

1 General Markov chains

Consider a finite state space Ω and a *transition kernel* $P : \Omega \times \Omega \rightarrow [0, 1]$ such that for every $x \in \Omega$, $\sum_{y \in \Omega} P(x, y) = 1$. The *Markov chain* corresponding to the kernel P is the sequence of random variables $\{X_0, X_1, X_2, \dots\}$ such that for every $t \geq 0$, we have $\mathbb{P}[X_{t+1} = y \mid X_t = x] = P(x, y)$. Note that we also have to specify a distribution for the initial state X_0 .

Corresponding to every such process, one can consider the (weighted) directed graph $D = (\Omega, A)$ with $A = \{(x, y) : P(x, y) > 0\}$ and edge weights $w(x, y) = P(x, y)$. Then the random process $\{X_t\}$ corresponds precisely to random walk on D : At every time step, one moves from the current vertex x to a neighbor y with probability $P(x, y)$.

Convergence to stationarity. For every $t \geq 0$, let $P^t(x, y) = \mathbb{P}[X_t = x \mid X_0 = y]$. The Markov chain described by P is said to be *irreducible* if for every $x, y \in \Omega$, there is some t such that $P^t(x, y) > 0$; in other words, there is always some way to reach any state from any other. This corresponds precisely to the digraph D being strongly connected. The chain is *aperiodic* if for every $x, y \in \Omega$,

$$\gcd(\{t : P^t(x, y) > 0\}) = 1.$$

Theorem 1.1 (Fundamental Theorem of Markov Chains). *If P is irreducible and aperiodic, then there is a unique probability measure $\pi : \Omega \rightarrow [0, 1]$ such that for every $x, y \in \Omega$, we have*

$$P^t(x, y) \rightarrow \pi(y) \quad \text{as } t \rightarrow \infty.$$

In other words, the Markov chain “forgets” where it started and converges to a unique limiting distribution. This is referred to as the *stationary measure* π .

Reversibility. A Markov chain is said to be *reversible with respect to the measure μ* if for every $x, y \in \Omega$, we have $\mu(x)P(x, y) = \mu(y)P(y, x)$. (These are called the “detailed balance conditions.”) The chain is said to be *reversible* if it is reversible with respect to some probability measure. Note that reversible chains correspond precisely to random walks on (weighted) *undirected graphs*.

Also, if P is irreducible and aperiodic—and hence has a unique stationary measure π by [Theorem 1.1](#)—then actually $\pi = \mu$. To see this, note that by the detailed balance conditions: For every $y \in \Omega$, we have

$$\sum_{x \in \Omega} \mu(x)P(x, y) = \sum_{x \in \Omega} \mu(y)P(y, x) = \mu(y) \sum_{x \in \Omega} P(y, x) = \mu(y). \quad (1.1)$$

The right-hand side can be interpreted as the probability of going to y in one step started from the measure μ . Now [Theorem 1.1](#) implies that if we start from distribution μ , then we converge to π ; on the other hand, (1.1) says that if we start distributed according to μ , then we stay that way under the chain. Thus $\mu = \pi$. This provides a nice local way to check that some measure is the stationary measure of the chain.

Remark 1.2. If P is irreducible, but not necessarily aperiodic, then there is still a unique stationary distribution, i.e. a probability π such that for every $x \in \Omega$, $\sum_{y \in \Omega} P(x, y)\pi(y) = \pi(x)$. But it may not be the case that the chain converges to π from some starting states.

For instance, if the chain is given by a directed graph with two nodes $\Omega = \{x, y\}$ and arcs (x, y) and (y, x) , then $\pi = (1/2, 1/2)$ is the unique stationary measure, but the chain does not converge to π when starting in either state x or y (because of periodicity).

For our purposes, aperiodicity is a rather weak obstruction to mixing. Given any chain P and number $\alpha \in (0, 1)$, we can consider the chain $P' = \alpha I + (1 - \alpha)P$. If P is irreducible, then so is P' . Moreover, for any such α , the chain P' is aperiodic (even if P was not). When measuring convergence to equilibrium, this α “self loop” probability does not slow down the chain too much.

1.1 Eigenvalues and mixing

Let us prove [Theorem 1.1](#) in the reversible case. To do this, we will think of P as an $\Omega \times \Omega$ matrix. If we also think about a probability measure $\mu \in \mathbb{R}^\Omega$ as a column vector, then $P\mu$ denotes the distribution that arises by starting at μ and taking one step of the chain associated to P .

If P is reversible with respect to π , then [\(1.1\)](#) implies that $P\pi = \pi$, i.e. π is an eigenvector with eigenvalue 1. To prove that $P^t(x, y) \rightarrow \pi(y)$ for every $x, y \in \Omega$, we will show that for every $x \in \Omega$,

$$\|\pi - P^t e_x\|_2 \rightarrow 0$$

where e_x is the vector with a 1 in the entry corresponding to x and zeros elsewhere.

Real eigenvalues. Note that P is not necessarily a symmetric matrix, but we can prove that P is similar to a symmetric matrix. Let D denote the diagonal matrix with $D_{xx} = \pi(x)$. Then

$$(\sqrt{D^{-1}}P\sqrt{D})_{xy} = \langle e_y, \sqrt{D^{-1}}P\sqrt{D}e_x \rangle = \langle \sqrt{D^{-1}}e_y, P\sqrt{D}e_x \rangle = \sqrt{\frac{\pi(x)}{\pi(y)}}P(x, y).$$

But by [\(1.1\)](#), this is equal to $\sqrt{\frac{\pi(y)}{\pi(x)}}P(y, x)$. Thus $\sqrt{D^{-1}}P\sqrt{D}$ is a real, symmetric matrix and hence has real eigenvalues. This implies that P also has real eigenvalues.

All eigenvalues in $[-1, 1]$. Now note that for any $v \in \mathbb{R}^\Omega$, we have

$$\|Pv\|_1 \leq \|P|v|\|_1 = \||v|\|_1, \tag{1.2}$$

where $|v|$ denotes the vector whose entries are the absolute value of the corresponding entries in v . This is simply because P is an averaging operator.

Now suppose that $Pv = \lambda v$. Then using [\(1.2\)](#)

$$|\lambda| \cdot \||v|\|_1 = \|Pv\|_1 \leq \||v|\|_1,$$

implying that $|\lambda| \leq 1$.

Unique eigenvector with eigenvalue 1. Suppose now that $Pv = v$ and consider the corresponding Laplacian matrix $L = D - PD$ (using our notation for “edge Laplacians,” this is $\frac{1}{2} \sum_{x,y} \pi(x)P(x, y)L_{\{x,y\}}$). One can check that this matrix is symmetric since PD is symmetric by the detailed balance conditions. As we saw in Lectures 14-15, for any vector w we have

$$w^T L w = \frac{1}{2} \sum_{x,y} \pi(x)P(x, y)(w_x - w_y)^2.$$

(The factor 1/2 is due to the fact that we are summing over all pairs x, y vs. all edges $\{x, y\}$.) Let $w = D^{-1}v$. Thus $Pv = v \implies Lw = 0 \implies w^T L w = 0$. Thus $w_x = w_y$ whenever $P(x, y) > 0$. But since the chain P is irreducible, we can connect every pair x, y by a chain of such implications, implying that $w = \alpha(1, 1, \dots, 1)$ is a multiple of the all-ones vector. But this implies that $v = Dw$ is a multiple of π . Since P is an averaging operator, it preserves the ℓ_1 norm, hence $\alpha = 1$ and $v = \pi$.

Not a bipartite graph. Now we claim that if P is aperiodic, -1 cannot be an eigenvalue of P . Suppose, for the sake of contradiction, that $Pv = -v$ for some $v \neq 0$. Again, let $|v|$ denote the vector whose entries are the absolute values of the corresponding entries in v . Then

$$\|v\|_2^2 = \|v\|^T Pv \leq |v|^T P|v| \leq \|v\|_2^2,$$

where the last inequality follows from the fact that all the eigenvalues of P lie in $[-1, 1]$. We conclude that $P|v| = |v|$, implying that $|v| = \pi$.

Finally, observe that $Pv = -v$ implies that for every x , one has $v_x = -(Pv)_x$, hence

$$\pi(x)\text{sgn}(v_x) = v_x = -(Pv)_x = -\sum_y P(y, x)v_y = -\sum_y P(y, x)\pi(y)\text{sgn}(v_y).$$

But by the detailed balance conditions, we have $\pi(x) = \sum_y P(y, x)\pi(y)$. Hence it must be that $\text{sgn}(v_x) = -\text{sgn}(v_y)$ whenever $P(y, x) > 0$.

Thus if we set $L = \{x : v_x < 0\}$ and $R = \{x : v_x > 0\}$, then $P(x, y) > 0$ implies x and y are on different sides of the bipartition. (Note that this is a bipartition since $|v| = \pi$ implies that $v_x \neq 0$ for any $x \in \Omega$.) But this implies that for x, y on the same side of the bipartition, we have $P^t(x, y) = 0$ when t is odd, contradicting the fact that P was assumed aperiodic.

Convergence to stationarity. Finally, consider any vector $w \in \mathbb{R}^\Omega$. Let $\lambda_1 = 1, \lambda_2, \dots, \lambda_n$ denote the eigenvalues of P arranged so that $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$, and denote the associated eigenvectors $\pi = v^{(1)}, v^{(2)}, \dots, v^{(n)}$. Recall that from the above reasoning, we have $|\lambda_i| < 1$ for $i > 1$. Write $w = \sum_{i=1}^n \alpha_i v^{(i)}$, and note that for any $t \geq 0$,

$$P^t w = \alpha_1 \pi + \sum_{i>2} \lambda_i^t \alpha_i v^{(i)}.$$

In particular, we have

$$\|P^t w - \alpha_1 \pi\|_2^2 = \sum_{i>2} \lambda_i^{2t} |\alpha_i|^2 \leq \lambda_2^{2t} \|w\|_2^2. \quad (1.3)$$

Since $|\lambda_2| < 1$, this implies that $\|P^t w - \alpha_1 \pi\|_2 \rightarrow 0$ as $t \rightarrow \infty$, showing that $P^t w \rightarrow \alpha_1 \pi$.

Note that if w has all non-negative entries, then since P is an averaging operator, we have $\|P^t w\|_1 = \|w\|_1$, hence $P^t w \rightarrow \|w\|_1 \pi$. Finally, observe that if $w = e_x$ then this implies $P^t e_x \rightarrow \pi$, which is exactly the claim of [Theorem 1.1](#).

In the next lecture, we will investigate how fast we converge to the stationary measure π in terms of the “inverse spectral gap” $1/(1 - |\lambda_2|)$.

1.2 Mixing times

Now we have seen that any irreducible, aperiodic Markov chain P on a finite state space Ω converges to a unique stationary measure π . We are not only concerned with convergence, but also the rate of convergence—we would like to be able to sample efficiently from π .

To this end, we first introduce a metric on the space of probability measures on Ω : For any two measures μ and ν on Ω , the *total variation distance* is defined by

$$d_{TV}(\mu, \nu) \stackrel{\text{def}}{=} \frac{1}{2} \|\mu - \nu\|_1 = \frac{1}{2} \sum_{x \in \Omega} |\mu(x) - \nu(x)|.$$

As an exercise, one can also show that $d_{TV}(\mu, \nu) = \max_{A \subseteq \Omega} |\mu(A) - \nu(A)|$.

For simplicity of notation, let us define $p_t^{(x)}$ to be the distribution given by $P^t e_x$ (i.e., the distribution of the chain started at x after t steps). For any $t \geq 0$ and $x \in \Omega$, define the quantity $\Delta_x(t) = d_{TV}(\pi, p_t^{(x)})$, and we set $\Delta(t) = \max_{x \in \Omega} \Delta_x(t)$. For $\varepsilon > 0$, we denote

$$\tau(\varepsilon) = \min\{t : \Delta(t) \leq \varepsilon\}.$$

In words, this is the first time t such that, starting from any initial state, the measure of the chain after t steps is within ε of the stationary measure. Finally, by convention, one takes $\tau_{\text{mix}} = \tau(1/2e)$ as the *mixing time* of the Markov chain P . Note that the precise value of ε is not so important; as the following lemma shows, once we have obtained the mixing time, further convergence to stationarity happens very fast. (We will skip the elementary proof.)

Lemma 1.3. *For every $t \geq 0$, we have*

$$\Delta(t) \leq \exp\left(-\left\lfloor \frac{t}{\tau_{\text{mix}}} \right\rfloor\right).$$

In particular, for every $\varepsilon > 0$, it holds that $\tau(\varepsilon) \leq \tau_{\text{mix}} \lceil \ln(1/\varepsilon) \rceil$.

Finally, we can use our proof of [Theorem 1.1](#) in the reversible case to give an upper bound on τ_{mix} in terms of the spectral gap of the chain.

Theorem 1.4. *Let P be a reversible and irreducible, aperiodic Markov chain on the state space Ω . Suppose that P has eigenvalues $1 = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$, and let $\lambda(P) = \max\{|\lambda_2|, |\lambda_n|\}$. Then*

$$\tau_{\text{mix}} \leq \left\lceil \frac{1 + \ln |\Omega|}{1 - \lambda(P)} \right\rceil.$$

Proof. Consider $x \in \Omega$ and $\varepsilon > 0$. Recalling [\(1.3\)](#) and using Cauchy-Schwarz, we have

$$d_{TV}(p_t^{(x)}, \pi)^2 = \frac{1}{4} \|p_t^{(x)} - \pi\|_1^2 \leq \frac{|\Omega|}{4} \|p_t^{(x)} - \pi\|_2^2 \leq \frac{|\Omega|}{4} \lambda(P)^{2t}.$$

Now setting $t = \lceil \frac{1}{1-\lambda(P)} \ln(|\Omega|/\varepsilon) \rceil$ and using the fact that $(1-\delta)^{1/\delta} \leq e^{-1}$ for $\delta > 0$ yields

$$d_{TV}(p_t^{(x)}, \pi)^2 \leq \frac{\varepsilon^2}{4},$$

implying $d_{TV}(p^t(x), \pi) \leq \varepsilon/2$. Setting $\varepsilon = 1/e$ and recalling the definition of τ_{mix} yields the desired result. \square

Finally, one should note that this bound is essentially tight up to the $O(\log |\Omega|)$ factor.

Theorem 1.5. *Under the assumption of [Theorem 1.4](#), we have*

$$\tau_{\text{mix}} \geq \frac{1}{1 - \lambda(P)} - 1.$$

Proof. Let v be an eigenvector of P with eigenvalue $\lambda = \lambda(P) \neq 1$. In that case, since π is also an eigenvector of P , we see that v is orthogonal to the stationary measure π , i.e. $\sum_{y \in \Omega} \pi(y) v_y = 0$. It follows that for $t \geq 0$ and any $x \in \Omega$,

$$|\lambda^t v_x| = |(P^t v)_x| = \left| \sum_y P^t(x, y) v_y - \pi(y) v_y \right| \leq \|v\|_\infty \sum_y |P^t(x, y) - \pi(y)| = 2\|v\|_\infty d_{TV}(p_t^{(x)}, \pi).$$

Now choose x so that $|v_x| = \|v\|_\infty$, yielding

$$d_{TV}(p_t^{(x)}, \pi) \geq \frac{1}{2} \lambda(P)^t.$$

Therefore $\lambda(P)^{\tau_{\text{mix}}} \leq 1/e$, implying that

$$\tau_{\text{mix}} \geq \frac{-1}{\log(1 - (1 - \lambda(P)))} \geq \frac{1}{1 - \lambda(P)} - 1,$$

where in the final line we have used that $\log(1 - a) \geq 1 + \frac{1}{a-1}$ for all $a \in [0, 1]$. □

So we see that up to a $\log |\Omega|$ factor, the spectral gap $1 - \lambda(P)$ controls the mixing time of the chain: If we set $\tau_{\text{rel}} = \frac{1}{1 - \lambda(P)}$ (commonly called the “relaxation time” of the chain), then

$$\tau_{\text{rel}} - 1 \leq \tau_{\text{mix}} \leq O(\log |\Omega|) \tau_{\text{rel}}.$$

1.3 Some Markov chains

One famous state space is the set of all permutations of n objects (for $n = 52$). In this case, $|\Omega| = n!$. Here are some shuffles:

1. **Random transposition.** At every step, we choose two uniformly random positions i and j (with replacement) and swap the cards at positions i and j .
2. **Top to random.** We take the top card and insert it at one of the n positions in the deck uniformly at random.
3. **Riffle shuffle.** We split the deck into two parts L and R uniformly at random, and then take a uniformly random interleaving of L and R .

And here’s a combinatorial example: Let $G = (V, E)$ be a graph with degree at most Δ , and suppose we have q colors with $q \geq \Delta + 1$ (so we are assured that G is q -colorable). Let Ω be the set of all q -colorings on G . Here is a natural Markov chain: Suppose we have a proper coloring $\chi : V \rightarrow [q]$. We choose a uniformly random $v \in V$ and a uniformly random color $c \in [q]$. If no neighbor of v in χ has color c , then we color v with c . Otherwise, we stay at the current coloring.

This example demonstrates the complex structure of Markov chains on combinatorial state spaces. For what values of q (depending on Δ) is the chain irreducible? It turns out that if $q \geq \Delta + 2$, then the chain is always irreducible, and the stationary measure is uniform on proper q -colorings. A huge open problem in MCMC (Markov chain Monte Carlo) is to resolve the following conjecture.

Conjecture 1.6. *For all $q \geq \Delta + 2$, this Markov chain has mixing time $O(n \log n)$, where $n = |V|$.*

The best bound (due to Vigoda, 1999) is that this holds for $q \geq \frac{11}{6} \Delta$. [Strictly speaking, this analysis holds for a slightly more complicated Markov chain with block dynamics—changing more than one color at a time. But I believe that by comparing the two chains, one can show rapid mixing for the chain above as well.]