cycle becomes unoriented again (Figure 7), then (6) is not sorting for the same reason why the second sorting reversal was not of type (6),

either. If the cycle is unoriented, then the next reversal must act on reality edges distinct from both $0, 2k$ and $2k + 1, 1$ of the unsigned permutation. This proves both Claim 5 and Claim 2, consequently, Claim 3.

After any of the reversals, every nontrivial component intersects the additional cycle (Claim 6). This is an easy observation: otherwise there would be a component consisting of numbers with identical signs, therefore this component would be a hurdle. But the permutation is hurdle-free after merging the two hurdles.

Claim 3 will separate signed permutations of optimal sorting sequences of Case 1 from that of Case 2: Claim 3 will not be true for permutations of sorting sequences of Case 2. Claim 6 will make the same for Case 1 and Case 3.

Case 2 After the specific first reversal of Case 2 we get

$$\begin{aligned} & (\text{a trivial cycle, } hurdle_1 \text{ reverted, a trivial cycle, } hurdle_2) = \\ & = (+1, \dots, +k, +(k+1), \dots). \end{aligned}$$

The next reversal is either a hurdle merging or a hurdle cutting. In both cases it decreases the number of hurdles, to 0 or to 1. By Lemma 5.5.3, using sorting reversals the number of hurdles is nonincreasing.

We claim that all forthcoming reversals leave (both the position and the sign of) $+1$ unchanged. Indirectly, if this was not the case, then the reversal which changes (the position or the sign of) $+1$ would decrease the number of cycles by one. By Lemma 5.5.3, two hurdles must have been disappeared. This could only happen if there were two hurdles still present, which is only possible right after the very first reversal. This is the permutation shown above. But in this case, no hurdle merging affects $+1$. This is a contradiction.

Case 3 After cutting one hurdle, there will be at most one hurdle in forthcoming signed permutations, because there is no fortress in our permutations.

If the additional cycle is destroyed during sorting, then Claim 6 is trivially false. If in any step the additional cycle still existed, and, indirectly, a desire

edge intersected the additional cycle, then this intersecting edge would have come into existence by a reversal that had acted on both sides of the additional cycle and had connected two different components. Therefore it had decreased the number of cycles. By Lemma 5.5.3 $\Delta c < 0$ implies $\Delta h = -2$, which is a contradiction because the number of hurdles is strictly less than two. ■

5.6 An application of the coupling method

In this section we give a coupling argument to prove that if $d = 2$, then the ParIS with window size 2, applied to sample from the set of simplified sequence alignments (or equivalently, to sample from paths on the square grid from bottom left to top right with steps \rightarrow or \uparrow) has a polynomial ($O(n^4)$) mixing time.

Pick the square in \mathbb{Z}^2 with edge lengths n , and vertices $(0;0)$, $(n;0)$, $(n;n)$ and $(0;n)$. We aim to sample uniform random elements from the set of paths connecting the bottom left corner $(0;0)$ with the top right one $(n;n)$ that start at $(0;0)$ and consist of steps from a point in \mathbb{Z}^2 to another one where each step is of the following two kinds: \rightarrow , \uparrow (or equivalently, $(+1;0)$, $(0;+1)$). Let I denote the set of such paths.

A path is composed of $n \rightarrow$ steps and $n \uparrow$ steps. We represent each path with a $2n$ -tuple: we write 1 instead of \rightarrow and 0 instead of \uparrow .

Since our main concern is a more complicated model of sequence alignments, we use a method which is admittedly not the most suitable for this one, but which can be generalized more easily. One should keep in mind that this problem can be solved by a direct method that is more efficient than MCMC simulation.

We define a Markov chain $X = (X_k)$ with state space I whose stationary distribution π is the uniform distribution on I .

Suppose that $\mathbf{x} = (x_1, \dots, x_{2n}) \in I$ (that is, $x_1, \dots, x_{2n} \in \{0, 1\}$ and $\sum_{i=1}^{2n} x_i = n$) is the current state of the chain. A transition of the chain is a possible swap of two neighbouring entries of \mathbf{x} :

Draw a pair of neighbouring coordinates of \mathbf{x} uniformly at random (there are $2n - 1$ possible choices) and swap them with probability $1/2$.

(Note that this is the ParIS with window size 2 on the set of paths.) If we get to a different state, say \mathbf{x}' , then we say that we made a *flip* on the path \mathbf{x} . For neighbouring states \mathbf{x} and \mathbf{x}' ($\mathbf{x} \neq \mathbf{x}'$),

$$p_{\mathbf{x}\mathbf{x}'} = \frac{1}{(2n - 1)2}.$$

Transition probabilities form the transition matrix P . P is symmetric, therefore π is the stationary distribution of this Markov chain indeed.

In Section 4.2 we defined the following function:

$$\tau_i(\varepsilon) := \min \{k_0 \in \mathbb{N} : d_V(\delta_i^T P^k, \pi^T) \leq \varepsilon \text{ for all } k \geq k_0\}.$$

We are now ready to formulate the main result of this section.

Proposition 5.6.1 *There exists some $c > 0$ such that*

$$\max_{i \in I} \tau_i(\varepsilon) \leq \frac{cn^4}{\varepsilon}.$$

PROOF We follow a similar argument to that of Aldous [1], and prove this bound via a coupling argument.

We construct two dependent versions (X^1, X^2) of X with arbitrary initial states and we show that $X_k^1 = X_k^2$ for all $k \geq \tau$ for some random time τ , and then we give an upper bound of τ .

One can choose X^2 to start from the stationary distribution. This implies that $(X_k^1)_{k \geq \tau}$ is also in equilibrium.

We start with constructing the chains. At time k

$$(X_k^1, X_k^2) = ((x_1^1, \dots, x_{2n}^1), (x_1^2, \dots, x_{2n}^2)).$$

In each position $i \in \{1, \dots, 2n\}$, $x_i^1 = x_i^2$, or not. If equality holds, we say that position i is *marked*:

$$\begin{array}{cccccccc} & & & & \text{M} & \text{M} & & \text{M} \\ X_k^1 = & (& \dots & 0 & 1 & 0 & 1 & 1 & 0 & \dots &), \\ X_k^2 = & (& \dots & 1 & 1 & 0 & 0 & 0 & 0 & \dots &). \end{array}$$

We specify a transition of the coupled chain (X_k^1, X_k^2) . Pick $i \in \{1, \dots, 2n-1\}$ uniformly at random (i.e. choose neighbours $i, i+1$) and swap x_i^1 for x_{i+1}^1 with probability $1/2$. Otherwise stay in state X_k^1 . (If $x_i^1 = x_{i+1}^1$, then a swap also results in remaining in X_k^1 .) If $x_i^1 \neq x_{i+1}^1$, then

$$P\left(X_{k+1}^1 = (x_1^1, \dots, x_{i-1}^1, x_{i+1}^1, x_i^1, x_{i+2}^1, \dots, x_{2n}^1) \mid X_k^1\right) = \frac{1}{(2n-1)2}.$$

We define the transition rule of X^2 as follows:

- (1) If at least one of $i, i+1$ is a marked position, then do the same transition in X^2 as in X^1 : swap the same two positions or do nothing. Note that in both cases the number of marks is preserved.
- (2) If none of $i, i+1$ is a marked position, then draw one of the neighbouring positions which are both unmarked, uniformly at random (there is at least one: $i, i+1$), say $j, j+1$, and
 - if we made a flip in this step in X^1 , then do nothing: $X_{k+1}^2 := X_k^2$;
 - if we did not, then make a flip in X^2 , that is, swap x_j^2 for x_{j+1}^2 .

In this case the number of marks remains the same or increases by 2.

This rule ensures that the number of marks is nondecreasing, moreover, no mark ever disappears. The more important observation is that each chain evolves according to the transition rule given by P .

Now we turn to bounding random time τ .

Obviously, the marking rule is introduced to indicate that $X_k^1 = X_k^2$: when there are $2n$ marks, then this is already the case. It is sufficient to count the number of marks only at positions where 1 is the common coordinate. There are $2n$ marks if and only if this latter quantity is n .

We would like to track the movement of 1's to see when they get marked. Let us introduce nonnegative integer random times $\tau_1, \tau_2, \dots, \tau_n$. We stress that they are lengths of time intervals, and they are not moments in time.

When $k = 0$, there may already be marked 1's, let us say, there are i of them. Then define $\tau_1 = \dots = \tau_i = 0$. Now we would like to track the movement (or, the random walk) of the leftmost unmarked 1. We call this the *candidate*, because it is a candidate to be the first 1 to become marked:

$$\begin{array}{cccccc} & \text{M} & \text{M} & \text{M} & \text{C} & & \text{M} \\ X_k^1 = & (& 0 & 1 & 1 & 0 & 0 & 1 & 1 & \dots &), \\ X_k^2 = & (& 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & \dots &). \end{array}$$

No matter which 1 will be the next to turn marked, we track this one. If this one becomes marked after k steps, we define $\tau_{i+1} := k$. At time k there are j ($j \geq i + 1$) marked 1's. If $j \geq i + 2$, then define $\tau_{i+2} = \dots = \tau_j = 0$. Now we search for the leftmost unmarked 1 and do the same procedure again and again.

If each 1 is marked, then we have

$$\sum_{i=1}^n \tau_i = \tau.$$

Let us investigate a related problem. We use arguments of Feller's book [9].

Consider a symmetric random walk on $\{1, 2, \dots, 2n\}$. Let us introduce the notation $p := 1/(2n - 1)2$. If $2 \leq i \leq 2n - 1$, the transition probabilities are given by

$$\begin{aligned} p_{i,i-1} &= p, \\ p_{i,i+1} &= p, \\ p_{ii} &= 1 - 2p, \\ p_{12} &= p, \\ p_{11} &= 1 - p, \\ p_{2n,2n-1} &= p, \\ p_{2n,2n} &= 1 - p. \end{aligned}$$

Let D_z ($z \in \{1, 2, \dots, 2n\}$) denote the expectation of the time to first reach state $2n$, if the random walk is started from z . D_z ($1 \leq z \leq 2n$)

We arranged in rule (2) that when there is a chance of turning marked, exactly one flip is made. We did it so to ensure that the candidate will not miss any other 1's, in other words, to avoid such situations:

$$\begin{array}{ccc} & \text{C} & \\ X_k^1 = & (\dots | 1 \ 0 | \dots) & \rightsquigarrow & X_{k+1}^1 = & (\dots | 0 \ 1 | \dots) \\ X_k^2 = & (\dots | 0 \ 1 | \dots) & & X_{k+1}^2 = & (\dots | 1 \ 0 | \dots) \\ & & & & \text{C} \end{array}$$

Therefore, for the appropriate i , $E(\tau_{i+1})$ is bounded above by a certain D_z , which is bounded by D_1 . By the coupling inequality (Theorem 4.12.2) and the well-known Markov-inequality

$$d_V(X_k^1, X_k^2) \leq P(\tau > k) \leq \frac{E(\tau)}{k}.$$

By using the fact that $\tau = \sum_{i=1}^n \tau_i$,

$$\frac{E(\tau)}{k} = \sum_{i=1}^n \frac{E(\tau_i)}{k} \leq \frac{nD_1}{k} = \frac{2n^2(2n-1)^2}{k},$$

which proves the proposition. ■

This proof has a remarkable weakness. Namely, we sum all $E(\tau_i)$'s to get an upper bound of $E(\tau)$. It seems possible to prove an $O(n^3)$ bound if we manage to take into account that random walks of 1's go on simultaneously.

5.7 Conclusions and future work

We investigated Markov chains that converge asymptotically to the uniform distribution on the set of all optimal sequences of sorting reversals of a signed permutation. We proved that, in general, the MIS is not a fast algorithm. Our results with the *Ladder* and the *Dragon's wing* suggest that the ParIS might be faster than the MIS. We proved that the whole sorting sequence must be allowed to be cut out as a window to ensure the irreducibility of the Markov chain.

We hope that the ParIS turns out to be a fast algorithm with an appropriate window size distribution. Further investigations are needed whether

some window size distribution exists that makes the ParIS applied to this problem a polynomially mixing MCMC method.

We proved that the ParIS with window size 2 has polynomial mixing time when sampling from the set of simplified sequence alignments, uniformly at random. We think that the proof is more important than the result. The proof may be a starting point for future investigations, because it might be developed to yield similar bounds for algorithms designed to solve problems of more biological relevance.

Acknowledgements

The author would like to thank the enormous amount of work Miklós István did while supervising the making of this thesis. His ideas and inspiring enthusiasm became fundamental to our joint work. Márkus László made useful comments on the whole thesis and helped to give it a final form. The author is grateful to Arató Miklós, Lovász László, Michaletzky György, Móri Tamás and Tusnádi Gábor for helpful discussions, for bringing different articles and books to the author's attention.

References

- [1] David J. Aldous, *Mixing time for a Markov chain on cladograms*, *Combin. Probab. Comput.* 9, pp. 191–204, 2000.
- [2] David J. Aldous, James A. Fill, *Reversible Markov Chains and Random Walks on Graphs* (monograph in preparation), available via <http://www.stat.berkeley.edu/~aldous>
- [3] David J. Aldous, *Some inequalities for reversible Markov chains*, *Journal of the London Mathematical Society* (2), 25, pp. 564–576, 1982.
- [4] D. Bader, B. Moret, M. Yan, *A linear-time algorithm for computing inversion distance between signed permutations with an experimental study*. In F. Dehne, J.-R. Sack, R. Tamassia, editors, *Algorithms and Data Structures: Seventh International Workshop, WADS 2001, Brown University, Providence, RI, August 8-10, 2001, Proceedings*, volume 2125 of *Lecture Notes in Computer Science*, pp. 365–376, Springer, 2001. Alternatively, *J. Comput. Biol.* 8, 5, pp. 483–491, 2001.
- [5] Anne Bergeron, *A very elementary presentation of the Hannenhalli-Pevzner Theory*. In A. Amir, G. M. Landau, editors, *Combinatorial Pattern Matching, 12th Annual Symposium, CPM 2001 Jerusalem, Israel, July 1-4, 2001 Proceedings*, volume 2089 of *Lecture Notes in Computer Science*, pp. 106–117, Springer, 2001.
- [6] Pierre Brémaud, *Markov Chains: Gibbs Fields, Monte Carlo Simulation, and Queues*, *Texts in Applied Mathematics*, Springer, New York, 1999.
- [7] Persi Diaconis, Daniel Stroock, *Geometric bounds for eigenvalues of Markov chains*, *The Annals of Applied Probability*, 1, 1, pp. 36–61, 1991.
- [8] Richard Durbin, Sean R. Eddy, Anders Krogh, Graeme Mitchison, *Biological sequence analysis, Probabilistic models of proteins and nucleic acids*, Cambridge University Press, 1998.

- [9] William Feller: *An Introduction to Probability Theory and its Applications*, Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, Inc., 1966.
- [10] Sridhar Hannenhalli, Pavel A. Pevzner, *Transforming cabbage into turnip (polynomial algorithm for sorting signed permutations by reversals)*, in Proceedings of the Twenty-Seventh Annual ACM Symposium on Theory of Computing, Las Vegas, Nevada, pp. 178–189, 29 May–1 June 1995. Full version in the JACM **46**(1), pp. 1–27, 1999.
- [11] W. K. Hastings, *Monte Carlo sampling methods using Markov chains and their applications*, Biometrika **57**(1), pp. 97–109, 1970.
- [12] Mark Jerrum, Alistair Sinclair, *Approximating the permanent*, SIAM Journal on Computing, **18**(6), pp. 1149–1178, 1989.
- [13] Haim Kaplan, Ron Shamir, Robert E. Tarjan, *A faster and simpler algorithm for sorting signed permutations by reversals*, SIAM Journal of Computing, **29**(3), pp. 880–892, 1999.
- [14] Bret Larget, Donald L. Simon, Joseph B. Kadane, *Bayesian phylogenetic inference from animal mitochondrial genome arrangements*, Journal of the Royal Statistical Society B **64**, Part 4, pp. 681–693, 2002.
- [15] Torgny Lindvall, *Lectures on the Coupling Method*, Wiley Series in Probability and Mathematical Statistics, Wiley-Interscience, 1992.
- [16] Jun S. Liu, *Metropolized independent sampling with comparisons to rejection sampling and importance sampling*, Statistics and Computing **6**(2), pp. 113–119, 1996.
- [17] Jun S. Liu, *Monte Carlo Strategies in Scientific Computing*, Springer Series in Statistics, Springer, New York, 2001.
- [18] Nicholas Metropolis, Arianna W. Rosenbluth, Marshall N. Rosenbluth, Augusta H. Teller, Edward Teller, *Equations of state calculations by fast*

- computing machines*, Journal of Chemical Physics **21**(6), pp. 1087–1091, 1953.
- [19] Pavel Pevzner, *Computational Molecular Biology*, MIT Press, Cambridge, Massachusetts, 2000.
- [20] Svetlozar Todorov Račev, *Probability metrics and the stability of stochastic models*, Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, Inc., 1991.
- [21] Jason Schweinsberg, *An $O(n^2)$ bound for the relaxation time of a Markov chain on cladograms*, Random Struct. Alg., 20, pp. 59–70, 2001.
- [22] Eugene Seneta, *Non-negative Matrices and Markov Chains*, Second Edition, Springer Series in Statistics, Springer-Verlag, New York, 1981.
- [23] Adam C. Siepel, *An algorithm to find all sorting reversals*, Proceedings of the 6th Annual International Conference on Research in Computational Biology (RECOMB '02), pp. 281–290, 2002.
- [24] Alistair Sinclair, *Improved bounds for mixing rates of Markov chains and multicommodity flow*, Combinatorics, Probability & Computing 1, pp. 351–370, 1992.

Abbreviations

gcd greatest common divisor

MCMC Markov chain Monte Carlo

MIS Metropolised Independent Sampler (or Sampling)

ParIS Partial Independent Sampler (or Sampling)

SLE second largest eigenvalue

SLEM second largest eigenvalue modulus