

The Mathematics of Card Shuffling

Johan Jonasson ^{*†‡}

April 2009

1 Introduction

A Markov chain is a random process $\{X_t\}$, where, in discrete time, $t \in \mathbb{Z}_+$ and in continuous time, $t \in \mathbb{R}_+$ and the X_t 's take values in some state space S , such that given the process up to time t , the distribution of the future of the process only depends on X_t . When S is the symmetric group, i.e. the set of permutations of a number of elements, we think of the Markov chain as a *card shuffling chain*.

In this course, all Markov chains will live on a finite state space S . By enumerating the states $1, 2, \dots, |S|$, we may identify S with $[|S|] = \{1, 2, \dots, |S|\}$.

A Markov chain in discrete time, $\{X_0, X_1, X_2, \dots\}$ is governed by the starting distribution $\mathbb{P}(X_0 \in \cdot)$ and the transition matrix $P = [p_{ij}]_{i,j \in S}$ where the p_{ij} 's are the transition probabilities

$$p_{ij} = \mathbb{P}(X_{t+1} = j | X_t = i).$$

In continuous time, we replace the transition matrix by the generator $Q = [q_{ij}]_{i,j \in S}$, where for $i \neq j$, q_{ij} is the intensity for a jump from i to j . More precisely, this means that when the process is in state i , the time to the next jump is exponentially distributed with intensity $\sum_{j \neq i} q_j$ and the next jump goes to j with probability $q_j / (\sum_{j \neq i} q_j)$. As a convention $q_{ii} = -\sum_{j \neq i} q_j$ so that all row sums of the generator are 0.

*Chalmers University of Technology

†Göteborg University

‡jonasson@math.chalmers.se

In discrete time, a *stationary distribution* is a probability distribution π on S , such that

$$\pi P = \pi.$$

In other words, a stationary distribution is a left eigenvector of P to the eigenvalue 1. In continuous time, the corresponding relation is

$$\pi Q = 0$$

and π is a left eigenvector of the eigenvalue 0. The following theorem is central for any first course in Markov chains.

Theorem 1.1 *Any irreducible and, in discrete time aperiodic, Markov chain $\{X_t\}$ has a unique stationary distribution π , such that $\mathbb{P}(X_t = i) \rightarrow \pi(i)$ for all i .*

The question in focus of this course is *how large t needs to be for $\mathbb{P}(X_t \in \cdot)$ to be close to π* . Then we must first of all decide on what to mean by “close”. There are several different criteria that are used for this. We will use only two of these.

Definition 1.1 *Let μ and π be two probability measures on S and let $p \in [1, \infty)$. Then the total variation distance between μ and π is given by the norm:*

$$\|\mu - \pi\|_{TV} = \frac{1}{2} \sum_{i \in S} |\mu(i) - \pi(i)| = \max_{A \subseteq S} (\mu(A) - \pi(A)).$$

The L^p -norm of the finite signed measure ν with respect to π is given by

$$\|\nu\|_p = \|\nu\|_{L^p(\pi)} = \left(\sum_{i \in S} \left| \frac{\nu(i)}{\pi(i)} \right|^p \pi(i) \right)^{1/p} = \left(\mathbb{E}_\pi \left| \frac{\nu(X)}{\pi(X)} \right|^p \right)^{1/p}.$$

Obviously $\|\mu - \pi\|_1 = 2\|\mu - \pi\|_{TV}$. By the Cauchy-Schwarz inequality, $\|\mu - \pi\|_p \leq \|\mu - \pi\|_q$ when $p \leq q$. Take it as an exercise to show this, as well as to show that

$$\|\mu - \pi\|_2^2 = \sum_i \frac{|\mu(i) - \pi(i)|^2}{\pi(i)} = \sum_i \frac{\mu(i)^2}{\pi(i)} - 1.$$

The most common criterion for closeness to the stationary distribution, and the one that will mainly focus on, is in terms of total variation distance. Later on, we will also work in terms of the L^2 -norm. We define τ_{mix} and $\hat{\tau}$ as

$$\tau_{\text{mix}} = \inf \left\{ t : \|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \leq \frac{1}{4} \right\}.$$

$$\hat{\tau} = \inf\{t : \|\mathbb{P}(X_t \in \cdot) - \pi\|_2 \leq \frac{1}{2}\}.$$

By the above observation, $\tau_{\text{mix}} \leq \hat{\tau}$, so convergence in L^2 is stronger than in total variation.

Usually when dealing with convergence rates problems, one lets $|S| \rightarrow \infty$ in some way. This means that one has MC:s $\{X_t^n\}$ with stationary distributions π^n . This is all very natural for card shuffling, where the number of cards increases, or for random walks on graphs such as \mathbb{Z}_n^d , \mathbb{Z}_2^n , K_n etc, where n tends to infinity. The results are almost always asymptotical in $|S|$, which means that we still cannot give a precise answer for a particular case. However we will have reached another level of confidence.

Definition 1.2 *For the convergence time of $\{X_t^n\}$, we say that $t = t^n$ is*

- *upper bound, if $\limsup_n \|\mathbb{P}(X_t^n \in \cdot) - \pi^n\| \leq 1/4$,*
- *lower bound, if $\liminf_n \|\mathbb{P}(X_t^n \in \cdot) - \pi^n\| \geq 1/4$,*
- *threshold, if, for any $a > 0$, $(1 - a)t$ and $(1 + a)t$ are lower and upper bounds respectively.*

Thresholds are very strong objects, but they surprisingly often exist. A threshold is of course on top of the wish-list. Second best is if we can find upper and lower bounds of the same order. If this also cannot be done, then of course any other information, such as e.g. only an upper bound, is also interesting.

On discrete versus continuous time. Since our primary focus is on card shuffling, we will basically only be interested in discrete time MC:s. Sometimes it is for technical reasons more convenient to work in continuous time. If one wants to do that, one can always “continuize” a discrete time MC, by letting times between updates be exponential with intensity 1, instead of being deterministically 1. However, one should be aware the there is no general way to translate results on mixing between corresponding discrete-time and continuous-time MC:s. Intuitively, the mixing time should be the same for both cases. This is often true, but not always and sometimes only for some criteria for mixing. The main reason for the discrepancy is that for continuous time, there is an extra randomness in the number of moves that have actually taken place at a certain point. Often this causes faster mixing for continuous time than for discrete time. An extreme case is when the discrete-time version MC is periodic and hence never mixes,

but the continuous version does. Sometimes however, for fast-mixing chains, the situation can be the opposite.

On the other hand, there is one general situation where the discrete versus continuous time problem vanishes, namely if the discrete-time MC is *lazy*, i.e. if there is a constant $a > 0$ independent of $|S|$ such that $p_{ii} \geq a$ for every i . In this case the same kind of randomness in the number of actual moves is present also in the discrete setting.

From now on, it is assumed that time is discrete until further notice.

Doesn't the second eigenvalue say all we need to know? Let P be the transition matrix and suppose that P has $n = |S|$ distinct eigenvalues. Let $1 = \lambda_1 > |\lambda_2| > \dots > |\lambda_n|$ where the λ 's are the eigenvalues of P , and write ϕ_i for the corresponding left eigenvectors (so that $\phi_1 = \pi$). Write the starting distribution in eigenvector base:

$$\mathbb{P}(X_0 \in \cdot) = \sum_i c_i \phi_i.$$

Then (since c_1 must be 1)

$$\mathbb{P}(X_t \in \cdot) = \mathbb{P}(X_0 \in \cdot)P^t = \pi + \sum_{i=2}^n c_i \phi_i P^t = \pi + \sum_{i=2}^n c_i \lambda_i^t \phi_i.$$

Hence

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} = \left\| \sum_{i=2}^n c_i \lambda_i^t \phi_i \right\|_{TV}$$

which behaves like $|\lambda_2|^t$ for large t . Thus it would seem that one only needs to know λ_2 in order to know the convergence rate. However, this is incorrect for several reasons. For example, the set of eigenvalues may be too dense for this approximation to work in the interesting range of t . The perhaps most serious problem, however, is that in general the transition is much too large and complicated to leave us any hope of determining the second eigenvalue.

In situations where the second eigenvalue can be calculated, this can sometimes indeed be made use of, in particular when the MC under study is *reversible*. Recall that a MC is said to be reversible if it satisfies the so called detailed balance equations:

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

for all i and j . In continuous time one replaces p_{ij} with q_{ij} . For a reversible MC, all eigenvalues are real, since the symmetric matrix

$$A = [\sqrt{\pi(i)/\pi(j)}p_{ij}]_{i,j \in S}$$

has the same eigenvalues as P with eigenvectors $[\sqrt{\pi(i)}\phi(i)]_{i \in S}$ when ϕ is a (right) eigenvector of P . The same goes for Q in the continuous time case, just replacing p_{ij} with q_{ij} .

Definition 1.3 *The relaxation time, τ_2 , for a reversible MC is given by $\tau_2 = 1/(1 - \lambda_2)$ in discrete time, and $\tau_2 = 1/\lambda_2$ in continuous time. In the continuous time case, λ_2 is the second smallest eigenvalue of $-Q$.*

Later on, we shall among other things, see that

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2 \leq \frac{1}{\pi_*} e^{-t/\tau_2}$$

where $\pi_* = \min_i \pi(i)$.

2 Coupling, strong stationary times and ad-hoc lower bounds

2.1 Techniques for lower bounds

Until very recently, there was no sophisticated general technique for finding lower bounds on τ_{mix} ; one had, and usually still has, to use an argument that works for the special case under study. The most common way to do this, is to find an event $A = A^n$ such that for $t = t^n$ one has

$$\mathbb{P}(X_t^n \in A) \rightarrow 1$$

and

$$\pi^n(A) \rightarrow 0$$

as $n \rightarrow \infty$. Then t is a lower bound for τ_{mix} .

There are also other ideas. See e.g. Aldous and Fill [3], chapter 4. For example, it is a general fact that $\tau_{\text{mix}} \geq \tau_2$. One can then e.g. determine τ_2 for the movement of a single card (in the case of card shuffling), thereby getting a lower bound for the whole MC.

Wilson's technique is a clever development of this fact. Section 3 is devoted to his technique. In this section, however, we will settle for ad-hoc lower bounds of the "classical" kind.

2.2 Techniques for upper bounds

The most basic and useful techniques for upper bounds are *coupling* and *strong stationary times*. Other, more advanced, techniques that we will see more of later on, are Nash inequalities, log-Sobolev inequalities and relative entropy technique. Another method uses so called *evolving sets*, see Morris and Peres [16]

Coupling. Let $\{Y_t\}$ be a MC with the same transition matrix (generator) as $\{X_t\}$, but starting from stationarity, i.e. $\mathbb{P}(Y_0 \in \cdot) = \pi$. Suppose also that the updates of the two MC's are synchronized in such a way that whenever $t < t'$ and $X_t = Y_t$, then $X_{t'} = Y_{t'}$. (Obviously, such a synchronization can always be implemented without tampering with the marginal distributions of the two chains.) Let

$$T = \inf\{t : X_t = Y_t\}.$$

The stopping time T is called the *coupling time*. The following important inequality is called the *coupling inequality*:

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \leq \mathbb{P}(T > t).$$

Proof. For any $A \subseteq S$,

$$\begin{aligned} \mathbb{P}(X_t \in A) - \pi(A) &= \mathbb{P}(X_t \in A) - \mathbb{P}(Y_t \in A) \\ &= \mathbb{P}(X_t \in A, T \leq t) - \mathbb{P}(Y_t \in A, T \leq t) \\ &\quad + \mathbb{P}(X_t \in A, T > t) - \mathbb{P}(Y_t \in A, T > t) \\ &\leq \mathbb{P}(T > t) \end{aligned}$$

where the inequality follows from that the events of the first two terms are the same. \square

Strong stationary times. Assume that T is a stopping time such that $\mathbb{P}(X_T \in \cdot) = \pi$ and X_T is independent of T . Then T is called a *strong stationary time*. Again, the coupling inequality holds:

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \leq \mathbb{P}(T > t).$$

Proof. In the discrete time case

$$\begin{aligned} \mathbb{P}(X_t \in A) &= \sum_{s=1}^t \mathbb{P}(X_t \in A | T = s) \mathbb{P}(T = s) + \mathbb{P}(X_t \in A, T > t) \\ &= \pi(A) \mathbb{P}(T \leq t) + \mathbb{P}(X_t \in A, T > t) \\ &\leq \pi(A) + \mathbb{P}(T > t). \end{aligned}$$

In the continuous time case, replace the sum with an integral from 0 to t . \square

The independence assumption is important. It is easy to construct examples of uniform, but useless, stopping times. Consider e.g. simple random walk on the triangle K_3 and let $T = 1$ with probability $1/3$ and $T = 2$ with probability $2/3$. Clearly X_T is uniform, but equally clearly, the coupling inequality is not satisfied.

Random walks on groups. Suppose that the state space S is a group. The MC $\{X_t\}$ is then called a random walk on S if there is a probability measure ν on S such that

$$p_{ij} = \nu(i^{-1}j).$$

Then all column sums of P are 1 and hence $\pi = U$, the uniform distribution. Most card shuffling chains we shall encounter are random walks on the symmetric group, but in some cases, we will study card shuffling MC's, where the updating permutation has a distribution that depends on the state the MC is in.

A useful general observation for random walks on groups is that the time-reversed MC, $\{\hat{X}_t\}$, converges to uniformity exactly as quickly as the $\{X_t\}$ itself. This follows from the simple fact that the updating measure $\hat{\nu}$ of the time-reversed process is given by $\hat{\nu}(g) = \nu(g^{-1})$, so that $\mathbb{P}(\hat{X}_t \in A) = \mathbb{P}(X_t \in A^{-1})$.

We are now ready to attack some card shuffling chains. Always when working with such MC's, it is assumed that $X_0 = id = (1\ 2\ 3\ \dots\ n)$.

2.3 The top-to-random shuffle

For each step of the shuffle, take the top card and insert it at a uniformly chosen random position. Formally, this is a random walk on S_n (assuming that there are n cards) and the updating measure ν is given by

$$\nu(k\ k-1\ k-2\ \dots\ 1) = \frac{1}{n}, \quad k = 1, 2, \dots, n.$$

The following result was proved by Aldous and Diaconis [2]:

Theorem 2.1 *A threshold for the top-to-random shuffle is given by*

$$t = n \log n.$$

Proof. Let $T_0 = 0$ and let T_1 be the first time that the top card gets inserted into position n . Then inductively define T_2, T_3, \dots, T_n by letting T_k be the first

time after T_{k-1} that the top card gets inserted into either position $n, n-1, \dots$, or $n-k+1$.

Then $T = T_n$ is obviously a strong stationary time. To estimate $\mathbb{P}(T > t)$, write $T = \sum_{i=1}^n V_i$, where $V_i = T_i - T_{i-1}$. The V_i 's are independent and V_i is geometrically distributed with parameter i/n . Thus $\mathbb{E}V_i = n/i$ and $\mathbb{V}arV_i = (n-i)/n^2 \leq n^2/i^2$. Hence $\mathbb{E}T = n \sum_{i=1}^n \frac{1}{i} = (1+o(1))n \log n$ and $\mathbb{V}arT \leq Cn^2$, where $C = \sum_i i_{-2} < \infty$. Now, for any $a > 0$, by Chebyshev's inequality,

$$\mathbb{P}(T > (1+a)n \log n) \leq \mathbb{P}(|T - \mathbb{E}T| > an \log n) \rightarrow 0.$$

Hence by the coupling inequality, $(1+a)n \log n$ is an upper bound for τ_{mix} .

For a lower bound, let A be the event that card n is among the bottom $n/\log n$ cards (assuming for simplicity of notation that $n/\log n$ is an integer). Then $\pi(A) = 1/\log n \rightarrow 0$. However

$$\mathbb{P}(X_t \in A) \geq \mathbb{P}(T_{n/\log n} > t)$$

and since $T_{n/\log n}$ has expectation $(1+o(1))n(\log n - \log \log n) = (1+o(1))n \log n$ and variance bounded by Cn^2 , we get by an analogous use of Chebyshev's inequality that $\mathbb{P}(X_t \in A) \rightarrow 1$ when $t = (1-a)n \log n$. This shows that $(1-a)n \log n$ is a lower bound for τ_{mix} . \square

2.4 The random transpositions shuffle

Here, the right and left hand independently each pick a uniformly chosen card. Then, unless both hands have chosen the same card, the two cards are interchanged. Formally, this means that the updating measure given by $\nu(id) = 1/n$ and $\nu(ij) = 2/n^2$ for $i \neq j$. The following sharp result was proven by Matthews [12]:

Theorem 2.2 *A threshold for the random transpositions shuffle is given by*

$$t = \frac{1}{2}n \log n.$$

Proof. Let us start with a lower bound $(1-a)\frac{1}{2}n \log n$. For this, let A be the event that at least $\log n$ cards are in their starting positions. Since the expected number of such cards at stationarity is 1, it follows from Markov's inequality that $\pi(A) \rightarrow 0$.

However, $\mathbb{P}(X_t \in A)$ is bounded from below by the probability that at least $\log n$ cards are untouched by time t . Since at most two cards are touched at each step, it follows by mimicking the arguments for the top-to-random shuffle, that at time $t = (1 - a)\frac{1}{2}n \log n$, $\mathbb{P}(X_t \in A) \rightarrow 1$.

For an upper bound $(1 + a)\frac{1}{2}n \log n$, we define a strong stationary time via the following *marking procedure*. Start with all cards unmarked and then mark cards according to the two rules:

- (1) If both cards chosen are unmarked, then mark the card chosen by the left hand. This also applies when both hands chose the same unmarked card.
- (2) If only one of the hands chose an unmarked card, then mark this card.

Then, letting T be the first time that all cards are marked, we have that T is a strong stationary time. This is proved by induction, noting that at each time a new card is marked, that card is in a uniform position among the cards that are marked so far.

Since each step involves two randomly chosen cards, one for the left hand and one for the right hand, the now familiar Chebyshev argument shows that

$$P(T > (1 + a)n \log n) \rightarrow 0.$$

□

Our next shuffle is much slower than two we have seen so far.

2.5 The transposing neighbors shuffle

At each step, do nothing with probability $1/2$ and with probability $1/2$ pick a position uniformly at random from $\{1, 2, \dots, n - 1\}$ and switch the card there with the card next to and below it. I.e.

$$\nu(id) = \frac{1}{2}, \quad \nu(i, i + 1) = \frac{1}{2(n - 1)}, \quad i = 1, \dots, n - 1.$$

Theorem 2.3 (Aldous [1]) *There is a constant $K < \infty$ (independent of n) such that $\tau_{\text{mix}} \leq Kn^3 \log n$. There is also a $c > 0$ such that $\tau_{\text{mix}} \geq cn^3$.*

Proof. A coupling will be used for the upper bound. Let $\{X_t\}$ be the ordinary deck and let $\{Y_t\}$ be another deck with the same transition rules, started from

uniformity. Say that a card i is coupled at time t if $X_t(i) = Y_t(i)$, i.e. if i is in the same position in the two decks at time t .

One way of describing the shuffling rules, is to say that one makes two uniformly random choices, the first being a position in $\{1, \dots, n-1\}$ and the second whether or not to swith the card at the chosen position with its neighbor.

Couple the transitions of the two decks by using the same first choice, i , for the Y -process as for the X -process. Then, if either the card at position i or the card at position $i+1$ is coupled, make the same second choice for the two decks. If not, make the opposite second choices for the Y -process as for the X -process.

Now clearly the marginal updating distributions are correct, and a given card in one deck cannot pass its copy in the other deck without being coupled. Hence, if T is the first time that every card has visited position 1 in both decks, then the two decks are coupled at time T and $\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \leq \mathbb{P}(T > t)$.

Fix a card i and let $Z_t = X_t(i)$ be i 's position at time t . Then $\{Z_t\}$ is a random walk, reflected in 1 and n , whose steps are -1 or 1 with probability $1/2(n-1)$ each and 0 with probability $1 - 1/(n-1)$. By Donsker's Theorem (see Durrett [9], p. 365),

$$\left\{\frac{1}{n}Z_{sn^3}\right\}_{s \geq 0} \rightarrow \{B_s\}_{s \geq 0}$$

where B_s is standard Brownian motion, reflected at 0 and 1, and convergence is in distribution. It is a well known result (see Durrett [9] again) that the probability that $\{B_s\}$ has not been reflected at 0 before time t is bounded by $Ce^{-\alpha t}$, where $C < \infty$ and $\alpha > 0$. Letting T_i be the first time that card i has visited 1 in both decks, this entails that

$$\mathbb{P}(T_i > b(n)n^3) \leq 2Ce^{-\alpha b(n)}$$

where the factor 2 refers to that we are considering two decks. With $K = 2/\alpha$ and $b(n) = K \log n$, the rhs is bounded by $2Cn^{-2}$. Finally, since $T = \max_i T_i$,

$$\mathbb{P}(T > Kn^3 \log n) \leq n\mathbb{P}(T_i > Kn^3 \log n) \leq \frac{2C}{n}.$$

For the lower bound, note that after time cn^3 , the position of card i has variance $(1+o(1))cn^2$, so for c small enough, this is far from the uniform distribution which has variance $(1+o(1))n^2/12$. We omit some details here, since this lower bound will be improved upon in the next section. \square

2.6 The overhand shuffle

This is our first model of a "real" shuffle, in the sense that people actually use it to mix a deck of cards.

The model has a probability parameter $p \in (0, 1)$. Given p , each of the $n - 1$ slots between successive cards is *marked* with probability p , independently of other slots. Then, given the marks, each pack of cards between two successive marks, is reversed. (To be true, a step of the shuffle described gives the reversed deck of what one gets from an actual overhand shuffle. Of course, every second time, one gets back the "real" deck. It is equally obvious that the convergence rate is not affected.)

Example. 123|45|6|789 results in 321546987. □

The following result was established by Pemantle [17].

Theorem 2.4 *There exists a constant $C = C(p) < \infty$ such that for the overhand shuffle, $\tau_{\text{mix}} \leq Cn^2 \log n$. On the other hand, τ_{mix} is at least of order n^2 .*

Proof. The lower bound part will only be sketched, since it will be improved upon later on. Here we simply note that each card makes a random walk which is very close to symmetric with step size variance of order 1. (Only very close since the boundaries of the deck have a slight repelling effect.) Hence, disregarding the fact that the random walk is not perfectly symmetric, the CLT entails that mixing time of a single card is of order n^2 .

For the upper bound, we use a coupling. Let X_t and Y_t be the original and the stationary deck respectively, as usual. For simplicity, let us consider a circular deck, i.e. we get an extra slot between positions 1 and n and cards may "come around the corner". (This means that a single cards makes a perfectly summetric random walk on \mathbb{Z}_n and so the lower bound argument above is complete as it stands.) Now we couple the decks by using the same nearest marked slots around each coupled card. More precisely, this means that when we flip coins for determining the status of a slot, we start with the slots closest to coupled cards.

After having found the nearest marked slots, the status of the rest of slots are determined independently for the two decks.

Now consider a single card, i , and let $\{V_t\}$ and $\{W_t\}$ denote its random walk in the the two decks respectively. Let us again simplify a bit, by assuming that $p = 1/2$, letting general p be an exercise for the reader. Then the step size distribution of these random walks is that the step size is $j \pmod n$ with probability

$\frac{1}{3}(\frac{1}{2})^{|j|}$, $j \in \mathbb{Z}$. Hence the step size distribution has step size variance of order 1. From this one can show that there is a constant c , such that within time $cn^2 - 1$, the distance between V_t and W_t , will at least once be at most 1. (Proving this properly is more complicated than one would think, since the two random walks are slightly dependent and this dependence is non-negligible when they are close to each other.) Now, given that the distance between card i 's position in the two decks is 1 at some time, then the probability that i gets coupled in the next step is at least $(1/2)^6$. This entails that the probability that i is coupled at time cn^2 is at least $(1/2)^7$. This means that there is a constant C such that at time $Cn^2 \log n$, card i is coupled with probability at least $1 - 1/n^2$, and thus

$$\mathbb{P}(T > Cn^2 \log n) \leq \frac{1}{n}.$$

□

2.7 The riffle shuffle

This is a model for the most commonly used shuffle, when people shuffle a real deck of cards, where the deck is divided in two packs of roughly $n/2$ cards each, and then these two packs are interleaved. Bayer and Diaconis [5] showed that for this shuffle, $\frac{3}{2} \log_2 n$ is a threshold and, via a detailed analysis, that for $n = 52$, 7 shuffles is sufficient to bring $\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV}$ down to about 0.3. This result was considered of such common interest, that it made it to the front page of the New York Times!

Here we shall only make the immediate observation that $\log_2 n$ is a lower bound and give the surprisingly simple argument of Aldous and Diaconis [2], that $2 \log_2 n$ is an upper bound.

First we describe the model. Each step of the shuffle is a function of a sequence $\xi_1, \xi_2, \dots, \xi_n$ of iid random variables where $\mathbb{P}(\xi_i = 0) = \mathbb{P}(\xi_i = 1) = 1/2$. Such a sequence is interpreted in the following way. The number of 0's of the sequence tells us how many cards that goes in the top pack when the deck is divided in two. Then the order of the 0's and 1's tells us in what order the cards fall during the interleaving part, where, of course, a 0 corresponds to a card from the top pack and a 1 to a card from the bottom pack.

Example. If the $n = 10$ cards are in the order

abcdefghij

and the random sequence is 0110111010, then the order of the cards after this step of the shuffle is

$$aefbg hicjd.$$

□

Since there are at most 2^n (in fact exactly $2^n - n - 1$) possible outcomes of each step of the shuffle and $n! > n^n/e^{2n} = 2^{(1+o(1))n \log_2 n}$ by Stirling's formula, we get that $\tau_{\text{mix}} \geq (1 + o(1)) \log_2 n$.

For the upper bound, we will use the fact that the time-reversed MC has the same mixing time as the original MC. One step of the time-reversed chain “undoes” the riffle shuffle, i.e. for the given sequence of 0's and 1's, it takes the cards at positions corresponding to a 0 and puts them on top of the other cards, without changing their relative order.

Example. The sequence 0011010001 takes $abcdefg hij$ to $abegh icdfj$. □

Let us say that a card in a position corresponding to a 0 is *marked* with a 0 in that step, and similarly with a 1. Then in each step of the shuffle, each card gets a new mark, so after k shuffles, each card has an iid sequence of k marks. Now make a coupling with a stationary deck by always letting each card get the same mark in the two decks.

Let T be the coupling time and let \hat{T} be the first time that all cards have distinct sequences of marks. When all cards have distinct mark sequences, the positions of the cards can be read off from the marks. This is so, since when comparing two cards, we have that the card which has a 0-mark at the latest time that the marks differ, is higher up in the deck. Hence $T \leq \hat{T}$. Now, since there are $\binom{n}{2}$ pairs of cards and since the probability that two given cards have the same mark sequences after k shuffles is $(1/2)^k$,

$$\mathbb{P}(\hat{T} > k) \leq \binom{n}{2} 2^{-k} < \frac{n^2}{2} 2^{-k}$$

which is less than $1/4$ when $k = 2 \log_2 n + 1$.

Remarks. The model for the riffle shuffle is fairly realistic, but it deviates from how people really shuffle in the following significant aspect. Since the 0/1-sequence that models a step is iid, it follows that for each stage of the interleaving part of the shuffle that, if there are k_r cards left in the right hand and k_l cards in the left hand, then the probability that the next card is dropped from the right hand is $k_r/(k_r + k_l)$, no matter which hand dropped the card before.

In practice, people have a slight tendency to drop the next card from the other hand. Thus the 0/1-sequence which on average changes from 0 to 1 or vice versa $n/2$ times, should, in order to be more realistic, be replaced by a dependent sequence that changes, say, $3n/4$ times. To analyse such a model is a difficult open problem. Note that the coupling above breaks down in this case.

The Thorp shuffle, to which a later section is devoted, is another variant of a finer riffle shuffle.

3 Wilson's technique for lower bounds

Wilson's technique is the first systematic technique for finding test functions that provide tight lower bounds on the mixing time. The technique tends to work well in situations where the motion of a single card is easily analyzed and when cards move fairly independently. Sometimes one can alternatively use some other simple MC's embedded in the MC under study. The method was introduced in [19] and [20]. The presentation here is based on those papers and on the extensions presented in [11] and [10].

3.1 The basic setup

Let as usual $\{X_t\}_{t=0}^{\infty}$ denote the (irreducible aperiodic) MC under study. The test function used in basic Wilson technique is a (right) eigenvector, ϕ , of the transition matrix, i.e. a function on the state space S such that

$$\mathbb{E}[\phi(X_{t+1})|X_t] = \lambda\phi(X_t)$$

almost surely. For simplicity, we will start with assuming that ϕ corresponds to a real eigenvalue $\lambda = 1 - \gamma$, so that ϕ itself is also real-valued. (This holds automatically e.g. when $\{X_t\}$ is reversible.) We also make the very mild assumption that $0 < \gamma < 1/2$. Let

$$R = \max_{s \in S} \mathbb{E}[(\phi(X_{t+1}) - \phi(X_t))^2 | X_t = s]$$

and assume that $\phi(X_0) > 0$. Fix $a > 0$ and let

$$T = \frac{\log \phi(X_0) - \frac{1}{2} \log \frac{4R}{\gamma^a}}{-\log(1 - \gamma)}.$$

The following result is the key to Wilson's technique.

Theorem 3.1 (Key Theorem) For all $t \leq T$,

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \geq 1 - a.$$

Usually when applying the key theorem, one takes a to be some specific value of at most $3/4$, e.g. $a = 1/2$ so that the ratio in the second term of the numerator of T becomes $8R/\gamma$. In almost all cases we will have $\gamma = \gamma(n) = o(1)$, so that, with $a = 1/2$,

$$T = (1 + o(1))\gamma^{-1}(\log \phi(X_0) - \frac{1}{2} \log(8R/\gamma)).$$

Proof. Let X be a random variable on S distributed according to π . By induction

$$\mathbb{E}\phi(X_t) = (1 - \gamma)^t \phi(X_0).$$

Hence, letting $t \rightarrow \infty$, we have $\mathbb{E}\phi(X) = 0$.

Write $\Delta\phi = \phi(X_{t+1}) - \phi(X_t)$. Since $\phi(X_{t+1})^2 = \phi(X_t)^2 + 2\phi(X_t)\Delta\phi + (\Delta\phi)^2$ and $\mathbb{E}[\Delta\phi|X_t] = -\gamma\phi(X_t)$,

$$\mathbb{E}[\phi(X_{t+1})^2|X_t] \leq (1 - 2\gamma)\phi(X_t)^2 + R.$$

Using this inductively, utilizing the assumption that $\gamma < 1/2$, it follows that for all t ,

$$\mathbb{E}[\phi(X_t)^2] \leq (1 - 2\gamma)^t \phi(X_0)^2 + \frac{R}{2\gamma}.$$

This entails

$$\begin{aligned} \text{Var}\phi(X_t) &= \mathbb{E}[\phi(X_t)^2] - (\mathbb{E}\phi(X_t))^2 \\ &\leq \left((1 - 2\gamma)^t - (1 - \gamma)^{2t} \right) \phi(X_0)^2 + \frac{R}{2\gamma} \\ &\leq \frac{R}{2\gamma}. \end{aligned}$$

Hence, by Chebyshev's inequality

$$\mathbb{P}\left(|\phi(X_t) - \mathbb{E}\phi(X_t)| \geq \sqrt{\frac{R}{\gamma a}}\right) \leq \frac{a}{2}$$

and by letting $t \rightarrow \infty$ and keeping in mind that $\mathbb{E}\phi(X) = 0$,

$$\mathbb{P}(|\phi(X)| \geq \sqrt{\frac{R}{\gamma a}}) \leq \frac{a}{2}.$$

This implies that if t is such that $\mathbb{E}\phi(X_t) \geq \sqrt{\frac{4R}{\gamma a}}$, then

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \geq 1 - a.$$

Finally we observe that

$$T = -\log_{1-\gamma} \frac{\phi(X_0)}{\sqrt{(4R)/(\gamma a)}}$$

and since $\mathbb{E}\phi(X_t) = (1 - \gamma)^t \phi(X_0)$, it follows that this holds for all $t \leq T$. \square

Example. Transposing neighbors, continued. Let us go back to the transposing neighbors shuffle, for which we established above that τ_{mix} is at least of order n^3 and at most of order $n^3 \log n$. We will now establish that $n^3 \log n$ is the correct order of τ_{mix} , following the analysis in [20]. For convenience we will make a slight adjustment of the shuffle and use the circular deck convention and skip the holding probability of $1/2$. Thus the updating measure is now given by

$$\nu(id) = \nu(12) = \nu(23) = \dots = \nu(n-1n) = \nu(n1) = \frac{1}{n+1}.$$

Let $Z_t = Z_t^i$ be the position of card i at time t . Then

$$\begin{aligned} \mathbb{E}[\cos \frac{2\pi Z_{t+1}}{n} | X_t] &= (1 - \frac{2}{n+1}) \cos \frac{2\pi Z_t}{n} + \frac{1}{n+1} \cos \frac{2\pi(Z_t + 1)}{n} \\ &\quad + \frac{1}{n+1} \cos \frac{2\pi(Z_t - 1)}{n} \\ &= (1 - \frac{2}{n+1}) \cos \frac{2\pi Z_t}{n} + \frac{2}{n+1} \cos \frac{2\pi Z_t}{n} \cos \frac{2\pi}{n} \\ &= (1 - \gamma) \cos \frac{2\pi Z_t}{n} \end{aligned}$$

where $\gamma = 2(1 - \cos(2\pi/n))/(n+1) = (1 + o(1))4\pi^2/n^3$. Hence $\cos(2\pi Z_t/n)$ is an eigenvector of the eigenvalue $1 - \gamma$. By linearity of expectation, this also holds for

$$\phi(X_t) = \sum_{i=1}^m \cos \frac{2\pi Z_t^i}{n}$$

where $m = \lfloor n/2 \rfloor$.

Let us now apply the key theorem. Since the starting order of the deck is irrelevant for the mixing time, we can start with the cards in an order that maximizes $\phi(X_0)$. Then $\phi(X_0) = \sum_{i=1}^m \sin \frac{2\pi i}{n} = Cn$ where $C = (1 + o(1))1/\pi$. Since $\phi(X_{t+1})$ and $\phi(X_t)$ differ only if step $t+1$ consists of a transposition of one of the cards $1, \dots, m$ and one of the other cards,

$$R \leq \left(\cos\left(\frac{\pi}{2} - \frac{2\pi}{n}\right) \right)^2 \leq \left(\frac{2\pi}{n} \right)^2.$$

Plugging all this into the key theorem (with $a = 1/2$), we get

$$\begin{aligned} \tau_{\text{mix}} \geq T &= \frac{n^3}{4\pi^2} \left(\log(Cn) - \frac{1}{2} \log \frac{32\pi^2/n^2}{4\pi^2/n^3} \right) \\ &= (1 + o(1)) \frac{1}{8\pi^2} n^3 \log n. \end{aligned}$$

□

Remark. It is not an essential problem to remove the circular deck assumption. When doing so, the given $(1 - \gamma, \phi)$ is not an eigenvalue/eigenvector pair, but only very close. Quantifying the error term in the conditional expectation, one can then use the extension of the key theorem in [11]. Alternatively one can use Newton's method to find a true eigenvalue/eigenvector pair close to $(1 - \gamma, \phi)$.

Example. Random transpositions. We already know that the threshold for random transpositions is $\frac{1}{2}n \log n$, but let us for illustration derive a lower bound via Wilson's technique.

Let again $Z_t = Z_t^i$ denote the position of card i at time t . Then, if n is even,

$$\begin{aligned} \mathbb{E}[(-1)^{Z_{t+1}} | X_t] &= \left(1 - \frac{2}{n}\right) (-1)^{Z_t} + \frac{1}{n} (-1)^{Z_t} - \frac{1}{n} (-1)^{Z_t} \\ &= \left(1 - \frac{2}{n}\right) (-1)^{Z_t} \end{aligned}$$

where the terms in the first equality correspond to the card i not being touched, card i touched and moved an even number of steps and card i touched and moved an odd number of steps, respectively. If n is odd, this equality is not exact, but the error is of order n^{-2} and will hence vanish in the limit.

Now, with $m = \lfloor n/2 \rfloor$, let

$$\phi(X_t) = \sum_{i=1}^m (-1)^{Z_t^i}.$$

Then $\mathbb{E}[\phi(X_{t+1})|X_t] = (1 - 2/n)\phi(X_t)$, so we may apply the key theorem with $\gamma = 2/n$. Starting with cards $1, \dots, m$ in even positions, we have $\phi(X_0) = m$ and since at most two cards move at a given step, $R \leq 2^2 = 4$. Hence

$$\tau_{\text{mix}} \leq T = (1 + o(1))\frac{n}{2} \left(\log m - \frac{1}{2} \log(16n) \right) = (1 + o(1))\frac{1}{4}n \log n.$$

□

Remark. Note that for random transpositions, Wilson's technique gave a lower bound which deviates from the true threshold by a factor 2. For transposing neighbors, it is also believed that the true mixing time is twice as large as the one we got from Wilson's technique and this deviation by a factor two seems to very common when applying Wilson. It would be interesting if one could find a general explanation for this fact and a remedy for it.

Example. The overhand shuffle, continued. We know from above that for the overhand shuffle, τ_{mix} is at most of order $n^2 \log n$. Here we will via Wilson's technique, following [11], establish a lower bound of the same order.

Assume to begin with, that a slot is marked with probability $1/2$. Use also the circular deck convention. (This can be removed along the lines indicated in a remark above, see [11].) As we observed in the previous section, each card makes a symmetric random walk on \mathbb{Z}_n , where the probability of moving j steps equals $(1/2)^{|j|}/3$, $j \in \mathbb{Z}$. Thus

$$\begin{aligned} \mathbb{E}\left[\cos \frac{2\pi Z_{t+1}}{n} | X_t\right] &= \frac{1}{3} \cos \frac{2\pi Z_t}{n} \\ &+ \frac{1}{3} \sum_{j=1}^{\infty} \frac{1}{2^j} \left(\cos \frac{2\pi(Z_t + j)}{n} + \cos \frac{2\pi(Z_t - j)}{n} \right) \\ &= \frac{1}{3} \cos \frac{2\pi Z_t}{n} \left(1 + 2 \sum_{j=1}^{\infty} \frac{\cos(2\pi j/n)}{2^j} \right) \\ &= (1 - \gamma) \cos \frac{2\pi Z_t}{n} \end{aligned}$$

where $\gamma = (1 + o(1))8\pi^2/n^2$; the value of γ follows from Taylor's formula:

$$\cos \frac{2\pi j}{n} = 1 - \frac{2\pi^2 j^2}{n^2} + O((j/n)^4)$$

and that $\sum_{j=1}^{\infty} j^2/2^j = 6$. Letting, with $m = \lfloor n/2 \rfloor$,

$$\phi(X_t) = \sum_{i=1}^m \cos \frac{2\pi Z_t^i}{n}$$

we have that ϕ is a eigenvector corresponding to $1 - \gamma$. Starting with the cards in a suitable order, we have, as for transposing neighbors, $\phi(X_0) = Cn$.

Now we must bound R . Let

$$Y_i = \cos \frac{2\pi Z_{t+1}^i}{n} - \cos \frac{2\pi Z_t^i}{n}.$$

Then $(\phi(X_{t+1}) - \phi(X_t))^2 = (\sum_{i=1}^m Y_i)^2$. Hence

$$\mathbb{E}[(\phi(X_{t+1}) - \phi(X_t))^2 | X_t] = \sum_i \sum_j \mathbb{E}[Y_i Y_j | X_t].$$

Consider a given (i, j) . Let Y'_i and Y'_j be distributed as Y_i and Y_j but independent. Then we can make a coupling such that $(Y_i, Y_j) = (Y'_i, Y'_j)$ on the event, A , that there are at least two marked slots between cards i and j for update $t + 1$ of the shuffle. If the distance between card i and card j at time t is at least, say, $10 \log n$, then $\mathbb{P}(A) = 1 - o(1/n^5)$ and hence

$$\mathbb{E}[Y_i Y_j | X_t] = \mathbb{E}[Y'_i | X_t] \mathbb{E}[Y'_j | X_t] + o(n^{-5}) = O(n^{-4}).$$

For other i and j ,

$$\mathbb{E}[Y_i Y_j | X_t] \leq \max_{s \in S} \mathbb{E}[Y_i^2 | X_t = s] = \frac{2}{3} \sum_{j=1}^{\infty} \frac{1}{2^j} \left(\sin \frac{2\pi j}{n} \right)^j < 50n^{-2}.$$

Hence

$$\begin{aligned} \sum_i \sum_j \mathbb{E}[Y_i Y_j | X_t] &< m^2 O(n^{-4}) + 20n \log n \cdot 50n^{-2} \\ &= O(n^{-2}) + 1000n^{-1} \log n. \end{aligned}$$

Thus, we may use $R = (1 + o(1))1000n^{-1} \log n$ in the key theorem, together with $\phi(X_0) = Cn$ and $\gamma = (1 + o(1))8\pi^2/n^2$ and get

$$\tau_{\text{mix}} \geq (1 + o(1)) \frac{1}{16\pi^2} n^2 \log n.$$

In the more general case where slots are marked with probability p (we just did the case $p = 1/2$), the step size distribution for the random walk described by single card is that one moves distance j with probability $p(1-p)^{|j|}/(2-p)$, $j \in \mathbb{Z}$. Recalculating the above with this step size distribution gives

$$\gamma = (1 + o(1))4\pi^2 \frac{1-p^2}{p^2(2-p)} n^{-2}$$

which in the end leads to

$$\tau_{\text{mix}} \geq \frac{p^2(2-p)}{8\pi^2(1-p^2)} n^2 \log n.$$

□

3.2 The move-to-front rule and the first extension

The basic setup for Wilson's technique has some restrictions that we now set out to relieve. One restriction is the assumption that eigenvalues be real. For a non-reversible MC, this is usually not the case, and the right order of mixing is connected to how close eigenvalues are to 1 in absolute value, rather than the distance in the complex plane between the eigenvalues and the point 1. Thus we would like a modification of the setup that enables us to confirm this. We will come back to this in the next subsection.

Another restriction is that in the basic setup, only one eigenvalue is used. In some situations, e.g. when different cards have different behavior, one may have to use a combination of different eigenvector/eigenvalue pairs. This is the issue that will now be addressed.

The first extension. Suppose that $\phi_1, \phi_2, \dots, \phi_k$ are eigenvectors of the transition matrix corresponding to the eigenvalues $\gamma_1, \gamma_2, \dots, \gamma_k$. Assume that

$$\text{Cov}(\phi_i(X_t), \phi_j(X_t)) \leq 0, \quad i \neq j$$

and that for all i ,

$$\max_{s \in S} \mathbb{E}[(\phi_i(X_{t+1}) - \phi_i(X_t))^2 | X_t = s] \leq \gamma_i.$$

(In fact, the last condition is just a matter of scaling, since an eigenvector can be scaled as one likes. What will be crucial is the relation between γ_i and $\phi(X_0)$ after

the scaling.) Now, plug in ϕ_i in the beginning of the proof of the key theorem to get

$$\text{Var}\phi_i(X_t) \leq \frac{1}{2}.$$

Let $\alpha_i, i = 1, \dots, k$, be positive numbers (or *weights*) such that $\sum_i \alpha_i^2 = 1$ and let

$$\phi(X_t) = \sum_i \alpha_i \phi_i(X_t).$$

Then, by the negative correlation condition,

$$\text{Var}\phi(X_t) \leq \frac{1}{2}.$$

Now, plugging in ϕ in the second part of the proof of the key theorem yields that $\|\mathbb{P}(X_t \in \cdot) - \pi\|_{TV} \geq 1 - a$ as long as $\mathbb{E}\phi(X_t) \geq \sqrt{4/a}$. Since

$$\mathbb{E}\phi(X_t) = \sum_i \alpha_i (1 - \gamma_i)^t \phi_i(X_0)$$

this holds e.g. when, for all i ,

$$\alpha_i (1 - \gamma_i)^t \phi_i(X_0) \geq \alpha_i^2 \sqrt{4/a},$$

i.e. when $(1 - \gamma_i)^t \phi_i(X_0) \geq \alpha_i \sqrt{4/a}$. This gives us

$$\tau_{\text{mix}} \geq \min_i \frac{\log \phi_i(X_0) - \frac{1}{2} \log(4\alpha_i^2/a)}{-\log(1 - \gamma_i)}.$$

Typically $\gamma_i = o(1)$, $\alpha_i = o(1)$ and $\phi_i(X_0) = O(1)$ for all i , and then we get

$$\tau_{\text{mix}} \geq (1 + o(1)) \min_i \gamma_i^{-1} \log \alpha_i^{-1}.$$

When using this result, the idea is of course to use weights α_i such that the ratios in the minimum are equal or approximately equal.

Example. The move-to-front rule. This is a random-to-top shuffle, where different cards have different probabilities of being picked. More precisely, to each card, i , a probability p_i is associated. At each step of the shuffle, a card is chosen at random according to these probabilities and the chosen card is moved from its present position to the top of the deck. (The standard metaphor for the

items is “books” rather than “cards”, but in practice, the “books” are usually data files.)

Note that this is not a random walk on S_n in the above sense, since the probability for a given updating permutation depends on the present order of the cards. Hence the stationary distribution is not necessarily uniform. Indeed, the stationary distribution is given by

$$\pi((c_1 c_2 c_3 \dots c_n)^{-1}) = \prod_{k=1}^n \frac{p_{c_k}}{1 - \sum_{i=1}^{k-1} p_i},$$

see e.g. [18].

Without loss of generality, we assume that $p_1 \leq p_2 \leq \dots \leq p_n$. For this shuffle, the mixing time may depend on the starting state. In such cases one is usually interested in the slowest possible mixing. Not surprisingly, a starting order which turns out to achieve the highest possible order of mixing, is to start, as we do, with the cards in the order $123 \dots n$.

To apply the extension of the key theorem, let, for each odd i ,

$$\phi_i(X_t) = \begin{cases} \frac{p_{i+1}}{p_i + p_{i+1}}, & Z_t^i < Z_t^{i+1} \\ -\frac{p_i}{p_i + p_{i+1}}, & Z_t^i > Z_t^{i+1} \end{cases}$$

Then it is readily checked, considering the two cases $Z_t^i < Z_t^{i+1}$ and $Z_t^i > Z_t^{i+1}$ separately, that $\mathbb{E}[\phi_i(X_{t+1})|X_t] = (1 - p_i - p_{i+1})\phi_i(X_t)$. Hence $\gamma_i = p_i + p_{i+1}$. Since $|\phi_i(X_{t+1}) - \phi_i(X_t)| \leq 1$ and the probability that ϕ_i changes at all in one step is bounded by p_{i+1} ,

$$\mathbb{E}[(\phi_i(X_{t+1}) - \phi_i(X_t))^2|X_t] \leq p_{i+1} \leq \gamma_i.$$

We must also check that ϕ_i and ϕ_j , $i \neq j$, are not positively correlated. However, given that card k or card $k + 1$ has at least once been touched by time t , the expectation of $\phi_k(X_t)$ is 0. Hence

$$\mathbb{E}\phi_k(X_t) = (1 - \gamma_k)^t \phi_k(X_0)$$

and similarly

$$\mathbb{E}[\phi_i(X_t)\phi_j(X_t)] = (1 - \gamma_i - \gamma_j)^t \phi_i(X_0)\phi_j(X_0)$$

from which the required negative correlation is immediate. Now let

$$\phi(X_t) = \sum_i \alpha_i \phi_i(X_t)$$

where the sum is over odd i 's and $\sum_i \alpha_i^2 = 1$. Apply the extension of the key theorem with $a = 1/2$ using that $\phi(X_0) \geq 1/2$ to get

$$\tau_{\text{mix}} \geq \min_i \frac{\frac{1}{2} \log(\alpha_i^{-2}/8) - \log 2}{-\log(1 - \gamma_i)}.$$

When $\gamma_i = o(1)$ and $\alpha_i = o(1)$, this becomes, as remarked above,

$$\tau_{\text{mix}} \geq (1 + o(1)) \min_i \gamma_i^{-1} \log \alpha_i^{-1}.$$

In order to optimize, we would like all $\gamma_i^{-1} \log \alpha_i^{-1}$ to be equal. This means to let α_i be proportional to $e^{-\gamma_i}$. However, in the special case, it is easier to pick the α_i 's only approximately optimal, via inspection of the p_i 's. We will now study four special cases.

Case A. Ordinary random-to-top, i.e. $p_i = 1/n$ for all i . Then $\gamma_i = 2/n$ and we can take $\alpha_i = (2/n)^{1/2}$. This gives

$$\tau_{\text{mix}} \geq (1 + o(1)) \frac{1}{4} n \log n.$$

Case B. Let $m = \lfloor n/2 \rfloor$ and let $p_i = 2/n(n+1)$ for $i \leq m$ and $p_i = 2/(n+1)$ for $i > m$. Then the corresponding γ_i 's are $(1 + o(1))4/n^2$ and $(1 + o(1))4/n$ respectively. The important cards are the first m cards, since these are most seldom touched. Taking $\alpha_i = (1 + o(1))2n^{-1/2}$, $i \leq m$ and $\alpha_i = 0$ otherwise, we get

$$\tau_{\text{mix}} \geq \frac{1}{8} n^2 \log n.$$

Case C. Let $p_i = (n+1-i)^{-1}/(\sum_{j=1}^n j^{-1}) = (1 + o(1))(n+1-i)^{-1}/\log n$. This gives $\gamma_i = (1 + o(1))2/(n+1-i) \log n$. Taking $\alpha_i = 2n^{-1/2}$ for $i \leq m$ and $\alpha_i = 0$ otherwise gives

$$\tau_{\text{mix}} \geq \frac{1}{8} n (\log n)^2.$$

Case D. Let $p_i = 2i/n(n+1)$. Here we need only put weight on the first few cards, e.g. $\alpha_1 = \alpha_3 = \dots = \alpha_{127} = 1/8$. This gives

$$\tau_{\text{mix}} \geq \gamma_{128}^{-1} \left(\frac{1}{2} \log 8 - \log 2 \right) = Cn^2$$

for a constant C . This turns out to be the correct order of mixing.

We would now like to match these lower bounds with good upper bounds. We do this via coupling. As usual $\{X_t\}$ denotes the MC under study and $\{Y_t\}$ a stationary copy of it. The coupling is very simple: always move the same card in the two decks. With T for the coupling time, we have that T must have occurred as soon as all but one card has been touched, in particular when all cards $2, 3, \dots, n$ have been touched. Hence

$$\mathbb{P}(T \geq t) \leq \sum_{i=2}^n (1 - p_i)^t.$$

Letting

$$\tau_u = \min\left\{t : \sum_{i=1}^n (1 - p_i)^t \leq \frac{1}{4}\right\}$$

we get $\tau_{\text{mix}} \leq \tau_u$. (One can also show that $\tau_{\text{mix}} \geq \tau_u/25$, see [10].) Let us now estimate τ_u for the above special cases.

Case A. $\sum_{i=2}^n (1 - p_i)^t = (n - 1)(1 - 1/n)^t$ which is less than $1/4$ when $t \geq \log(4(n - 1))/\log(n/(n - 1)) = (1 + o(1))n \log n$. Thus

$$\tau_{\text{mix}} \leq n \log n.$$

Case B. Here

$$\sum_{i=2}^n (1 - p_i)^t = (1 + o(1)) \frac{n}{2} \left(\left(1 - \frac{2}{n}\right)^t + \left(1 - \frac{2}{n^2}\right)^t \right)$$

which is less than $1/4$ for $t \geq (1 + o(1))\frac{1}{2}n^2 \log n$. Hence

$$\tau_{\text{mix}} \leq (1 + o(1))\frac{1}{2}n^2 \log n.$$

Case C. In this case, with $t = cn(\log n)^2$,

$$\sum_{i=2}^n (1 - p_i)^t = (1 + o(1)) \sum_{i=1}^{n-1} \left(1 - \frac{1}{i \log n}\right)^{cn(\log n)^2}$$

which clearly tends to 0 for $c > 1$ (and tends to 1 for $c < 1$, you may take it as an exercise to show this). Hence

$$\tau_{\text{mix}} \leq n(\log n)^2.$$

Case D. Taking $t = cn^2$, we get

$$\begin{aligned} \sum_{i=2}^n (1 - p_i)^t &= (1 + o(1)) \sum_{i=2}^n \left(1 - \frac{2i}{n^2}\right)^{cn^2} \\ &= (1 + o(1)) \sum_{i=2}^n e^{-2ci} \\ &= (1 + o(1)) \frac{e^{-4c}}{1 - e^{-2c}} \end{aligned}$$

which is less than $1/4$ when $c > \frac{1}{2} \log(8/(\sqrt{17} - 1))$. Hence

$$\tau_{\text{mix}} < 0.471n^2.$$

□

3.3 GR-shuffles and the second extension

Similar in spirit to the move-to-front rule are the shuffles where one instead of moving *card* i to the top with probability p_i , one moves the card in *position* i to the top with probability p_i . We will call such shuffles *Generalized Rudvalis shuffles* or, in short *GR-shuffles*, since they are generalization of the so called Rudvalis shuffle, where $p_{n-1} = p_n = 1/2$. They are also generalizations of the random-to-top shuffle (the time-reversal of top-to-random) for which $p_i = 1/n$ for all i .

Unlike the move-to-front rule, the GR-shuffles are random walks on the symmetric group. The updating measure ν is given by $\nu(1\ 2 \dots i) = p_i, i = 1, 2, \dots, n$. Thus the stationary distribution is uniform. They are however non-reversible, so when using Wilson's technique for lower bounds, we face the problem of dealing with complex eigenvalues/eigenvectors. Wilson [19] deals with this via an extension of the state space. Here we will take a slightly different route.

The second extension. Suppose that ϕ is an eigenvector of the transition matrix of the MC $\{X_t\}$, that corresponds to the eigenvalue $(1 - \gamma)e^{i\theta}$, where $\gamma \in (0, 1/2)$ and $\theta \in [0, \pi]$. Let

$$R \geq \max_{s \in S} \mathbb{E}[|e^{-i\theta} \phi(X_{t+1}) - \phi(X_t)|^2 | X_t = s].$$

Then $\tau_{\text{mix}} \geq 1 - a$ for $t \leq T$, where

$$T = \frac{\log |\phi(X_0)| - \frac{1}{2} \log(4R/\gamma a)}{-\log(1 - \gamma)}.$$

Proof. This is a fairly straightforward modification of the proof for the basic setup.

Let $X \in S_n$ be uniform. Then $\mathbb{E}\phi(X_t) = (1 - \gamma)^t e^{it\theta} \phi(X_0)$ and consequently $\mathbb{E}\phi(X) = 0$. Now let $\psi_t = e^{-it\theta} \phi$. Then $\mathbb{E}\psi_t(X_t) = (1 - \gamma)^t \phi(X_0)$. Moreover,

$$\mathbb{E}[|\psi_{t+1}(X_{t+1}) - \psi_t(X_t)|^2 | X_t] = |e^{-it\theta}| \mathbb{E}[|e^{-i\theta} \phi(X_{t+1}) - \phi(X_t)|^2 | X_t] \leq R.$$

by assumption.

Now mimic the proof for the basic setup, replacing $\phi(X_t)^2$ with $|\psi_t(X_t)|^2$ and keeping in mind that for a complex-valued random variable, Y , $\text{Var}Y = \mathbb{E}[|Y - \mathbb{E}Y|^2] = \mathbb{E}[|Y|^2] - |\mathbb{E}Y|^2$. \square

When it comes to upper bounds for GR-shuffles, we will rely on coupling. It will turn out, for upper as well as lower bounds, that we will achieve useful results only in some simple special cases. The only case where we will be able to determine the precise order of τ_{mix} is the *bottom-to-top shuffle*, presented below.

Lower bounds. As in many of the examples above, we will use the movement of a single card to find eigenvalue/eigenvector pairs to use. Even so, it turns out that even in the simplest special cases, it is virtually impossible to determine these exactly. Then the following approximation lemma is very useful.

Lemma 3.1 (*Approximation Lemma*) *Let \mathbb{D} be the closed unit disc in the complex plane. Assume that $f : \mathbb{D} \rightarrow \mathbb{C}$ is analytic (e.g. a polynomial), $f(0) = 1$ and $|f'(z)| \geq 1$ for all z . Then there exists $z_0 \in \mathbb{D}$ such that $f(z_0) = 0$.*

The real-valued analog of the approximation lemma is obvious. A proof of the lemma can be found in [10]. It requires, however, a background in complex analysis that is not contained in this course.

The way the approximation lemma will be used, is to draw the conclusion that if $f(z_1) = a$ and $|f'(z)| \geq b$ in a sufficiently large neighborhood of z_1 , then f has a root within distance a/b of z_1 .

Writing $P^{(1)}$ for the transition matrix of the MC described by a single card, writing the eigenvalue/eigenvector relation $P^{(1)}\xi = \lambda\xi$ coordinatewise and letting $m_k = \sum_{j < k} p_j$ and $M_k = \sum_{j > k} p_j$ yields

$$\lambda\xi(k) = p_k\xi(1) + m_k\xi(k) + M_k\xi(k + 1),$$

$k = 1, 2, \dots, n$. For general p_i 's, this system of equations is intractable. We will consider two special cases: the bottom-to-top shuffle and the *overlapping cycles shuffle*.

Example. The bottom-to-top shuffle. This shuffle is given by $p_{n-k+1} = p_{n-k+2} = \dots = p_n = 1/k$ for some fixed $k = k(n)$. The eigenvalue/eigenvector equations then simplify to $\xi(1) = 1$ (assumed without loss of generality),

$$\lambda\xi(j) = \xi(j+1),$$

$j = 1, \dots, n-k$ and

$$\lambda\xi(n-k+j) = \frac{1}{k} + \frac{j-1}{k}\xi(n-k+j) + \frac{k-j}{k}\xi(n-k+j+1),$$

$j = 1, \dots, k$.

The first $n-k$ equations imply that $\xi(n-k+1) = \lambda^{n-k}$. Solving for $\xi(n)$ in the last equation yields

$$\xi(n) = \frac{1}{k\lambda - k + 1}$$

(or $\lambda = (k-1)/k$ which is not a useful eigenvalue for our purposes). Inserting into equation $n-1$ and solving for $\xi(n-1)$, then leads to

$$\xi(n-1) = \frac{1 + \xi(n)}{\lambda - k + 2} = \frac{1}{k\lambda - k + 1} = \xi(n).$$

Insert this into equation $n-2$ and solve for $\xi(n-2)$ to get

$$\xi(n-2) = \frac{1 + 2\xi(n)}{\lambda - k + 3} = \frac{1}{k\lambda - k + 1} = \xi(n).$$

Carrying on like throughout the last k equations yields

$$\xi(n) = \phi(n-1) = \dots = \xi(n-k+1) = \frac{1}{k\lambda - k + 1}.$$

Equating the two expressions for $\xi(n-k+1)$ gives us the following characteristic equation for λ :

$$g(\lambda) := \lambda^{n-k+1} - \frac{k-1}{k}\lambda^{n-k} - \frac{1}{k} = 0.$$

Assume now for a while that $k = o(n)$. Let $w = 2\pi/n$. We will now try to “guess” a root of $g(\lambda)$, estimate the error and the derivative of g and use the

approximation lemma. We will be slightly sketchy about this, referring to [10] for more details.

It is easy to see, using Taylor's formula for the sine and cosine functions, that taking $\lambda = e^{iw}$ gives an $\Im g(\lambda) = O(k^3 n^{-3})$ and $\Re g(\lambda) = O(k^2 n^{-2})$. The worst of these errors is the one for the real part, so in order to adjust for this, our next guess is

$$\lambda_0 = (1 - \gamma_0)e^{iw}$$

for a small order real γ_0 . Inserting this we get

$$\begin{aligned} \Re g(\lambda_0) &\approx (1 - n\gamma_0)\left(1 - \frac{(k-1)^2 w^2}{2}\right) - \frac{k-1}{k}(1 - n\gamma_0)\left(1 - \frac{k^2 w^2}{2}\right) - \frac{1}{k} \\ &\approx -\frac{n}{k}\gamma_0 + \frac{k-1}{2}w^2 \end{aligned}$$

which vanishes if $\gamma_0 = k(k-1)w^2/(2n)$. Since γ_0 is of order $k^2 n^{-3}$ it readily seen, using one more term in Taylor's formula for sine and cosine, that $\Im g(\lambda_0) = O(k^2 n^{-3})$ and $\Re g(\lambda_0) = O(k^3 n^{-4})$. Hence

$$g(\lambda_0) = O(k^2 n^{-3}).$$

Since

$$g'(\lambda) = (n - k + 1)\lambda^{n-k} + \frac{k-1}{k}(n-k)\lambda^{n-k-1} = (1 + o(1))\frac{n}{k}$$

within distance, say, $\gamma_0/2$ of λ_0 , the approximation lemma entails that g has a root within distance $O(k^3 n^{-4})$ of λ_0 , in particular, there is an eigenvalue $\lambda = (1 - \gamma)e^{i\theta}$, where $\gamma = (1 + o(1))k(k-1)w^2/(2n)$ and $\theta = (1 + o(1))w$.

Now we apply the second extension. As usual, Z_t^i denotes the position of card i at time t . Let $\phi^i(X_t) = \xi(Z_t^i)$ and let

$$\phi = \sum_{i=1}^m \phi^i,$$

where $m = \lfloor n/2 \rfloor$. Then (λ, ϕ) is an eigenvalue/eigenvector pair. Since $|\xi(i)| = 1 + O(k^2 n^{-2})$ for all i , starting with the cards in order gives $|\phi(X_0)| = Cn$. Next we must control the quadratic change of ϕ for one step of the shuffle. By the triangle inequality

$$|e^{-i\theta}\phi(X_{t+1}) - \phi(X_t)| \leq \sum_i |e^{-i\theta}\xi(Z_{t+1}^i) - \xi(Z_t^i)|.$$

For the at most $m \leq n - k$ cards in a position r among the top $n - k$ positions,

$$|e^{-i\theta}\xi(Z_{t+1}^i) - \xi(Z_t^i)| \leq |\lambda|^r \gamma \leq \gamma = O(k^2 n^{-3}).$$

For the card taken to the top, $e^{-i\theta}\xi(Z_{t+1}^i)$ and $\xi(Z_t^i)$ differ by $|1 - e^{-i\theta}\lambda^{n-k}| = O(kn^{-1})$. For the k cards in the bottom of the deck, the worst case is when $k - 1$ of them stay put and for such a card

$$|e^{-i\theta}\xi(Z_{t+1}^i) - \xi(Z_t^i)| = |\lambda^{n-k}(1 - e^{-i\theta})| = O(n^{-1}).$$

Summing up

$$\begin{aligned} |e^{-i\theta}\phi(X_{t+1}) - \phi(X_t)| &\leq mO(k^2 n^{-3}) + O(kn^{-1}) + (k - 1)O(n^{-1}) \\ &= O(kn^{-1}). \end{aligned}$$

Hence we can take $R = O(k^2 n^{-2})$. Taking $a = 1/2$ in the second extension and recalling that $w = 2\pi/n$ now gives

$$\begin{aligned} \tau_{\text{mix}} &\geq (1 + o(1))\gamma^{-1}(\log(Cn) - \frac{1}{2}\log(8R/\gamma)) \\ &= (1 + o(1))\frac{1}{4\pi^2 k(k-1)}n^3 \log n. \end{aligned}$$

Note that a special case is the Rudvalis shuffle for which we get the lower bound $(1/(8\pi^2))n^3 \log n$, the same lower bound as Wilson found in [19] via a state space extension. The result deserves to be stated as a theorem.

Theorem 3.2 *For the bottom-to-top shuffle taking a card from the bottom $k = o(n)$ positions,*

$$\tau_{\text{mix}} \geq \frac{1}{4\pi^2 k(k-1)}n^3 \log n.$$

□

Example. The overlapping cycles shuffle. Here $k = k(n)$ is again a predetermined fixed number and $p_{n-k} = p_n = 1/2$. To avoid periodicity, we must stipulate that k be odd. The name ‘‘overlapping cycles shuffle’’ was proposed by Angel, Peres and Wilson in [4], where a detailed analysis of the eigenvalues for the single card MC is carried out; it turns out that these describe two overlapping

cycles in the unit disc of the complex plane. The equations for (λ, ξ) for the single card chain now become

$$\lambda\xi(j) = \xi(j+1),$$

$$j = 1, 2, \dots, n-k-1,$$

$$\lambda\xi(n-k) = \frac{1}{2}\xi(1) + \frac{1}{2}\xi(n-k+1),$$

$$\lambda\xi(j) = \frac{1}{2}\xi(j) + \frac{1}{2}\xi(j+1),$$

$$j = n-k+1, \dots, n-1 \text{ and}$$

$$\lambda\xi(n) = \frac{1}{2}\xi(1) + \frac{1}{2}\xi(n).$$

Here we will be very sketchy, referring to [10] and [4] for all details.

Solving forward and backward to get two expressions for $\xi(n-k)$ gives the characteristic equation

$$(2\lambda - 1)^k(2\lambda^{n-k} - 1) - 1 = 0.$$

Some guessing work leads to the eigenvalue candidate $\lambda_0 = (1 - \gamma_0)e^{i2\pi/n}$, where $\gamma_0 \leq 2\pi^2k(k+1)/n^3$. Estimating the error and the derivative of the characteristic polynomial and using the approximation lemma, tells us that there is a true eigenvalue $\lambda = (1 - \gamma)e^{i\theta}$ where $\gamma = (1 + o(1))\gamma_0$ and $\theta = (1 + o(1))2\pi/n$.

An analogous use of the second extension to that for the bottom-to-top shuffle, then yields

$$\tau_{\text{mix}} \geq \frac{1}{4\pi^2k(k+1)}n^3 \log n.$$

We conjecture that for $k = O(1)$, this is the true order of mixing.

For other k however, this is probably not the case. The case $k = n/2$ was treated in [10], where it was shown that τ_{mix} is at least of order n^2 . This was done via a ‘‘classical’’ argument, i.e. by finding an event which is highly probable for the actual chain but highly unlikely for the uniform deck.

Angel, Peres and Wilson [4] found (which I failed to do in [10]) that when k is of the same order as n , there are actually eigenvalues that confirm this. E.g. when $k = n/3$, there is an eigenvalue very close to $(1 - 3\pi^2/n^2)e^{i3\pi/n}$. Hence

$$\tau_{\text{mix}} \geq \tau_2 \geq \tau_2^{(1)} = (1 + o(1))\frac{2}{3\pi^2}n^2,$$

where $\tau_2^{(1)}$ is the relaxation time for the single card chain.

What Angel, Peres and Wilson found, was the following surprising result.

Theorem 3.3 *Let $k = \lfloor \alpha n \rfloor$, $\alpha \in (0, 1)$. For almost every α , $\tau_2^{(1)}$ is of order $n^{3/2}$. However, when $\alpha = p/q$, where p and q are two relatively prime integers, then*

$$\tau_2^{(1)} = \begin{cases} (1 + o(1)) \frac{2}{\pi^2 pq} n^2, & p \text{ and } q \text{ both odd} \\ (1 + o(1)) \frac{8}{\pi^2 pq} n^2, & \text{otherwise} \end{cases}$$

Unfortunately, due to the strong dependence between cards, Wilson's technique cannot increase this result by the usual $\log n$ -factor for the whole deck. It is still a wide open question, what the mixing time of the whole deck is.

However, the ideas of [4] can be used to make a slight improvement (but still probably not tight) of the lower bounds when k is of larger order than $n^{2/3}$ and lower order than n . Assume e.g. that $2k|n - k$. (This corresponds to p and q odd in the above result.) Then, using the approximation lemma, one can check that there is an eigenvalue of the single card chain,

$$\lambda = (1 - \gamma)e^{i\theta}$$

where $\gamma = (1 + o(1))\pi^2/(2nk)$ and $\theta = (1 + o(1))\pi/k$. The second extension then gives

$$\tau_{\text{mix}} \geq (1 + o(1)) \frac{1}{\pi^2} nk \log \frac{n}{k}.$$

Assuming instead that $n - k$ is divisible by k but not $2k$ (corresponding to p or q even), gives $\gamma = (1 + o(1))2\pi^2/(nk)$ and $\theta = (1 + o(1))2\pi/k$ and in the end the lower bound

$$\tau_{\text{mix}} \geq (1 + o(1)) \frac{1}{4\pi^2} nk \log \frac{n}{k}.$$

All of this works fine for all $k = o(n)$, but gives no improvement over earlier results for $k = O(n^{2/3})$.

In analogy with the above theorem, we conjecture that for most $k = o(n)$ (i.e. where $n - k$ is not divisible by k), the relaxation time for the single card chain is of order $1/(n\sqrt{k})$. \square

Upper bounds. This will only be done for the bottom-to-top shuffle. We assume $k = o(n)$, leaving $k = \Theta(n)$ to the reader.

Let $\{Y_t\}$ denote a stationary copy of the MC under study. Let A_t and be the set of cards, that at time t are in one of the bottom k positions in the $\{X_t\}$ -process and let B_t be the corresponding set for $\{Y_t\}$. Couple $\{Y_t\}$ to $\{X_t\}$ by, for the $t + 1$ 'th step for each t , when a card in $A_t \cap B_t$ is moved to the top on one deck, letting the same card be moved to the top in the other deck too. When the card

moved to the top for $\{X_t\}$ is in $A_t \setminus B_t$, then pick a uniformly chosen card from $B_t \setminus A_t$ for $\{Y_t\}$. Note about this coupling that

- A card, c , can never pass its copy in the other deck unless $c \in A_t \cap B_t$.
- If $c \in A_t \cap B_t$, it will be coupled as soon as it is moved to the top.

Consequently, the decks will be coupled as soon as every card c has been in $A_t \cap B_t$ and, after the first time this happens, all cards that are at that time in one of the bottom k positions, have been moved to the top. Let T as usual denote the coupling time, let

$$T_0(c) = \min\{t : c \in A_t \cap B_t\}$$

and let $T_0 = \max_c T_0(c)$. Since $T - T_0$ is bounded by the time taken to pick to the top, every card in A_t , $\mathbb{E}[(T - T_0)^+] \leq k \log k$ be the coupon collector's problem. Hence

$$\mathbb{P}(T - T_0 \geq f(n)k \log k) = o(1)$$

for any $f(n) \rightarrow \infty$. It remains to analyze T_0 .

For this, it will be more convenient to denote the positions $0, 1, \dots, n - 1$ rather than $1, 2, \dots, n$ and consider a transformation of the decks:

$$X'_t = (n - 1 \ n - 2 \ \dots \ 1 \ 0)^t \circ X_t$$

$$Y'_t = (n - 1 \ n - 2 \ \dots \ 1 \ 0)^t \circ Y_t.$$

In words, the shuffle acting on the transformed decks behaves in the following way. An “active layer” of k cards starts from the bottom of the deck. At each step of the shuffle, a uniformly chosen card in the active layer is moved to the bottom of the layer, mod n . Then the active layer is moved one step up the deck mod n .

Fix a card c . Let τ_j be the j 'th time c leaves A_t and let ν_j be the j 'th time it leaves B_t , $j = 1, 2, \dots$. Let J be the smallest J such that $\tau_j = \nu_j$ or $\tau_j = \nu_{j+1}$. Then $\tau_J = T_0(c)$. We now want to estimate J .

Let $\tau_j - \tau_{j-1}$ be called the j 'th cycle for c in $\{X_t\}$. During a cycle, c moves $k - G$ steps down the transformed deck mod n , where G is geometric with parameter $1/k$. Also, up to cycle J , c moves independently in the two decks. Hence, letting

$$U_j = X'_{\tau_j}(c) - Y'_{\nu_j}(c),$$

$\{U_j\}$ is a simple random walk on \mathbb{Z}_n with step size mean 0 and step size variance $2k(k - 1)$, starting from somewhere in $(0, n - 1)$ and not passing 0 until cycle J .

Thus, for estimating J , we can equally well identify $\{U_j\}$ with a SRW, $\{V_j\}$, on \mathbb{Z} starting from somewhere in $(0, n)$. Then J coincides with the first time this SRW leaves $(0, n)$ and is thus dominated by the first time it leaves $(-n/2, n/2)$ under the assumption $V_0 = 0$. Let

$$W_j = \frac{1}{\sqrt{2k(k-1)}} V_j.$$

Then W_j has step size variance 1 and leaves $(-n/(2(2k(k-1))^{1/2}), n/(2(2k(k-1))^{1/2}))$ when V_j leaves $(-n/2, n/2)$. By Donsker's Theorem, $M^{-1}W_{M^s}$ converges in distribution to a standard Brownian motion B_s as $M \rightarrow \infty$. Taking

$$M = \frac{n}{2\sqrt{2k(k-1)}}$$

we get

$$\begin{aligned} \mathbb{P}\left(\forall s \leq s_0 : V_{M^2 s} \in \left(-\frac{n}{2}, \frac{n}{2}\right)\right) &= (1 + o(1)) \mathbb{P}\left(\forall s \leq s_0 : |B_s| < 1\right) \\ &\leq (1 + o(1)) \frac{4}{\pi} e^{-\pi^2 s_0 / 8} \end{aligned}$$

where the last inequality can be found in [9, Section 7.8] (and was used for the coupling of the transposing neighbors shuffle in an earlier section). Letting $s_0 = (8/\pi^2) \log n + \log \log n$, the right hand side is $o(1)$. Hence, with probability $1 - o(1)$, $T_0(c)$ does not exceed $(1 + o(1))(8/\pi^2)M^2 \log n = (1 + o(1))(n^2/\pi^2 k(k-1)) \log n$ cycles.

Since cycle times are independent and distributed like $n - k + G$, it follows from Chernoff bounds for binomial random variables, that the time taken for these $C(n^2/k^2) \log n$ cycles exceeds $(1 + a)C(n^3/k^2) \log n$ is $o(1)$ for any $a > 0$.

Maximizing over c and adding $T - T_0$ noting that $f(n)k \log k$ is of smaller order than $(n^3/k^2) \log n$ for suitable $f(n)$ and taking n large enough gives

$$\mathbb{P}\left(T > (1 + o(1)) \frac{n^3}{\pi^2 k(k-1)} \log n\right) = o(1).$$

Comining this with the lower bound above gives the following theorem.

Theorem 3.4 *For bottom-to-top shuffling with $k = o(n)$,*

$$(1 + o(1)) \frac{n^3}{4\pi^2 k(k-1)} \log n \leq \tau_{\text{mix}} \leq (1 + o(1)) \frac{n^3}{\pi^2 k(k-1)} \log n.$$

4 Advanced L^2 -techniques

In this section we will go through the L^2 -theory that provides the background for Nash and log-Sobolev technique, and leads up to the comparison technique of Diaconis and Saloff-Coste [6]. The L^2 -theory is by far most useful for reversible Markov chains. Hence it will throughout this section be assumed that MC's under study are reversible, unless otherwise stated.

4.1 The basics

We take this from [3, Chapter 8].

Recall that the relaxation time τ_2 for a continuous time Markov chain $\{X_t\}$ with generator Q and stationary distribution π , is given by $\tau_2 = 1/\lambda_2$. Here λ_2 is the second smallest eigenvalue of $-Q$. (The smallest eigenvalue is of course 0.) In discrete time $\tau_2 = 1/(1 - \lambda_2)$ where λ_2 is the second largest eigenvalue of the transition matrix.

Recall also the definition

$$\hat{\tau} = \inf \{t : \|\mathbb{P}(X_t \in \cdot) - \pi\|_2 \leq \frac{1}{2}\}$$

and the basic fact that $\tau_{\text{mix}} \leq \hat{\tau}$. The p -norm or, more correctly, the norm in $L^p(\pi)$ of a signed measure ν on S , the state space of the MC, was earlier defined by

$$\|\nu\|_p^p = \sum_{i \in S} \left| \frac{\nu(i)}{\pi(i)} \right|^p \pi(i).$$

For a function, g on S , we define the p -norm differently:

$$\|g\|_p^p = \sum_{i \in S} |g(i)|^p \pi(i) = \mathbb{E}_\pi[|g(X_0)|^p].$$

(Equivalently, one could define the p -norm of the signed measure ν as the p -norm of the function g given by $g(i) = \nu(i)/\pi(i)$.) The *Dirichlet form*, $\mathcal{E}(g, g)$, of a function g on S is given by

$$\mathcal{E}(g, g) = - \sum_i \sum_j \pi(i) q_{i,j} g(i) g(j)$$

in continuous time and

$$\mathcal{E}(g, g) = \sum_i \sum_j \pi(i) p_{ij} g(i) (g(i) - g(j))$$

in discrete time. In probabilistic terms, this can be written as

$$\mathcal{E}(g, g) = \lim_{t \downarrow 0} \frac{1}{t} \mathbb{E}_\pi [g(X_0)(g(X_0) - g(X_t))] = \frac{1}{2} \lim_{t \downarrow 0} \frac{1}{t} \mathbb{E}_\pi [(g(X_0) - g(X_t))^2],$$

where the second inequality uses reversibility. The discrete-time versions are

$$\mathcal{E}(g, g) = \mathbb{E}_\pi [g(X_0)(g(X_0) - g(X_1))] = \frac{1}{2} \mathbb{E}_\pi [(g(X_0) - g(X_1))^2].$$

The following lemma is known as the extremal characterization of relaxation time.

Lemma 4.1 *The relaxation time satisfies*

$$\tau_2 = \sup \left\{ \frac{\|g\|_2^2}{\mathcal{E}(g, g)} : g \neq 0, \mathbb{E}_\pi g(X_0) = 0 \right\}.$$

Proof. We start with continuous time. We first need the *extremal characterization of eigenvalues*. Let A be a symmetric $n \times n$ -matrix. Then A has the real eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and corresponding pairwise orthogonal eigenvectors, v_1, v_2, \dots, v_n . Then

$$A = VDV^T$$

or, equivalently,

$$D = V^T AV$$

where $V = [v_1 \ v_2 \ \dots \ v_n]$ and $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. For a unit vector x , write $y = P^T x$. Then y is also a unit vector and

$$\sum_i \sum_j x_i a_{ij} x_j = x^T A x = y^T D y = \sum_i \lambda_i y_i^2.$$

The right hand side is minimized over y when $y = e_1$, and then takes on the value λ_1 . Since $y = e_1$ is equivalent to $x = v_1$, we have that

$$\lambda_1 = \min \left\{ \frac{\sum_i \sum_j x(i) a_{ij} x(j)}{\sum_i x(i)^2} : x \neq 0 \right\}$$

and that a vector x for which the minimum is attained is an eigenvector to λ_1 .

Now repeat the procedure, but now only minimizing the quadratic form over x 's that are orthogonal to v_1 . Then the minimum is λ_2 , which is attained for $x = v_2$, i.e.

$$\lambda_2 = \min \left\{ \frac{\sum_i \sum_j x(i) a_{ij} x(j)}{\sum_i x(i)^2} : x \neq 0, \sum_i x(i) v_1(i) = 0 \right\}$$

and a vector x for which the minimum is attained is an eigenvector to λ_2 .

Next, we apply this to the symmetric matrix $-A$, where

$$A = \left[\sqrt{\frac{\pi(i)}{\pi(j)}} q_{ij} \right]_{i,j \in S}.$$

Then $-A$ and $-Q$ have the same eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$, and if $x = (x(1), x(2), \dots, x(n))^T$ is an eigenvector of Q , then

$$(\sqrt{\pi(1)}x(1), \sqrt{\pi(2)}x(2), \dots, \sqrt{\pi(n)}x(n))^T$$

is the corresponding eigenvector of A . Since the first eigenvector of $-Q$ is the constant vector, a first eigenvector of A is $[\sqrt{\pi(i)}]_{i \in S}$. Hence

$$\lambda_2 = \frac{1}{\tau_2} = \inf \left\{ \frac{\sum_i \sum_j x(i) (-\sqrt{\pi(i)/\pi(j)} q_{ij}) x(j)}{\sum_i x(i)^2} : \sum_i x(i) \sqrt{\pi(i)} = 0 \right\}.$$

Substituting $x(i) = \sqrt{\pi(i)}g(i)$, the right hand side becomes

$$\inf \left\{ \frac{-\sum_i \sum_j \pi(i) q_{ij} g(i) g(j)}{\sum_i \pi_i g(i)^2} : \sum_i \pi(i) g(i) = 0 \right\}$$

which equals

$$\inf \left\{ \frac{\mathcal{E}(g, g)}{\|g\|_2^2} : \mathbb{E}_\pi g(X_0) = 0 \right\}.$$

Inverting this now finishes the proof.

The discrete time proof is analogous, with Q replaced by $P - 1$. □

Note that an alternative way to express the extremal characterization is:

$$\tau_2 = \sup \left\{ \frac{\mathbb{V}\text{ar}_\pi f(X_0)}{\mathcal{E}(f, f)} : f \neq 0 \right\}.$$

From now on, we assume until further notice that time is continuous.

Lemma 4.2 For each t , let

$$f_t(j) = \frac{\mathbb{P}(X_t = j)}{\pi(j)}.$$

Then

$$\frac{d}{dt} \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 = -2\mathcal{E}(f_t, f_t).$$

Proof. By the forward equations (which are essentially the definition of the q_{ij} 's, using the convention that $-q_{jj}$ is the intensity for jumping away from j),

$$\frac{d}{dt} \mathbb{P}(X_t = j) = \sum_i \mathbb{P}(X_t = i) q_{ij}.$$

Therefore, by the chain rule,

$$\begin{aligned} \frac{d}{dt} \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 &= \frac{d}{dt} \left(\sum_j \frac{\mathbb{P}(X_t = j)^2}{\pi(j)} - 1 \right) \\ &= \sum_j \sum_i \frac{2}{\pi(j)} \mathbb{P}(X_t = j) \mathbb{P}(X_t = i) q_{ij} \\ &= 2 \sum_j \sum_i \pi_i q_{ij} f_t(i) f_t(j) = -2\mathcal{E}(f_t, f_t). \end{aligned}$$

□

The next result is the L^2 -contraction Lemma.

Lemma 4.3 For all t ,

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2 \leq e^{-t/\tau_2} \|\mathbb{P}(X_0 \in \cdot) - \pi\|_2.$$

Proof. Let f_t be as in Lemma 4.2 and let $h(t) = \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2$. By Lemma 4.2 and the third expression for the Dirichlet form,

$$\frac{d}{dt} h(t) = -2\mathcal{E}(f_t, f_t) = -2\mathcal{E}(f_t - 1, f_t - 1).$$

Since $\mathbb{E}_\pi f_t(X_0) = 1$, the extremal characterization of τ_2 entails that the right hand side is bounded above by $-2\|f_t - 1\|_2^2/\tau_2$, which in turn equals

$$-\frac{2}{\tau_2} \sum_i \pi_i \left(\frac{\mathbb{P}(X_t = i)}{\pi_i} - 1 \right)^2 = -\frac{2}{\tau_2} \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 = -\frac{2}{\tau_2} h(t).$$

Integrating, using integrating factor, now gives

$$h(t) \leq e^{-2/\tau_2} h(0).$$

Taking square roots now gives the result. \square

Operator norms. Let A be an $n \times n$ -matrix. Pick $q_1, q_2 \in [1, \infty]$. One may think of A as a linear operator from $L^{q_1}(\pi)$ to $L^{q_2}(\pi)$ defined by $A(f) = Af$. One can then define the following operator norm for A :

$$\|A\|_{q_1 \rightarrow q_2} = \sup\{\|Af\|_{q_2} : \|f\|_{q_1} = 1\}.$$

When $A = P$, a transition matrix for a reversible MC with stationary distribution π , one can use the symmetry of the matrix $[\sqrt{\pi(i)/\pi(j)}p_{ij}]_{i,j \in S}$ to show that

$$\|P\|_{q_1 \rightarrow q_2} = \sup\{\|\nu P\|_{q_2} : \|\nu\|_{q_1} = 1\}$$

where the supremum is now over signed measures ν on S . This also holds for the difference of such matrices, such as e.g. $P_s - P_t$, where P_t is the transition matrix for t units of time of a continuous time reversible MC. All our applications will be on transition matrices or differences of the above kind, so we do not have to distinguish between the two expressions for the operator norm.

We will be mostly concerned with the special cases $\|P\|_{2 \rightarrow 2}$ and $\|P\|_{2 \rightarrow \infty}$. For $\|\cdot\|_{\infty}$ we use the usual conventions

$$\|f\|_{\infty} = \max_i |f(i)|$$

for functions f and

$$\|\nu\|_{\infty} = \max_i \left| \frac{\nu(i)}{\pi(i)} \right|$$

for signed measures ν .

For continuous time MC's, we will use the notation

$$N(s) = \|P_s\|_{2 \rightarrow \infty}.$$

Then it is readily seen that $N(0) = \|I\|_{2 \rightarrow \infty} = \pi_*^{-1/2}$, where $\pi_* = \min_i \pi(i)$. It is also evident from Lemma 4.6 below that $N(s) \rightarrow 1$ as $s \rightarrow \infty$.

By definition,

$$\|BA\|_{q_1 \rightarrow q_3} \leq \|A\|_{q_1 \rightarrow q_2} \|B\|_{q_2 \rightarrow q_3}$$

for any A, B, q_1, q_2 and q_3 .

Lemma 4.4 For any time t and any function f on S , the following holds.

(a)

$$\frac{d}{dt} \|P_t f\|_2^2 = -2\mathcal{E}(P_t f, P_t f) \leq -\frac{2}{\tau_2} \text{Var}_\pi P_t f(X_0),$$

(b)

$$\|P_t - P_\infty\|_{2 \rightarrow 2} \leq e^{-t/\tau_2},$$

where $P_\infty = \lim_{t \rightarrow \infty} P_t$ is the transition matrix for which every row equals π .

Proof. The inequality in (a) follows immediately from the extremal characterization of τ_2 . For the equality, use the backwards equations to get

$$\frac{d}{dt} p_{ik}(t) = \sum_j q_{ij} p_{jk}(t).$$

This entails that

$$\begin{aligned} \frac{d}{dt} P_t f(i) &= \frac{d}{dt} \sum_k p_{ik}(t) f(k) = \sum_k \sum_j q_{ij} p_{jk}(t) f(k) \\ &= \sum_j q_{ij} P_t f(j). \end{aligned}$$

Hence

$$\begin{aligned} \frac{d}{dt} \|P_t f\|_2^2 &= \frac{d}{dt} \sum_i \pi_i (P_t f(i))^2 \\ &= 2 \sum_i \sum_j \pi_i q_{ij} P_t f(i) P_t f(j) \\ &= -2\mathcal{E}(P_t f, P_t f). \end{aligned}$$

For part (b), note that for any f , $P_\infty f(i) = \mathbb{E}_\pi P_t f(X_0)$ for all i and t . Hence

$$\|P_t(f - P_\infty f)\|_2^2 = \text{Var}_\pi P_t f(X_0).$$

Therefore, by (a),

$$\begin{aligned} \frac{d}{dt} \|(P_t - P_\infty)f\|_2^2 &= \frac{d}{dt} \|P_t(f - P_\infty f)\|_2^2 \\ &\leq -\frac{2}{\tau_2} \|P_t(f - P_\infty f)\|_2^2 \\ &= -\frac{2}{\tau_2} \|(P_t - P_\infty)f\|_2^2. \end{aligned}$$

Integrating then gives

$$\begin{aligned}\|(P_t - P_\infty)f\|_2^2 &\leq e^{-2t/\tau_2} \|(P_0 - P_\infty)f\|_2^2 = e^{-2t/\tau_2} \text{Var}_\pi f(X_0) \\ &\leq e^{-2t/\tau_2} \|f\|_2^2.\end{aligned}$$

Taking square roots and using the definition of $\|\cdot\|_{2 \rightarrow 2}$ now completes the proof. \square

By definition of $\|\cdot\|_{2 \rightarrow 2}$ and Lemma 4.4(b),

$$\begin{aligned}\|\mathbb{P}(X_{s+t} \in \cdot) - \pi\|_2 &\leq \|\mathbb{P}(X_s \in \cdot)\|_2 \|P_t - P_\infty\|_{2 \rightarrow 2} \\ &\leq \|\mathbb{P}(X_s \in \cdot)\|_2 e^{-t/\tau_2}.\end{aligned}$$

Taking e.g. $s = 0$, one gets

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2 \leq \|\mathbb{P}(X_0 \in \cdot)\|_2 e^{-t/\tau_2} \leq \frac{1}{\pi_*} e^{-t/\tau_2}.$$

However, we will not apply this result here. Our main result of this section will be the following.

Lemma 4.5 *For any s and t ,*

$$\begin{aligned}\|\mathbb{P}(X_{s+t} \in \cdot) - \pi\|_2 &\leq \|P_{s+t} - P_\infty\|_{2 \rightarrow \infty} \\ &\leq \|P_s\|_{2 \rightarrow \infty} \|P_t - P_\infty\|_{2 \rightarrow 2} \\ &\leq \|P_s\|_{2 \rightarrow \infty} e^{-t/\tau_2}.\end{aligned}$$

Proof. The last inequality is Lemma 4.4(b). The second inequality follows from

$$\|P_{s+t} - P_\infty\|_{2 \rightarrow \infty} = \|P_s(P_t - P_\infty)\|_{2 \rightarrow \infty} \leq \|P_t - P_\infty\|_{2 \rightarrow 2} \|P_s\|_{2 \rightarrow \infty}.$$

For the first inequality, define the function f as

$$f(j) = \frac{\mathbb{P}(X_u = j)}{\pi(j)} - 1.$$

Then

$$\|f\|_2^2 = \sum_j \pi(j) \left(\frac{\mathbb{P}(X_u = j)}{\pi(j)} - 1 \right)^2 = \|\mathbb{P}(X_u \in \cdot) - \pi\|_2^2.$$

On the other hand

$$\begin{aligned}
\|f\|_2 &= \max\{|\mathbb{E}_\pi[fg]| : \|g\|_2 = 1\} \\
&= \max\{|\sum_j f(j)g(j)\pi(j)| : \|g\|_2 = 1\} \\
&= \max\{|\sum_j (\mathbb{P}(X_u = j) - \pi(j))g(j)| : \|g\|_2 = 1\} \\
&\leq \max\{\max_i |(P_u - P_\infty)g(i)| : \|g\|_2 = 1\} \\
&= \max\{\|(P_u - P_\infty)g\|_\infty : \|g\|_2 = 1\} \\
&= \|P_u - P_\infty\|_{2 \rightarrow \infty},
\end{aligned}$$

where the inequality follows on conditioning on X_0 . Now take $u = s + t$ to finish the proof. \square

Taking $s = 0$ in Lemma 4.5 again leads to the inequality preceding the lemma. However, for MC's on large state spaces, such as card shuffling chains, taking $s = 0$ will not produce sharp results. The essence of Nash- and log-Sobolev technique is to make a better choice of s in Lemma 4.5. For this, we need some more information about $N(s) = \|P_s\|_{2 \rightarrow \infty}$. Let $g_i(j) = p_{ij}(s)/\pi(j)$. Then

$$\begin{aligned}
\max_i \sqrt{\sum_j \frac{p_{ij}(s)^2}{\pi(j)}} &= \max_i \|g_i\|_2^2 \\
&= \max_i \left(\max\{|\sum_j p_{ij}(s)f(j)| : \|f\|_2 = 1\} \right) \\
&= \max\{\max_i |P_s f(i)| : \|f\|_2 = 1\} \\
&= \max\{\|P_s f\|_\infty : \|f\|_2 = 1\} = N(s).
\end{aligned}$$

On the other hand, by reversibility,

$$\sum_j \frac{p_{ij}(s)^2}{\pi(j)} = \sum_j \frac{p_{ij}(s)p_{ji}(s)}{\pi(i)} = \frac{p_{ii}(2s)}{\pi(i)}.$$

We have shown

Lemma 4.6 *The function $N(s)$ satisfies*

$$N(s) = \max_i \sqrt{\frac{p_{ii}(2s)}{\pi(i)}}.$$

Example. Let $\{X_t\}$ be a continuous time random walk on \mathbb{Z}_2 , i.e. from any of the two states, 0 and 1, one jumps to the other state with intensity 1. Equivalently, at the times of a Poisson process with intensity 2, a fair coin is flipped to decide whether to move or not. Thus

$$p_{00}(s) = p_{11}(s) = e^{-2s} + \frac{1}{2}(1 - e^{-2s}) = \frac{1}{2}(1 + e^{-2s})$$

and hence

$$N(s) = \sqrt{1 + e^{-4s}}.$$

□

Let us now use the result of this example to bound $\hat{\tau}$ for a random walk on the hypercube \mathbb{Z}_2^d for large d .

Note first the following general fact. Let $X_t = (X_t^1, X_t^2)$, where $\{X_t^1\}$ and $\{X_t^2\}$ are two independent MC's. Then by Lemma 4.6, in the obvious notation,

$$N(s) = N^1(s)N^2(s).$$

Example. Random walk on \mathbb{Z}_2^d . This is the product MC

$$\{X_t\} = \{(X_t^1, X_t^2, \dots, X_t^d)\},$$

where the coordinate processes are random walks on \mathbb{Z}_2 , with jump intensity $1/d$ (so that the total jump intensity is 1). This means that

$$N^k(s) = \sqrt{1 + e^{-4s/d}}$$

and hence

$$N(s) = (1 + e^{-4s/d})^{d/2}.$$

Since the intensity for “stationarization” for the coordinate processes is $2/d$, the natural coupling will, via a coupon-collector analysis, lead to

$$\tau_{\text{mix}} \leq \frac{1}{2}d \log d.$$

We will now see that Lemma 4.5 with $s = \frac{1}{4}d \log d$ leads to an improvement by a factor 2 (apart from the fact that convergence in L^2 is stronger than in total variation). This in fact gives the right mixing time. Intuitively this is seen the following

way. At stationarity, the coordinates of the walker are iid fifty/fifty Bernoulli random variables. By the CLT, this will typically mean that the number of 1's will deviate from $n/2$ by order $n^{1/2}$. To achieve this, starting from $(0, 0, \dots, 0)$, it suffices to stationarize all but $cn^{1/2}$ coordinates (but not less). By the coupon-collector analysis, this takes time $(1 + o(1))\frac{1}{4}d \log d$.

Now, let's get back to work. Recall first that if Q is the generator of a continuous time MC and (λ, f) is an eigenvalue/eigenvector pair of Q , this is equivalent to

$$\lim_{t \downarrow 0} \frac{1}{t} \mathbb{E}[f(X_t) - f(X_0) | X_0 = i] = \lambda f(i)$$

for all $i \in S$. From this, it is readily checked that the eigenvalues of a product MC are the sums of eigenvalues of the coordinate chains and the eigenvectors are the corresponding products of coordinate eigenvectors.

For a random walk on \mathbb{Z}_2 with jump intensity $1/d$, the eigenvalues of the generator are 0 and $-2/d$. Thus, for the random walk on \mathbb{Z}_2^d , the second smallest eigenvalue of $-Q$ is $0+0+\dots+0+2/d = 2/d$, i.e. $\tau_2 = d/2$. Taking $s = \frac{1}{4}d \log d$, we have $N(s) = (1 + o(1))e^{1/2}$, so by Lemma 4.5,

$$\|\mathbb{P}(X_{s+t} \in \cdot) - \pi\|_2 \leq (1 + o(1))e^{1/2-2t/d}$$

which is bounded by $(1 + o(1))e^{-1} < 1/2$ (for large d) as soon as $t \geq 3d/4$. Hence

$$\hat{\tau} \leq \frac{1}{4}d \log d + \frac{3}{4}d = (1 + o(1))\frac{1}{4}d \log d.$$

□

4.2 The comparison technique

The comparison technique was introduced by Diaconis and Saloff-Coste in [6] for symmetric random walks on a group, G , i.e. for random walks generated by a symmetric probability measure ν on G , (i.e. such that $\nu(x) = \nu(x^{-1})$ for every $x \in G$). This was later generalized by the same authors to general reversible Markov chains, see [7].

In this presentation, we stick to the simpler original setting with random walks on groups. Hence, let G be a finite group and let ν be a symmetric probability measure on G . Let $n = |G|$. For discrete time let $1 = \kappa_1 \geq \kappa_2 \geq \dots \geq \kappa_n \geq -1$ be the eigenvalues of the transition matrix $P = [\nu(x^{-1}y)]_{x,y \in G}$ and $0 = \lambda_1 \leq \lambda_2 \leq \dots$ be the eigenvalues of $-Q$ for the corresponding continuous time random

walk. Note that since $q_{xy} = p_{xy}$ for $x \neq y$, we have $\lambda_i = 1 - \kappa_i$. Let π be the stationary distribution, i.e. the uniform distribution on G .

Lemma 4.7 *Let $\{X_t\}$ be a discrete time random walk on G generated by ν and let $\{Y_t\}$ be its continuous time version. Then*

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 = \sum_{i=2}^n \kappa_i^{2t}.$$

and

$$\|\mathbb{P}(Y_t \in \cdot) - \pi\|_2^2 = \sum_{i=2}^n e^{-2\lambda_i t}.$$

Proof. Assume discrete time. Since the stationary distribution is uniform, P is symmetric and hence has orthonormal left eigenvectors ϕ_1, \dots, ϕ_n , where ϕ_1 is the uniform vector $n^{-1/2}\mathbf{1}$. Since $X_0 = id$, it is easy to see that $\mathbb{P}(X_0 \in \cdot) = n^{-1/2} \sum_{i=1}^n \phi_i$. Hence

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 = \left\| \sum_{i=2}^n \frac{1}{\sqrt{n}} \kappa_i^t \phi_i \right\|_2^2 = \sum_{i=2}^n \kappa_i^{2t}.$$

The continuous time case is analogous. □

From Lemma 4.7 we see that negative eigenvalues are important in discrete time, whereas they play an insignificant role in continuous time. This is due to the possible parity problems that a discrete time MC may have, but which vanish in continuous time. Because of this, it is considerably easier to transfer results from discrete time to continuous time than vice versa. Here are two general results for transferring.

Lemma 4.8 *Let $\{X_t\}$ and $\{Y_t\}$ be the discrete and continuous time version respectively of a random walk on G generated by ν . Then*

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 \leq n\kappa_n^{2t} + \|\mathbb{P}(Y_t \in \cdot) - \pi\|_2^2$$

and

$$\|\mathbb{P}(Y_{2t} \in \cdot) - \pi\|_2^2 \leq ne^{-2t} + \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2.$$

Proof. Since $\kappa_i = 1 - \lambda_i$ and $1 - x \leq e^{-x}$, we have for every i that $\kappa_i^{2t} \leq e^{-2\lambda_i t}$. Use this together with Lemma 4.7 to get

$$\begin{aligned} \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 &\leq n\kappa_n^{2t} + \sum_{i:\kappa_i > 0} \kappa_i^{2t} \leq n\kappa_n^{2t} + \sum_{i:\lambda_i < 1} e^{-2\lambda_i t} \\ &\leq \kappa_n^{2t} + \sum_{i=2}^n e^{-2\lambda_i t} = \kappa_n^{2t} + \|\mathbb{P}(Y_t \in \cdot) - \pi\|_2^2. \end{aligned}$$

The second case is similar; since $e^{-2x} \leq 1 - x$ for $x \leq 1/2$, we have $e^{-4\lambda_i t} \leq \kappa_i^{2t}$ for $\lambda_i \leq 1/2$ and so

$$\begin{aligned} \|\mathbb{P}(Y_{2t} \in \cdot) - \pi\|_2^2 &= \sum_{i=2}^n e^{-4\lambda_i t} \leq ne^{-2n} + \sum_{i:\lambda_i < 1/2} \kappa_i^{2t} \\ &= ne^{-2n} + \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2. \end{aligned}$$

□

The idea of the comparison technique is to bound mixing time in L^2 of a difficult random walk by comparing it with another random walk for which the L^2 -mixing time is known. The theory of this comparison goes via a comparison of the eigenvalues of the two MC's, which in turn goes via a comparison of the Dirichlet forms.

Recall from above the extremal characterization of eigenvalues (in the proof of the extremal characterization of relaxation time). One way of rephrasing that is the following.

$$\lambda_i = \max_W \min_{g \in W} \frac{\mathcal{E}(g, g)}{\|g\|_2^2},$$

where the maximum is taken over linear subspaces of \mathbb{R}^G of dimension $n + 1 - i$. Now let μ be another symmetric measure on G and consider the random walks (i.e. discrete time and continuous time) generated by μ . Denote these $\{X_t^b\}$ and $\{Y_t^b\}$ respectively (where “b” stands for “benchmark”; these are the MC's we want to compare with). Let κ_i^b and λ_i^b denote the corresponding eigenvalues and \mathcal{E}^b the corresponding Dirichlet form. Then an immediate consequence of the extremal characterization of eigenvalues is the following lemma.

Lemma 4.9 *Assume that A is a constant such that $\mathcal{E}^b(g, g) \leq A\mathcal{E}(g, g)$ for every g . Then, for every $i \in [n]$,*

$$\lambda_i \geq \lambda_i^b / A.$$

Equivalently

$$\kappa_i \leq 1 - \frac{1 - \kappa_i^b}{A}.$$

Next we introduce a new quadratic form:

$$\begin{aligned} \mathcal{F}(g, g) &= \sum_{x, y \in G} \pi(x)p(x, y)g(x)(g(x) + g(y)) \\ &= \mathbb{E}_\pi[g(X_0)(g(X_0) + g(X_1))] \\ &= \frac{1}{2}\mathbb{E}_\pi[(g(X_0) - g(X_1))^2], \end{aligned}$$

where the last equality uses reversibility. Then the extremal characterization applied to $I + P$ under the condition $\mathcal{F}^b \leq A\mathcal{F}$ gives $1 + \kappa_i \geq (1 + \kappa_i^b)/A$, i.e.

Lemma 4.10 *If $\mathcal{F}^b \leq A\mathcal{F}$, then*

$$\kappa_i \geq \frac{1 + \kappa_i^b}{A} - 1.$$

Combining the results so far will now give the first key result of this section.

Lemma 4.11 *If $\mathcal{E}^b \leq A\mathcal{E}$, then*

$$\begin{aligned} \|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 &\leq n\kappa_n^{2t} + \|\mathbb{P}(Y_t^b \in \cdot) - \pi\|_2^2 \\ &\leq n\kappa_n^{2t} + ne^{-t/A} + \|\mathbb{P}(X_{\lfloor t/2A \rfloor}^b \in \cdot) - \pi\|_2^2 \end{aligned}$$

and

$$\|\mathbb{P}(Y_t \in \cdot) - \pi\|_2^2 \leq \|\mathbb{P}(Y_{t/A}^b \in \cdot) - \pi\|_2^2.$$

If also $\mathcal{F}^b \leq A\mathcal{F}$, then

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 \leq ne^{-t/A} + \|\mathbb{P}(X_{\lfloor t/2A \rfloor}^b \in \cdot) - \pi\|_2^2.$$

Proof. The second inequality of the first statement is the second part of Lemma 4.8. The first part follows from mimicking the proof Lemma 4.8 with the extra ingredient of using that $\kappa_i = 1 - \lambda_i \leq 1 - \lambda_i^b/A \leq e^{-\lambda_i^b/A}$ for i 's such that $\kappa_i > 0$. The second statement follows from Lemma 4.7 and the first statement of Lemma 4.9. For the third statement, note that Lemma 4.9 and Lemma 4.10 together imply that $1 - |\kappa_i| \geq (1 - |\kappa_i^b|)/A$. Use the inequality $x \leq e^{-(1-x)}$, $x > 0$, to see that $\kappa_i^{2t} \leq e^{-2t(1-|\kappa_i|)} \leq e^{-2t(1-|\kappa_i^b|)/A}$ which is bounded by $e^{-t/A}$

for those κ_i^b which are less than $1/2$ in absolute value. For the other κ_i^b 's, use that $e^{-2(1-x)} \leq x$ for $1/2 \leq x \leq 1$ to get an upper bound of $|\kappa_i^b|^{t/A}$. Now Lemma 4.7 implies

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 \leq ne^{-t/A} + \sum_{i:|\kappa_i^b| \geq 1/2} |\kappa_i^b|^{t/A} \leq \|\mathbb{P}(X_{[t/2A]}^b \in \cdot) - \pi\|_2^2.$$

□

The next step is now to find useful ways to relate \mathcal{E} to \mathcal{E}^b and \mathcal{F} to \mathcal{F}^b . This can be done via *path counting*. Let E be a symmetric set of generators of G contained in the support of ν . For each $y \in G$, pick a representation $y = x_1 x_2 \dots x_k$, $x_j \in E$. Write $|y| = k$. For each $x \in E$, let $N(x, y)$ be the number of times that x appears in the chosen representation of y .

Lemma 4.12 *Let*

$$A = \max_{x \in E} \frac{1}{\nu(x)} \sum_{y \in G} |y| N(x, y) \mu(y).$$

Then $\mathcal{E}^b \leq A\mathcal{E}$.

Proof. Pick $g \in \mathbb{R}^G$ and $z, y \in G$ and let $y = x_1 \dots x_k$ be the representation of y chosen above. Write

$$\begin{aligned} g(z) - g(zy) &= (g(z) - g(zx_1)) + (g(zx_1) - g(zx_2)) + \dots \\ &\quad + (g(zx_1 \dots x_{k-1}) - g(zx_1 \dots x_k)). \end{aligned}$$

Square both sides and use the inequality $(\sum_{j=1}^k a_j)^2 \leq k \sum_{j=1}^k a_j^2$ (by Cauchy-Schwarz) to bound the right hand side. This gives

$$(g(z) - g(zy))^2 \leq |y| \sum_{j=1}^k (g(zx_1 \dots x_{j-1}) - g(zx_1 \dots x_j))^2.$$

Every term in the sum on the right hand side can be written as $(g(z') - g(z'x_j))$ on setting $z' = zx_1 \dots x_{j-1} \in G$. Summing over $z' \in G$, each term appears at most $N(x_j, y)$ times. Hence

$$\sum_{z \in G} (g(z) - g(zy))^2 \leq |y| \sum_{z \in G} \sum_{x \in E} (g(z) - g(zx))^2 N(x, y).$$

Now multiply both sides with $\mu(y)$ and sum over $y \in G$. Then the left hand side becomes $n\mathcal{E}^b(g, g)$. The right hand side becomes

$$\begin{aligned} \sum_{z \in G} \sum_{x \in E} (g(z) - g(zx))^2 \sum_{y \in G} |y| N(x, y) \mu(y) &\leq A \sum_{z \in E} \sum_{x \in E} (g(z) - g(zx))^2 \nu(z) \\ &\leq An\mathcal{E}(g, g). \end{aligned}$$

□

To get a similar result for the relation between \mathcal{F} and \mathcal{F}^b , the constant A needs only the slight modification that we demand all presentations $y = x_1 \dots x_k$ of elements $y \in G$ to have odd length. This redefines the $|y|$'s and the $N(x, y)$'s and thereby A . Then the result of Lemma 4.12 holds with \mathcal{E} replaced by \mathcal{F} . The only difference from the proof of Lemma 4.12 is in how one writes $g(z) + g(zx)$ as a telescoping sum. (This works only when all $|y|$'s are odd.) We leave this as an exercise to the reader.

Now, before we can use any of this, we need some benchmark card shuffle to compare with. Unfortunately, bounds in L^2 for random walks on S_n are difficult to find. Luckily, Diaconis and Shashahani [8] managed to do this for the random transpositions shuffle, $\mu(id) = 1/n$, $\mu(i j) = 2/n^2$ for $i \neq j$. via Fourier analysis on S_n . Their result is very sharp and states that there is a universal constant β such that at time $(1/2)n(\log n + c)$, the square of the L^2 -distance to stationarity is bounded by βe^{-2c} . In particular $\hat{\tau} \leq (1 + o(1))(1/2)n \log n$ for the random transpositions shuffle.

Example. Transposing neighbors. Consider the transposing neighbors shuffle generated by $\nu(i i + 1) = \nu(id) = 1/n$, $i = 1, \dots, n - 1$. This is slightly different from the version studied earlier, which had a holding probability of $1/2$ in order for the disigned coupling there to work out. We compare this with the random transpositions shuffle. We have $E = \{(i i + 1) : i = 1, \dots, n - 1\} \cup \{id\}$. Let $y = (i j)$, $i > j$, be an arbitrary transposition. Then write

$$y = (i i + 1)(i + 1 i + 2) \dots (j - 1 j)(j - 2 j - 1) \dots (i i + 1).$$

From this it follows that $|y| \leq 2n$ and $N(x, y) \leq 2$ for every x and y . Hence $A \leq 2n^2$. Now we only need an estimate for $\kappa_{n!}$, the smallest eigenvalue of the transposing neighbors shuffle. The following lemma is useful.

Lemma 4.13 *For a discrete time random walk on G generated by ν ,*

$$\kappa_{|G|} \geq -1 + 2\nu(id).$$

Proof. The result is trivial when $\nu(id) = 0$, so assume $\nu(id) > 0$. Consider then the random walk generated by the measure ν' where $\nu'(x) = \nu(x)/(1 - \nu(id))$, $x \neq id$ and $\nu'(id) = 0$. This random walk has the eigenvalues $\kappa'_i = (\kappa_i - \nu(id))/(1 - \nu(id))$. Since $\kappa'_i \geq -1$, the result follows. \square

By Lemma 4.13, $\kappa_{n!} \geq -1 + 1/n$. By Lemma 4.12 and Lemma 4.11,

$$\|\mathbb{P}(X_t \in \cdot) - \pi\|_2^2 \leq \left(1 - \frac{1}{n}\right)^{2t} + n!e^{-t/(2n^2)} + \|\mathbb{P}(X_{t/(4n^2)}^b \in \cdot) - \pi\|_2^2.$$

This gives $\hat{\tau} \leq (1 + o(1))2n^3 \log n$, the same order as for total variation. \square

Example. Distorted random transpositions. Consider random transpositions where a few particular transpositions are not allowed. For concreteness, consider the shuffle generated by $\nu(id) = 1/(n-1)$, $\nu(i\ j) = 2/(n(n-1))$, $|i-j| \not\equiv 1 \pmod n$, i.e. random transpositions with neighbor transpositions (including $(1\ n)$) forbidden. Again compare with ordinary random transpositions. Writing $(i\ i+1) = (i\ i+3)(i+1\ i+3)(i\ i+3)$, it follows that $|y| = 3$ for every neighbor transposition y . Of course, $|y| = 1$ for all other y . For all non-neighbor transpositions x , we have $N(x, y) \leq 2$ and $N(x, y) = 0$ for all but three different y 's. This gives $A \leq 9$. Hence the comparison gives $\hat{\tau} \leq (1 + o(1))9n \log n$. \square

Example. The symmetrized Rudvalis shuffle. Recall that the Rudvalis shuffle is generated by the measure giving probability $1/2$ to each of the two generators $d_{n-1} = (n-1\ n-2\ \dots\ 1)$ and $d_n = (n\ n-1\ \dots, 1)$. This set of generators is not symmetric, so let us instead consider the *additive symmetrization* of the Rudvalis shuffle, i.e. the shuffle generated by $\nu(d_n) = \nu(d_{n-1}) = \nu(u_n) = \nu(u_{n-1}) = 1/4$, where $u_j = d_j^{-1}$. This shuffle has holding probability 0, so there is no easy bound on $\kappa_{n!}$. The easiest way out is to see to that all $|y|$ are odd and then use the third statement of Lemma 4.11.

Again compare with random transpositions. How to choose odd representations depends on whether n is even or odd. We assume here that n is even and leave the other case to the reader. First note that for all earlier examples, we implicitly used the obvious empty representation of $y = id$. However, since that representation has length 0, this does not work now. Use instead the representation $id = d_{n-1}^{n-1}$. Now assume that $y = (i\ j)$ with $i < j$. Then we may write $|y|$ as $u_n^{n-j} u_{n-1}^{j-i-1} u_n u_{n-1}^{n+i-j-1} d_n^{n-j}$. This representation has the odd length $3n - 2j - 1 < 3n$. Using also the trivial bound $N(x, y) \leq |y|$ gives $A \leq 36n^2$.

Now the third statement of Lemma 4.11 gives

$$\hat{\tau} \leq (1 + o(1))36n^3 \log n.$$

This is also the same order as for total variation. \square

Example. The symmetrized overlapping cycles shuffle. The overlapping cycles shuffle gives probability $1/2$ to each of u_n and $u_m = (1\ 2\ \dots\ m)$ for some $m > n$. This turned out to be a difficult shuffle to handle. Above, some effort was paid to the case $m = n/2$, n even, m odd. It was shown that the relaxation time for a single card is $\Theta(n^2)$. Here we compare its symmetrized version with the random transpositions shuffle to give an upper bound in L^2 of order $n^3 \log n$. Let $\nu(u_n) = \nu(d_n) = \nu(u_m) = \nu(d_m) = 1/4$. We have $id = u_m^m$ an odd-length representation of length $n/2$. No take $y = (i\ j)$, $i < j$. Write $x = d_n^m u_m u_n^{m-1} u_m$. Then $|x| = n + 1$ which is odd and $x = (1\ n)$. If $j \leq m$, then $y = vxv^{-1}$ where $v = u_m^{m-j+1} d_n u_m^{j-i+1}$. Hence $|y| = 2n + 5$ which is odd. If $j > m$, add a prefix to v in order to bring i and j to the upper half of deck. This takes a length of at most $m + 1$ giving $|y| \leq 3n + 7 < 4n$ and still odd. Mimicking the above example gives $A \leq 64n^2$ and

$$\hat{\tau} \leq (1 + o(1))64n^3 \log n.$$

\square

Example. Random walk on \mathbb{Z}_n . Let $a \leq \sqrt{n}$ and consider the random walk on \mathbb{Z}_n generated by $\nu(j) = 1/(2a + 1)$, $j = -a, \dots, a$, i.e. steps are uniform over $\{-a, \dots, a\}$. Compare with the trivial random walk that gets uniform over \mathbb{Z}_n in a single step. Identify each element $y \in \mathbb{Z}_n$ with $y \in [n]$ and represent k as

$$y = 0 + a + 1 + (a - 1) + 2 + (a - 2) + \dots + a + 0 + (a - 1) + 1 + \dots + b.$$

Here $b = y - a\lfloor y/a \rfloor$. Then $|y| \leq 2y/a + 1 \leq 2n/a + 1 < 3n/a$ and for each $x \in \{-a, \dots, a\}$, $N(x, y) \leq 4y/a^2 + 1 < 6n/a^2$. Hence

$$A < (2a + 1) \frac{3n}{a} \frac{6n}{a^2} < \frac{54n^2}{a^2}.$$

Since $\kappa_2^b = 0$, Lemma 4.9 reveals that $\kappa_2 \leq 1 - a^2/(54n^2)$, so that the relaxation time $O(a^2/n^2)$. This is the correct order. Since $\kappa_n \geq -1 + 2/(2a + 1)$ by Lemma 4.13, it follows from Lemma 4.11 that $\hat{\tau} = O((n^2/a^2) \log n)$. \square

5 The Thorp shuffle

The Thorp shuffle, proposed by E. Thorp in 1973, is known as the longest standing open problem in card shuffling, having resisted numerous attacks from prominent mathematicians for over three decades. The Thorp shuffle is a very fine riffle shuffle on a deck of an even number of cards: cut the deck in two equal sized packs, one in each hand. Now drop the two bottom cards in an order determined by a fair coin flip. Then repeat this for the two new bottom cards (i.e those two cards that were second to bottom to begin with). Keep on repeating this until the packs are empty.

The conjecture is that the mixing time should be of order $\log n$, or, in the very least, $O((\log n)^2)$. However, until very recently, the best known upper bound was of order n . The first upper bound polynomial in $\log n$, was given by Morris [15], 2005, and was of order $O((\log n)^{44})$, under the condition that $n = 2^d$ for integer d . Montenegro and Tetali [13] were able to build on this to improve the bound to $O((\log n)^{29})$, under the same condition. Then in 2008, representing the present state of the art, Morris [14] made significant progress and were able to remove the $n = 2^d$ condition and took the upper bound down to $O((\log n)^4)$. In this section, we will follow the latter paper.

Write $n = 2h$ for the number of cards. Formally, the Thorp shuffle can be described by first, for each of the pairs $(1, h + 1), (2, h + 2), \dots, (h, 2h)$, make a transposition of the two cards in the pair with probability $1/2$, independently of other pairs, and second, composing with the permutation moving the card in position i to position $2i - 1$ and the card in position $h + i$ to position $2i$, $i = 1, \dots, h$.

Morris makes use of an advanced entropy technique. For this, we first need some theoretical background.

Definition 5.1 *Let ν and π be two probability measures on the finite space S . The relative entropy of ν with respect to π is given by*

$$\text{ENT}(\nu||\pi) = \sum_{x \in S} \nu(x) \log \frac{\nu(x)}{\pi(x)}.$$

In the special case when π is the uniform distribution on S , one simply speaks of the relative entropy of ν , denoted $\text{ENT}(\nu)$. When X is a random variable on S distributed according to ν , one also writes $\text{ENT}(X)$ for $\text{ENT}(\nu)$. We have

$$\text{ENT}(\nu) = \sum_{x \in S} \nu(x) \log(|S|\nu(x)).$$

Here are a few observations:

- $\text{ENT}(\nu||\pi) \geq 0$, with equality if and only if $\nu = \pi$.
- $\text{ENT}(\nu) = \log |S| - H(\nu)$, where $H(\nu) = - \sum_{x \in S} \nu(X) \log \nu(X)$ is the usual absolute entropy.
- $\text{ENT}(\nu||\pi) = \mathbb{E}_\nu[\log \frac{\nu(X)}{\pi(X)}]$.

The two last observations are immediate. The first one follows from the last, using Jensen's inequality on the convex function $-\log(\pi(X)/\nu(X))$. The following lemma relates total variation distance to relative entropy.

Lemma 5.1 *Let π be the uniform measure. Then*

$$\|\nu - \pi\|_{TV} \leq \sqrt{\frac{1}{2} \text{ENT}(\nu)}.$$

Proof. By Schwarz' inequality,

$$\|\nu - \pi\|_{TV} = \sum_i \frac{1}{2} \left| \nu(i) - \frac{1}{|S|} \right| \leq \sqrt{\frac{1}{4} |S| \sum_i \left| \nu(i) - \frac{1}{|S|} \right|^2}.$$

Hence it suffices to show that

$$\frac{1}{2} |S| \sum_i \left| \nu(i) - \frac{1}{|S|} \right|^2 - \sum_i \nu(i) \log(|S| \nu_i) \leq 0.$$

This is a standard optimization problem over the $\nu(i)$'s, that can e.g. be solved using a Lagrange multiplier or setting $\nu(|S|) = 1 - \nu(1) - \dots - \nu(|S| - 1)$. \square

The conditional entropy of a random variable X given $Y = y$, denoted

$$\text{ENT}(X|Y = y)$$

is of course the entropy of the distribution of X given $Y = y$. By $\text{ENT}(X|Y)$ one means the corresponding random variable (a function of Y). Now let μ and ν be two probability measures on $S \times S$. Think of these as distributions of a pair of random variables. Write

$$\mu_1(i) = \sum_{j \in S} \mu(i, j), \quad \mu_2(j) = \sum_{i \in S} \mu(i, j)$$

for the marginal distributions with respect to μ , and correspondingly for ν . Write also

$$\mu_{2|1}(j|i) = \frac{\mu(i, j)}{\mu_1(i)}$$

as usual for the conditional distribution, and correspondingly for ν . Then we have the following chain rule for relative entropies:

$$\text{ENT}(\mu \parallel \nu) = \text{ENT}(\mu_1 \parallel \nu_1) + \sum_i \sum_j \mu(i, j) \log \frac{\mu_{2|1}(j|i)}{\nu_{2|1}(j|i)}.$$

This follows from a straightforward manipulation of the right hand side according to the definitions of the involved quantities. In the special case when ν is uniform on $S \times S$, we have that μ_1 and $\mu_{2|1}$ are both uniform on S , and the chain rule becomes

$$\begin{aligned} \text{ENT}(\mu) &= \text{ENT}(\mu_1) + \sum_i \sum_j \mu(i, j) \log(|S| \mu_{2|1}(j|i)) \\ &= \text{ENT}(\mu_1) + \sum_i \mu_1(i) \sum_j \mu_{2|1}(j|i) \log(|S| \mu_{2|1}(j|i)). \end{aligned}$$

If (X, Y) is distributed according to μ , then this in turn becomes

$$\begin{aligned} \text{ENT}(X, Y) &= \text{ENT}(X) + \mathbb{E}[\text{ENT}(Y|X)] \\ &= \text{ENT}(X) + \mathbb{E}[\text{ENT}(X, Y|X)]. \end{aligned}$$

An inductive generalization to higher dimensions leads to that, for any $i \in [n]$,

$$\begin{aligned} \text{ENT}(X_1, X_2, \dots, X_n) &= \mathbb{E}[\text{ENT}(X_1, X_2, \dots, X_n | X_i, X_{i+1}, \dots, X_n)] \\ &\quad + \sum_{k=i}^n \mathbb{E}[\text{ENT}(X_k | X_{k+1}, X_{k+2}, \dots, X_n)]. \end{aligned}$$

(Note that the last term in the sum is $\text{ENT}(X_n)$.) The situation in which the chain rule will be used here, is of course for random permutations. Let ν be a random permutation in S_n , let $\mathcal{F}_j = \sigma(\nu^{-1}(j), \nu^{-1}(j+1), \dots, \nu^{-1}(n))$, the σ -algebra generated by the identities of the cards in positions $j, j+1, \dots, n$, and let

$$E_j = \mathbb{E}[\text{ENT}(\nu^{-1}(j) | \mathcal{F}_{j+1})],$$

$j = 1, 2, \dots, n$. Then the chain rule takes on the form

$$\text{ENT}(\nu) = \mathbb{E}[\text{ENT}(\nu|\mathcal{F}_i)] + \sum_{k=i}^n E_k.$$

In order to bound relative entropies, we need to introduce yet another concept of “distance” between probability measures. Let d be a function on pairs of nonnegative numbers given by

$$d(x, y) = \frac{1}{2}x \log x + \frac{1}{2}y \log y - \frac{x+y}{2} \log \frac{x+y}{2}.$$

Since the function $x \rightarrow x \log x$ is strictly convex, Jensen’s inequality (applied to a random variable taking on the values x and y with probability $1/2$ each) entails that $d(x, y) \geq 0$ with equality if and only if $x = y$. Also, regarding d as a function of only one of its arguments and differentiating twice, shows that

Lemma 5.2 *The functions $d(x, \cdot)$ and $d(\cdot, y)$ are convex for all x and y .*

Definition 5.2 *Let μ and ν be two probability measures on S . Then*

$$d(\mu, \nu) = \sum_{i \in S} d(\mu(i), \nu(i)).$$

It should be noted that this notion of distance is not a proper metric, since it does not satisfy the triangle inequality.

An alternative and sometimes more convenient expression for d is

$$d(x, y) = \frac{x+y}{2} f\left(\frac{x-y}{x+y}\right),$$

where $f(u) = \frac{1}{2}(1+u) \log(1+u) + \frac{1}{2}(1-u) \log(1-u)$. The function f is well defined on $(-1, 1)$ and by differentiating twice, it follows that f is convex. The following lemma reveals that lumping decreases d .

Lemma 5.3 *Let X and Y be distributed according to μ and ν respectively. Let $g : S \rightarrow \mathbb{R}$ and let M and N be the distributions of $g(X)$ and $g(Y)$ respectively. Then*

$$d(M, N) \leq d(\mu, \nu).$$

Proof. Partition S by letting $S_i = \{x \in S : g(x) = i\}$, $i \in g(S)$. Then

$$\begin{aligned} d(\mu, \nu) &= \sum_i \sum_{x \in S_i} d(\mu(x), \nu(x)) = \sum_i \sum_{x \in S_i} \frac{\mu(x) + \nu(x)}{2} f\left(\frac{\mu(x) - \nu(x)}{\mu(x) + \nu(x)}\right) \\ &= \sum_i \frac{M(i) + N(i)}{2} \sum_{x \in S_i} \frac{\mu(x) + \nu(x)}{2} \frac{2}{M(i) + N(i)} f\left(\frac{\mu(x) - \nu(x)}{\mu(x) + \nu(x)}\right) \\ &= \sum_i \frac{M(i) + N(i)}{2} \mathbb{E}f(Z_i), \end{aligned}$$

where Z_i is a random variable that takes on the value $(\mu(x) - \nu(x))/(\mu(x) + \nu(x))$ with probability $(\mu(x) + \nu(x))/(M(i) + N(i))$, $x \in S_i$. By the convexity of f and Jensen's inequality,

$$\mathbb{E}f(Z_i) \geq f(\mathbb{E}Z_i) = f\left(\frac{M(i) - N(i)}{M(i) + N(i)}\right).$$

Inserting into the above gives

$$d(\mu, \nu) \geq \sum_i \frac{M(i) + N(i)}{2} f\left(\frac{M(i) - N(i)}{M(i) + N(i)}\right) = d(M, N).$$

□

Here is the result that captures how d relates to relative entropy.

Lemma 5.4 *There is a constant c , independent of μ and S , such that*

$$\text{ENT}(\mu) \leq c \log |S| d(\mu, \pi),$$

where π is uniform on S .

Proof. Let $n = |S|$ and let $g(x) = x \log x - x + 1$. Then

$$\text{ENT}(\mu) = \frac{1}{n} \sum_i n\mu(i) \log(n\mu(i)) = \frac{1}{n} \sum_i g(n\mu(i))$$

where the second equality follows from the fact that $\sum_i (-n\mu(i) + 1) = 0$. Since $d(ax, ay) = ad(x, y)$,

$$\begin{aligned} d(\mu, \pi) &= \sum_i d(\mu(i), \frac{1}{n}) = \frac{1}{n} \sum_i d(n\mu(i), 1) \\ &= \frac{1}{n} \sum_i \frac{n\mu(i) + 1}{2} f\left(\frac{n\mu(i) - 1}{n\mu(i) + 1}\right). \end{aligned}$$

Hence it suffices to show that for each i ,

$$g(n\mu(i)) \leq c \log n \frac{n\mu(i) + 1}{2} f\left(\frac{n\mu(i) - 1}{n\mu(i) + 1}\right).$$

Letting

$$R(x) = \frac{g(x)}{\frac{x+1}{2} f\left(\frac{x-1}{x+1}\right)},$$

this amounts to showing that $R(x) \leq c \log n$ for $x \in [0, n]$. Using Taylor's formula, one can show that $\lim_{x \rightarrow 1} R(x)$ exists. Hence R extends to a continuous function. As a consequence, $\sup_{x \in [0, 2]} R(x) < \infty$. Writing

$$f\left(\frac{x-1}{x+1}\right) = \log\left(\frac{2x}{x+1}\right) - \frac{1}{x+1} \log x,$$

we see that the denominator of $R(x)$ is increasing on $[2, n]$, so that on $[2, n]$, the denominator exceeds $\frac{1}{2}x f(1/3)$. Since $g(x) < x \log x \leq x \log n$, we have that on $[2, n]$,

$$R(x) < \frac{2}{f(1/3)} \log n.$$

□

Let $\mu \in S_n$ be a random permutation. Then, if ν is a fixed permutation, clearly

$$\text{ENT}(\mu\nu) = \text{ENT}(\mu),$$

since the effect of composing with ν is just to rename the elements of S_n .

Remark. Here one should beware of the risk of possible confusion; earlier we used greek letters for probability measures and capitals X, Y etc, for the corresponding random variables. From now on, greek letters will denote random permutations, i.e. random variables in S_n , and, when needed, the probability measure corresponding to μ , say, will be denoted $\mathcal{L}(\mu)$. Recall then that we identify $\text{ENT}(\mu)$ with $\text{ENT}(\mathcal{L}(\mu))$.

Lemma 5.5 *If μ and ν are two random permutations, then*

$$\text{ENT}(\mu\nu) \leq \text{ENT}(\mu).$$

Proof. Recall the general form of the chain rule for relative entropies:

$$\text{ENT}(M\|N) = \text{ENT}(M_1\|N_1) + \sum_{i \in S} \sum_{j \in S} M(i, j) \log \left(\frac{M_{2|1}(j|i)}{N_{2|1}(j|i)} \right),$$

where M and N are two probability measures on $S \times S$, S finite. Note that the second term is

$$\begin{aligned} -\mathbb{E}_M \log \frac{N_{2|1}(j|i)}{M_{2|1}(j|i)} &\geq -\log \mathbb{E}_M \frac{N_{2|1}(j|i)}{M_{2|1}(j|i)} \\ &= -\log \sum_i \sum_j M(i, j) \frac{N(i, j)/N_1(i)}{M(i, j)/M_1(i)} \\ &= -\log 1 = 0, \end{aligned}$$

where the inequality is Jensen's inequality and the equalities are easy algebraic manipulation. Hence $\text{ENT}(M\|N) \geq \text{ENT}(M_i\|N_i)$, $i = 1, 2$. Now apply this with $M = \mathcal{L}(\mu, \mu\nu)$ and $N = \mathcal{L}(\phi, \pi\nu)$, where π is a uniform random permutation. Then clearly N_1 and N_2 are both uniform on S_n . Since $M_{2|1}(j|i) = N_{2|1}(j|i)$ for all (i, j) , the second term in the chain rule vanishes and we have that

$$\text{ENT}(M\|N) = \text{ENT}(M_1\|N_1) = \text{ENT}(\mu).$$

On the other hand

$$\text{ENT}(M\|N) \geq \text{ENT}(M_2\|N_2).$$

The result follows. \square

We are now ready to state the key theorem of this section. Let $a, b \in [n]$ and let $c(a, b)$ denote the random permutation that equals id with probability $1/2$ and $(a b)$ with probability $1/2$; we will call such a random permutation a *collision* of the positions a and b .

Let ν be a random permutation and suppose that ν is written on the form

$$\nu = \theta c(a_1, b_1) c(a_2, b_2) \dots c(a_k, b_k)$$

where θ is a random or fixed permutation, the a_i 's and b_i 's are all distinct and the collisions are mutually independent given θ . Let ν_1, ν_2, \dots be iid copies of ν . Write

$$\nu_{(t)} = \nu_1 \nu_2 \dots \nu_t.$$

Say that the cards x and y collide at time t if $\nu_{(t)}^{-1}(i) = x$, $\nu_{(t)}^{-1}(j) = y$ for some positions i and j , and ν_t contains the collision $c(i, j)$.

Fix t and let $T \in [t]$ be a random variable independent of the ν_i :s. For each card x , let $b(x)$ be the first card that x collides with in the time interval $[T, t]$, with $b(x) = x$ if no such card exists. Let

$$m(x) = b(x)$$

if $b(b(x)) = x$, with $m(x) = x$ otherwise. (I.e. $m(x) = y$ and $m(y) = x$ if x and y collide in the time interval $[T, t]$ and none of x and y collide with any other card earlier in that time interval.) For each x , let A_x be the largest number such that $\mathbb{P}(m(x) = y) \geq A_x/x$ for all $y \leq x$.

Theorem 5.1 *Let μ be a random permutation independent of $\nu_{(t)}$. Then, with the above notation,*

$$\text{ENT}(\mu\nu_{(t)}) - \text{ENT}(\mu) \leq -\frac{C}{\log n} \sum_{k=1}^n A_k E_k$$

where the constant C is independent of n , μ , t , T and the ν_i :s and

$$E_k = \mathbb{E}[\text{ENT}(\mu^{-1}(k)|\mathcal{F}_{k+1})]$$

as usual.

Proof. Let $\mathcal{M} = (m(1), m(2), \dots, m(n))$ and let ρ be the random permutation

$$\rho = \prod_{i:m(i) \leq i} c(i, m(i)),$$

independent of μ and $\nu_{(t)}$, given \mathcal{M} . Then $\rho\nu_{(t)}$ and $\nu_{(t)}$ have the same distribution. Hence

$$\text{ENT}(\mu\nu_{(t)}) - \text{ENT}(\mu) = \text{ENT}(\mu\rho\nu_{(t)}) - \text{ENT}(\nu).$$

Observe that

$$\text{ENT}(\mu\rho\nu_{(t)}|\mathcal{M}, \nu_{(t)}) = \text{ENT}(\mu\rho|\mathcal{M}, \nu_{(t)}) = \text{ENT}(\mu\rho|\mathcal{M}).$$

Now let

$$\rho_k = \prod_{i:m(i) \leq i \leq k} c(i, m(i)),$$

$k = 0, 1, 2, \dots, n$, and note that $\rho_0 = id$ and $\rho_n = \rho$. Since μ is independent of \mathcal{M} , $\text{ENT}(\mu|\mathcal{M}) = \text{ENT}(\mu)$ a.s. so

$$\begin{aligned} \text{ENT}(\mu\rho\nu_{(t)}|\mathcal{M}, \nu_{(t)}) - \text{ENT}(\mu) &= \text{ENT}(\mu\rho|\mathcal{M}) - \text{ENT}(\mu|\mathcal{M}) \\ &= \text{ENT}(\mu\rho_n|\mathcal{M}) - \text{ENT}(\mu\rho_0|\mathcal{M}) \\ &= \sum_{k=1}^n \left(\text{ENT}(\mu\rho_k|\mathcal{M}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{M}) \right). \end{aligned}$$

By Lemma 5.5, $\text{ENT}(\mu\rho\nu_{(t)}|\mathcal{M}, \nu_{(t)}) \geq \text{ENT}(\mu\rho\nu_{(t)})$, so we will be done if we can show that for each k ,

$$\mathbb{E}[\text{ENT}(\mu\rho_k|\mathcal{M}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{M})] \leq -\frac{C}{\log n} A_k E_k.$$

On the event $\{m(k) > k\}$ we have that $\rho_k = \rho_{k-1}$ and so on this event

$$\text{ENT}(\mu\rho_k|\mathcal{M}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{M}) = 0.$$

On $\{m(k) \leq k\}$,

$$\mathcal{L}(\mu\rho_k) = \frac{1}{2}\mathcal{L}(\mu\rho_{k-1}) + \frac{1}{2}\mathcal{L}(\mu\rho_{k-1}(k \ m(k))).$$

Fix $i \leq k$, let $\lambda = \mu\rho_{k-1}$ and $\xi = \lambda(i \ k)$. The difference between the permutations λ and ξ is thus that $\lambda^{-1}(i) = \xi^{-1}(k)$ and vice versa. Note also that $\lambda^{-1}(j) = \xi^{-1}(j) = \mu^{-1}(j)$ for all $j \geq k+1$. Let $\mathcal{G}_j = \sigma(\mathcal{F}_j, \mathcal{M})$, $j \in [n]$. Since $\text{ENT}(\lambda|\mathcal{G}_{k+1}) = \text{ENT}(\xi|\mathcal{G}_{k+1})$, we have that

$$\begin{aligned} &\text{ENT}(\lambda c(i, k)|\mathcal{G}_{k+1}) - \text{ENT}(\lambda|\mathcal{G}_{k+1}) \\ &= \text{ENT}\left(\frac{1}{2}\mathcal{L}(\lambda|\mathcal{G}_{k+1}) + \frac{1}{2}\mathcal{L}(\xi|\mathcal{G}_{k+1})\right) - \frac{1}{2}\text{ENT}(\lambda|\mathcal{G}_{k+1}) - \frac{1}{2}\text{ENT}(\xi|\mathcal{G}_{k+1}) \\ &= -d(\mathcal{L}(\lambda|\mathcal{G}_{k+1}), \mathcal{L}(\xi|\mathcal{G}_{k+1})) \end{aligned}$$

by the definitions of relative entropy and d . By Lemma 5.3,

$$\begin{aligned} d(\mathcal{L}(\lambda|\mathcal{G}_{k+1}), \mathcal{L}(\xi|\mathcal{G}_{k+1})) &\geq d(\mathcal{L}(\lambda^{-1}(k)|\mathcal{G}_{k+1}), \mathcal{L}(\xi^{-1}(k)|\mathcal{G}_{k+1})) \\ &= d(\mathcal{L}(\lambda^{-1}(k)|\mathcal{G}_{k+1}), \mathcal{L}(\lambda^{-1}(i)|\mathcal{G}_{k+1})). \end{aligned}$$

Applying this with $i = m(k)$ on $\{m(k) \leq k\}$ yields

$$\begin{aligned}
\text{ENT}(\mu\rho_k|\mathcal{G}_{k+1}) &- \text{ENT}(\mu\rho_{k-1}|\mathcal{G}_{k+1}) \\
&= \text{ENT}(\lambda c(k, m(k))|\mathcal{G}_{k+1}) - \text{ENT}(\lambda|\mathcal{G}_{k+1}) \\
&\leq -d(\mathcal{L}(\lambda^{-1}(k)|\mathcal{G}_{k+1}), \mathcal{L}(\lambda^{-1}(m(k))|\mathcal{G}_{k+1})) \\
&= -\sum_{i=1}^k \mathbf{1}_{\{m(k)=i\}} d(\mathcal{L}(\mu^{-1}(k)|\mathcal{G}_{k+1}), \mathcal{L}(\mu^{-1}(i)|\mathcal{G}_{k+1})) \\
&= -\sum_{i=1}^k \mathbf{1}_{\{m(k)=i\}} d(\mathcal{L}(\mu^{-1}(k)|\mathcal{F}_{k+1}), \mathcal{L}(\mu^{-1}(i)|\mathcal{F}_{k+1}))
\end{aligned}$$

where the second last equality follows from the fact that $\lambda = \mu\rho_{k-1}$ does not have the collision $c(k, m(k))$ and the last equality follows from that μ is independent of \mathcal{M} . Taking expectations now gives

$$\begin{aligned}
&\mathbb{E}\left[\text{ENT}(\mu\rho_k|\mathcal{G}_{k+1}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{G}_{k+1})\right] \\
&\leq -\sum_{i=1}^k \mathbb{P}(m(k) = i) \mathbb{E}\left[d(\mathcal{L}(\mu^{-1}(i)|\mathcal{F}_{k+1}), \mathcal{L}(\mu^{-1}(k)|\mathcal{F}_{k+1}))\right] \\
&\leq -\frac{A_k}{k} \sum_{i=1}^k \mathbb{E}\left[d(\mathcal{L}(\mu^{-1}(i)|\mathcal{F}_{k+1}), \mathcal{L}(\mu^{-1}(k)|\mathcal{F}_{k+1}))\right] \\
&\leq -A_k \mathbb{E}\left[d\left(\frac{1}{k} \sum_{i=1}^k \mathcal{L}(\mu^{-1}(i)|\mathcal{F}_{k+1}), \mathcal{L}(\mu^{-1}(k)|\mathcal{F}_{k+1})\right)\right] \\
&= -A_k \mathbb{E}\left[d(U, \mathcal{L}(\mu^{-1}(k)|\mathcal{F}_{k+1}))\right],
\end{aligned}$$

where U is the uniform distribution on $[n] \setminus \{\mu^{-1}(k+1), \mu^{-1}(k+2), \dots, \mu^{-1}(n)\}$. Here the second inequality follows from the assumption $\mathbb{P}(m(k) = i) \geq A_k/k$ and the third inequality from the convexity of $d(\cdot, y)$. By Lemma 5.4,

$$\begin{aligned}
d(U, \mathcal{L}(\mu^{-1}(k)|\mathcal{F}_{k+1})) &\geq \frac{C}{\log(n-k)} \text{ENT}(\mu^{-1}(k)|\mathcal{F}_{k+1}) \\
&= \frac{CE_k}{\log(n-k)} \geq \frac{CE_k}{\log n}.
\end{aligned}$$

Inserting into the above yields

$$\mathbb{E}\left[\text{ENT}(\mu\rho_k|\mathcal{G}_{k+1}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{G}_{k+1})\right] \leq \frac{-CA_k E_k}{\log n}.$$

This result differs from what we want to prove only in that the conditioning is on \mathcal{G}_{k+1} and not on \mathcal{M} as desired. However, by the chain rule, for any $i \leq k$,

$$\text{ENT}(\mu\rho_i|\mathcal{M}) = \mathbb{E}[\text{ENT}(\mu\rho_i|\mathcal{G}_{k+1})|\mathcal{M}] + \sum_{j=k+1}^n \mathbb{E}[\text{ENT}((\mu\rho_i)^{-1}(j)|\mathcal{G}_{j+1})|\mathcal{M}]$$

and the last sum of this expression is independent of i , since $(\mu\rho_i)^{-1}(j) = \mu^{-1}(j)$ for all j 's of the sum. Hence, applying this for $i = k - 1$ and $i = k$ and taking the difference gives

$$\begin{aligned} & \mathbb{E} \left[\text{ENT}(\mu\rho_k|\mathcal{G}_{k+1}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{G}_{k+1}) \right] \\ &= \mathbb{E} \left[\text{ENT}(\mu\rho_k|\mathcal{M}) - \text{ENT}(\mu\rho_{k-1}|\mathcal{M}) \right]. \end{aligned}$$

□

Next we attack the Thorp shuffle. We will use that the inverse shuffle has the same mixing time as the shuffle itself. For convenience, denote the positions in the deck $0, 1, 2, \dots, n - 1$. Then one step of the inverse Thorp shuffle can be written as

$$\nu = \theta c(0, h) c(1, h + 1) \dots c(h - 1, 2h - 1)$$

where $h = n/2$ and θ is the fixed permutation taking the card in position $2i$ to position i and the card in position $2i + 1$ to position $h + i$, $i = 0, 1, \dots, h - 1$.

Partition the set of positions into intervals such that $I_0 = \{0\}$ and

$$I_m = \{2^{m-1}, 2^{m-1} + 1, \dots, 2^m - 1\} \cap [n - 1],$$

$m = 1, 2, \dots, \lceil \log_2 n \rceil$. Fix such an m arbitrarily. We will now show that the conditions of Theorem 5.1 are satisfied with $t = \lceil \log_2 n \rceil$, the random variable T having distribution $\mathbb{P}(T = 0) = 1/2^m$ and

$$\mathbb{P}(T = r) = \frac{1}{2^{m+1-r}},$$

$r = 1, 2, \dots, m$ (i.e. $m + 1 - T$ is geometric with parameter $1/2$, truncated at m) and $A_i = 1/8$ for $i \in I_m$. This means to show that

$$\mathbb{P}(m(i) = j) \geq \frac{1}{8i}$$

for all $i \in I_m$ and $j = 0, 1, 2, \dots, i - 1$. Fix such i and j .

Let $f(t) = \lfloor t/2 \rfloor$ and let $X_s(j)$ be the position of card j at time s . Then

$$X_s(j) = f(X_{s-1}(j)) + Z_s(j)$$

where the $Z_s(j)$'s, $s = 1, 2, \dots, j = 0, 1, 2, \dots, n - 1$, are independent random variables taking on the value 0 or h , each with probability $1/2$. An elementary observation about the function f is that for $x > y$,

$$f(x) - f(y) \leq \frac{x - y}{2} + \frac{1}{2}$$

with equality if and only if x is even and y is odd. Hence

$$\lceil \log_2(f(x) - f(y)) \rceil \leq \lceil x - y \rceil - 1$$

unless x is even and $y = x - 1$ in which case the left hand side and the right hand side both equal 1. Write (uniquely)

$$h = 2^k l$$

where l is odd.

Lemma 5.6 *With the above notation,*

$$\mathbb{P}(X_m(j) \text{ is even}) \geq \frac{1}{2}.$$

Proof. Since $j < i < 2^m$, $f^m(j) = 0$. Hence $X_m(j)$ is completely determined by $Z_1(j), Z_2(j), \dots, Z_m(j)$. If $m \leq k$, we simply have that

$$X_m(j) = \sum_{r=0}^{k-1} 2^{-r} Z_{m-r}(j).$$

Since every $Z_r(j)$ is either 0 or $2^k l$, every term of the sum is either 0 or $2^{k-r} l$ which is even since $r < k \leq m$. Hence $X_m(j)$ is deterministically even if $m \leq k$.

Now consider the case $m > k$. For $s = m - k, m - k + 1, \dots, m$, let $X'_s(j)$ be the position that card j would have had at time s , had Z_{m-k} been different. Then obviously $X'_s(j)$ and $X_s(j)$ have the same distribution. We have that

$$|X_s(j) - X'_s(j)| = 2^{m-s} l, \quad s = m - k, m - k + 1, \dots, m$$

and, in particular, $|X_m(j) - X'_m(j)| = l$, which is odd. Hence exactly one of $X_m(j)$ and $X'_m(j)$ is even, and so

$$\mathbb{P}(X_m(j) \text{ is even}) = \frac{1}{2}.$$

□

Now, let $y_0 = i = X_0(i)$ and recursively define

$$y_s = f(y_{s-1}) + X_s(j),$$

i.e. let y_s be the position that card i would have had at time s , had $Z_s(i)$ equalled $Z_s(j)$ for every s . Let τ be the first time s that $|y_s - X_s(j)| = 1$ and $X_s(j)$ is even.

Since $|i - j| < 2^m$ and $j < i$, there is an $s \leq m$ such that $y_s - X_s(j) = 1$. By the elementary observations about f above, $y_u - X_u(j) = 1$ for all $u = s, s+1, \dots, \tau$, since $X_u(j)$ is odd for all these u . Since Lemma 5.6 tells us that $X_m(j)$ is even with probability at least $1/2$, we get

$$\mathbb{P}(\tau \leq m) \geq \frac{1}{2}.$$

Clearly

$$\mathbb{P}(X_0(i) = y_0, X_1(i) = y_1, \dots, X_r(i) = y_r | \tau = r) = \frac{1}{2^r}.$$

Since $\mathbb{P}(T = r) \geq 1/2^{m+1-r}$ and T is independent of τ and all $Z_s(j)$'s, this entails that

$$\begin{aligned} \mathbb{P}(m(i) = j) &\geq \sum_{r=0}^m \mathbb{P}(T = r, \tau = r, X_0(i) = y_0, X_1(i) = y_1, \dots, X_r(i) = y_r) \\ &\geq \sum_{r=0}^m 2^{-(m+1-r)} 2^{-r} \mathbb{P}(\tau = r) \\ &= 2^{-(m+1)} \mathbb{P}(\tau \leq m) \\ &\geq \frac{1}{2^{m+2}} \geq \frac{1}{8i} \end{aligned}$$

where the last inequality follows from that $i \geq 2^{m-1}$.

We can now use the key theorem with the given t , T and A_k . Doing so will finally lead to the following result for the Thorp shuffle.

Theorem 5.2 *For the Thorp shuffle,*

$$\tau_{\text{mix}} = O((\log n)^4).$$

Proof. Continue with the same notation as above. By the chain rule,

$$\sum_{i=1}^n E_i = \sum_m \sum_{i \in I_m} E_i = \text{ENT}(\mu)$$

Since there are at most $\log_2 n + 2 < C \log n$ intervals I_m , we must have

$$\sum_{i \in I_{m^*}} E_i \geq \frac{1}{C \log n} \text{ENT}(\mu),$$

where m^* is the index that maximizes the inner sum. Therefore, by the key theorem and the above with $m = m^*$,

$$\begin{aligned} \text{ENT}(\mu\nu_{(t)}) &\leq \text{ENT}(\mu) - \frac{C'}{\log n} \sum_{k=1}^n A_k E_k \\ &\leq \text{ENT}(\mu) - \frac{C'}{\log n} \sum_{i \in I_{m^*}} A_i E_i \\ &\leq \left(1 - \frac{C'}{8C(\log n)^2}\right) \text{ENT}(\mu). \end{aligned}$$

Now iterate this result with $\mu = \nu_{(Bt(\log n)^3)}$, $B = 0, 1, 2, \dots$, and get

$$\begin{aligned} \text{ENT}(\nu_{(Bt(\log n)^3)}) &\leq \left(1 - \frac{c}{(\log n)^2}\right)^{B(\log n)^3} \text{ENT}(id) \\ &\leq n^{-Bc} \log(n!) \\ &\leq n^{1-Bc} \log n \leq \frac{1}{8} \end{aligned}$$

for large n as soon as, say, $B \geq 2/c$. By Lemma 5.1, for such B and n ,

$$\|\mathbb{P}(\nu_{(Bt(\log n)^3)} \in \cdot) - \pi\| \leq \sqrt{\frac{1}{2} \text{ENT}(\nu_{(Bt(\log n)^3)})} \leq \frac{1}{4}.$$

Since $t = (1 + o(1)) \log_2 n = O(\log n)$, this means that

$$\tau_{\text{mix}} = O((\log n)^4).$$

□

REFERENCES

1. D. Aldous, Random walks on finite groups and rapidly mixing Markov chains, *Seminaire de Probabilites*, Lecture notes in Math. **986** (1983), 243-297.
2. D. Aldous and P. Diaconis, Shuffling cards and stopping times, *Amer. Math Monthly* **93** (1986), 333-348.
3. D. Aldous and J. Fill, "Reversible Markov chains and random walks on graphs," draft book at www.stat.berkeley.edu/~aldous/RWG/book.html
4. O. Angel, Y. Peres and D. B. Wilson, Card shuffling and diophantine approximation, *Ann. Appl. Probab.* **18** (2008), 1215-1231.
5. D. Bayer and P. Diaconis, Trailing the dove-tail shuffle to its lair, *Ann. Appl. Probab.* **2** (1992), 294-313.
6. P. Diaconis and L. Saloff-Coste, Comparison techniques for random walk on finite groups, *Ann. Probab.* **21** (1993), 2131-2156.
7. P. Diaconis and L. Saloff-Coste, Comparison theorems for reversible Markov chains, *Ann. Appl. Probab.* **3** (1993), 696-730.
8. P. Diaconis and M. Shashahani, Generating a random permutation with random transpositions, *Z. Wahrsch. Verw. Gebiete* **57** (1981), 159-179.
9. R. Durrett, "Probability: Theory and Examples," Wadsworth & Brooks/Cole, 1991.
10. J. Jonasson, Biased random-to-top shuffling, *Ann. Appl. Probab.* **16** (2006), 1034-1058.
11. J. Jonasson, The overhand shuffle mixes in $\Theta(n^2 \log n)$ steps, *Ann. Appl. Probab.* **16** (2006), 231-243.
12. P. Matthews, A strong uniform time for random transpositions, *J. Theoret. Probab.* **4** (1988), 411-423.
13. R. Montenegro and P. Tetali, "Mathematical aspects of mixing time in Markov chains," Foundations and Trends in Theoretical Computer Science, NOW Publishers 2006.

14. B. Morris, Improved mixing time for the Thorp shuffle and L-reversal chain, to appear in *Ann. Probab.*
15. B. Morris, The mixing time of the Thorp shuffle, *SIAM Journal on Computing*, *STOC 2005 special issue*.
16. B. Morris and Y. Peres, Evolving sets, mixing and heat kernel bounds, *Probab. Th. Rel. Fields* **133** (2005), 245-266.
17. R. Pemantle, Randomization time for the overhand shuffle, *J. Theoret. Probab.* **2** (1989), 37-49.
18. E. Rodrigues, Convergence to stationarity for a Markov move-to-front scheme, *J. Appl. Probab.* **32** (1995), 768-776.
19. D. B. Wilson, Mixing time of the Rudvalis shuffle, *Electron. Commun. Probab.* **8** (2003), 77-85.
20. D. B. Wilson, Mixing times of lozenge tiling and card shuffling Markov chains, *Ann. Appl. Probab.* **14** (2004), 274-325.