

Math-Stat-491-Fall2014-Notes-III

Hariharan Narayanan

October 28, 2014

1 Introduction

We will be closely following the book

”Essentials of Stochastic Processes”, 2nd Edition, by Richard Durrett,

for the topic ‘Finite Discrete time Markov Chains’ (FDTM). This note is for giving a sketch of the important proofs. The proofs have a value beyond what is proved - they are an introduction to standard probabilistic techniques.

2 Markov Chain summary

The important ideas related to a Markov chain can be understood by just studying its graph, which has nodes corresponding to states and edges corresponding to nonzero entries in the transition matrix.

Figure 1 helps us to summarize key ideas.

The first part of this figure shows an irreducible Markov chain on states A, B, C . The graph in this case is strongly connected, i.e., one can move from any node to any other through directed paths. Such a Markov chain has a unique stationary distribution.

This Markov chain is also ‘aperiodic’. If you start from any node you can return to it in $2, 3, 4, 5, \dots$ steps. So the GCD of all these loop lengths is 1. For such Markov chains if you take a sufficiently large power P^n of the transition matrix P it will have all entries positive. (In this case however P itself has this property.) If you start from any probability distribution π' and run an irreducible aperiodic Markov chain for ‘infinite time’ $\pi'^T P^n$ will converge to the unique stationary distribution. The value of this distribution will be positive for each state.

Next consider the second Markov chain on A', B', C', D' . Here we can see that from D' we can reach A', B', C' , but not the other way about. Further if you restrict the Markov chain to A', B', C' you will get an irreducible chain. The Markov chain on A', B', C', D' is not irreducible but has a unique stationary distribution. However it takes zero value on some states. The general rule is the following. If from a given state X you can reach some other state Y but cannot return from Y to X , then the stationary distribution will take value zero on X . We call such states ‘transient’. If you start from any probability distribution π' and run this Markov chain indefinitely, $\pi'^T P^n$ will converge to the unique stationary distribution. The value of this distribution will be positive for each state in R_1 but zero for D, D' .

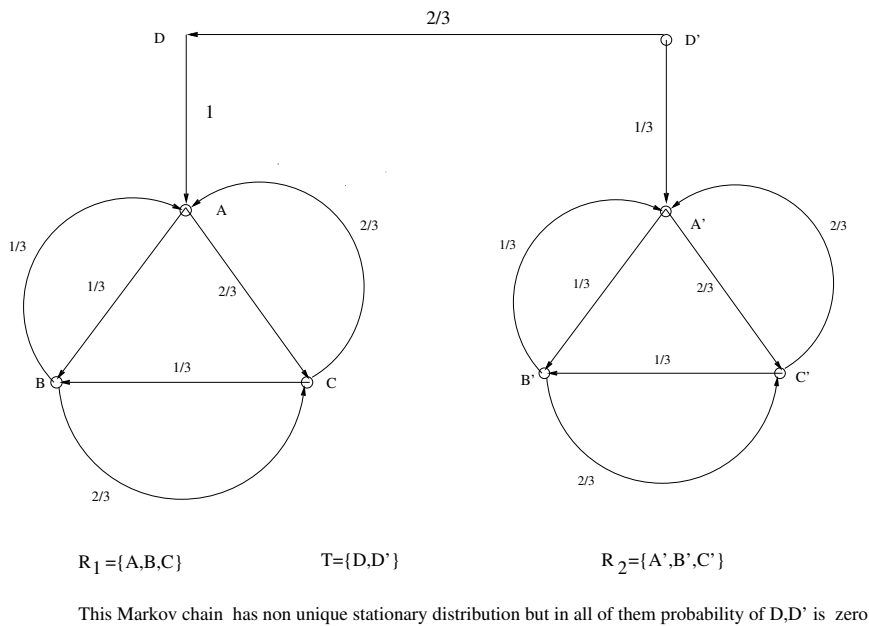
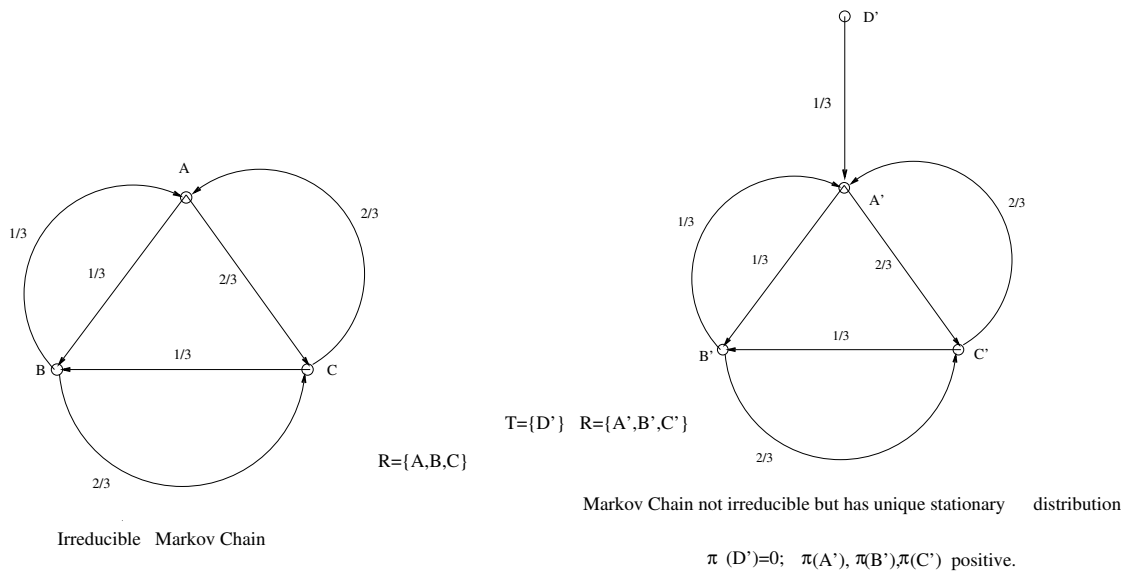


Figure 1: Markov Chain summary

In general just by examining the graph we can partition the set of states into T which are transient, and a number of R_i whose states are 'recurrent', i.e., not transient. The R_i s each have the property : there are *no directed paths* leaving the set and there are *directed paths from any node to any other in the set*.

The third Markov chain has states partitioned into $T \equiv \{D, D'\}$, $R_1 \equiv \{A, B, C\}$, $R_2 \equiv \{A', B', C'\}$. Note that from D, D' we can go to states in R_1, R_2 , but not return. If we restrict the Markov chain to either R_1 or R_2 , we will get an irreducible Markov chain. But on the set of states $R_1 \cup R_2$ the Markov chain

is not irreducible since you cannot go from a state in R_1 or R_2 to the other. In this case we can get two primitive stationary distributions one which is nonzero only on R_1 and zero on all the others and a second which is nonzero only on R_2 but zero on all the others. We can show that all stationary distributions are convex combinations of these two primitive distributions, i.e., $\pi^T = \lambda\pi^{1T} + (1 - \lambda)\pi^{2T}, 0 \leq \lambda \leq 1$.

The first important result

Theorem 2.1. *A finite, irreducible Markov chain X_n has a unique stationary distribution $\pi(\cdot)$.*

Remark: It is not claimed that this stationary distribution is also ‘steady state’, i.e., if you start from any probability distribution π' and run this Markov chain indefinitely, $\pi'^T P^n$ may not converge to the unique stationary distribution. That happens only if the irreducible Markov chain is aperiodic, i.e., the GCD of the length of loops starting from any node is 1. Figure 2 shows a Markov chain of period 2. This has a unique stationary distribution but $\pi'^T P^n$ does not converge to it.

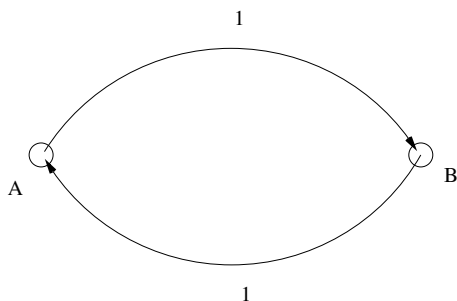


Figure 2: Markov chain with period 2

Recall that an irreducible Markov chain is one whose directed graph has a directed path from node i to node j for every possible node pair.

Now every stationary distribution is a nonnegative row eigenvector of the transition matrix corresponding to the eigenvalue 1.

Thus the statement of the theorem involves only ideas from linear algebra and graphs. We will first sketch an algebraic proof of this result and then a probabilistic proof.

3 Algebraic Proof of Theorem 2.1

Case 1. The transition matrix has no zero entries.

We know that if $\pi(\cdot)$ is a stationary distribution, when we write it as a row vector π^T , it satisfies $\pi^T P = \pi^T$, i.e., π^T is a row eigenvector for the eigen value 1.

We claim that every row eigenvector corresponding to eigenvalue 1, has only entries which are all of the same sign.

Suppose not. Let some entries be positive and some negative. Consider the equation $\pi^T P = \pi^T$. The j^{th} entry of the RHS of this equation is obtained by taking the ‘dot product’ of π^T with the j^{th} column of the matrix P . This dot product $\pi(j)$ is the sum of some positive and some negative terms since the j^{th} column of P is fully positive while π^T has both positive and negative entries. Formally,

$\sum_i \pi(i) P_{ij} = \pi(j)$. Let π^T be the nonnegative vector obtained by changing the sign of the negative entries of π^T , i.e., $\pi'(i) = |\pi(i)| \forall i$. Clearly, $|\sum_i \pi(i) P_{ij}| < \sum_i |\pi(i)| P_{ij}$, since the right side ‘dot product’ has no cancellation while the left side one has. Thus, $\pi'(j) = |\pi(j)| = |\sum_i \pi(i) P_{ij}| < \sum_i |\pi(i)| P_{ij} = \sum_i \pi'(i) P_{ij}$.

We now have $\sum_j (\sum_i \pi'(i) P_{ij}) > \sum_j |\sum_i \pi(i) P_{ij}| = \sum_j |\pi(j)|$. We will show however that the left side actually sums up to the same expression as the right side, i.e., to $\sum_j |\pi(j)|$. We have, $\sum_j (\sum_i \pi'(i) P_{ij}) = \sum_i (\sum_j \pi'(i) P_{ij}) = \sum_i \pi'(i) (\sum_j P_{ij}) = \sum_i \pi'(i) (\sum_j 1) = \sum_i \pi'(i) = \sum_i |\pi(i)|$.

We conclude therefore that every row eigen vector π^T of P corresponding to eigen value 1 has entries of the same sign. Wlog let us take this to be the positive sign.

Next, since P has only positive entries, $\pi^T P$ has only positive entries. But $\pi^T P = \pi^T$, so that π^T has only positive entries.

Thus without loss of generality, we may take a row eigen vector of P corresponding to eigen value 1 to have only positive entries.

We will now show that we cannot have two independent row eigen vectors of P corresponding to eigen value 1.

Suppose x^T, y^T are two such independent vectors. Then there exists a linear combination z^T , which has a zero entry but is not fully zero. But z^T is a row eigen vector of P corresponding to eigen value 1 and therefore, by the above proof, must be fully nonzero, a contradiction.

We conclude that P has a unique row eigen vector corresponding to eigen value 1 which is non negative and whose row sum is 1 and this vector is fully positive. But such a row vector is a stationary distribution of P .

Therefore P has a unique stationary distribution π and $\pi(y) > 0 \forall y$.

Case 2. The transition matrix P corresponds to a general irreducible Markov chain and has zero entries.

Observe that if $\pi^T P = \pi^T$, then $\pi^T (I + P)/2 = \pi^T$. But if P is a transition matrix (i.e., rows of P add up to 1), so is $(I + P)/2$, since its rows also add up to 1. The Markov chain corresponding to $(I + P)/2$ looks like that corresponding to P except that there are additional self loops of probability 1/2 and all the old edges have half the original probability. Therefore in the new graph too there are directed paths from any node to any other node. We conclude that the Markov chain corresponding to $(I + P)/2$ is also irreducible.

Let $Q \equiv (I + P)/2$. Observe that, because of the self loop at each node (state), $Q(i, i) \neq 0 \forall i$ and $(Q)^n(i, i) \neq 0 \forall i, n$. Using Chapman-Kolmogorov theorem we conclude from the irreducibility of the Markov chain corresponding to Q , that for each i, j there exists some n such that $Q^n(i, j) > 0$, and further that for each positive integer $Q^{n+k}(i, j) \geq (Q)^k(i, i) \times Q^n(i, j) > 0$. Thus for a large enough positive integer m we must have Q^m fully positive. Observe that if Q has all row sums as 1 so will Q^m have. It is clear that if $\pi^T P = \pi^T$, we also have $\pi^T Q = \pi^T$ and $\pi^T Q^m = \pi^T$. Now Q^m is a transition matrix of a Markov chain that is fully positive and any row eigen vector of P corresponding to eigen value 1 is also a row eigen vector of Q^m corresponding to eigen value 1. In particular any stationary distribution of P is also a stationary distribution of Q^m .

By discussion of case 1 above we know that Q^m has a unique stationary distribution π and $\pi(y) > 0 \forall y$.

It follows that P has a unique stationary distribution π and $\pi(y) > 0 \forall y$.

3.1 An additional result

Theorem 3.1. *The highest magnitude eigen value for a transition matrix P corresponding to an irreducible Markov chain is 1 and is unique.*

We use the following lemma which can be proved by plane geometry

Lemma

Let c_1, c_2, \dots, c_k be complex numbers. Let c_j have the maximum magnitude among all these numbers. Let $\sum_i p_i c_i = d$ be a convex combination of the c_i , with p_j , the coefficient for c_j , greater than zero. Then $|d| \leq |c_j|$ and further, $|d| = |c_j|$ iff whenever $p_i > 0$, we have $c_i = c_j$.

Proof of Theorem 3.1: Let λ be a maximum magnitude eigen value of P and let x be a right eigen vector corresponding to this eigen value. We have $Px = \lambda x$. Let x_j have the maximum magnitude among entries of x . We then have $\sum_k p_{jk} x_k = \lambda x_j$. Noting that $\sum_k p_{jk} x_k$ is a convex combination of the entries of x , and using the lemma we conclude that $|\lambda| \leq 1$, and further $|\lambda| = 1$ can happen only if $\sum_k p_{jk} x_k = x_j$, and therefore $\lambda = 1$.

4 Probabilistic proof of Theorem 2.1

In order to describe the probabilistic ideas technically we introduce some notation.

4.1 Notation

1. $T_y \equiv \min\{n \geq 1 : X_n = y\} \equiv$ time of first *return* to y or, equivalently, the time of first *visit* to y after $n = 0$. Here even if we don't start from $X_0 = y$, we still use the term 'return'. Note that, when the starting state is defined, this is a random variable, since the event $T_y = k$ has a probability associated with it.

2. $T_y^k \equiv \min\{n > T_y^{k-1} : X_n = y\} \equiv$ time of k^{th} *return* to y . Note that $T_y = T_y^1$.

The probability associated with $T_y = n$ starting from x is denoted by $Pr_x\{T_y = n\}$. The expected value of T_y starting from x is denoted by $E_x(T_y)$ and given by $\sum_n n \times Pr_x\{T_y = n\}$.

Remark: Observe that if the expected value of T_z (starting from z) is greater than the expected value of T_y (starting from y), you return to z less often than you return to y so we expect $\pi(z) < \pi(y)$. We will show that $\pi_y = 1/E_y(T_y)$.

3. ρ_{xy} denotes $Pr_x\{T_y < \infty\}$, i.e., the probability of reaching from x to y in finite time, given that the starting state is x . The event of 'reaching from x to y in finite time', is the event of all finite sequences $X_0 = x, X_1 = x_1, \dots, X_{k-1} = x_{k-1}, X_k = y$, where none of the X_i except X_k is equal to y . The probability associated with such a sequence is $p(x, x_1) \times p(x_1, x_2) \cdots \times p(x_{k-1}, x_k)$. The probability associated with the event is the sum of probabilities of all such sequences (which represent disjoint events).

4. ρ_{yy} denotes $Pr_y\{T_y < \infty\}$. Equivalently, ρ_{yy} is the probability of returning to y starting from y using a finite number of steps, i.e.,

$$\rho_{yy} = \sum_n Pr\{X_n = y \mid X_0 = y, X_j \neq y, 0 < j < n\}.$$

5. $\rho_{yy}^{(k)}$ denotes $Pr_y\{T_y^k < \infty\}$.
6. $N_y \equiv$ number of visits to y . This again is a random variable once the starting state is specified.
7. We say a state is *recurrent* if $\rho_{yy} = 1$. It is *transient* if $\rho_{yy} < 1$.

4.2 Preliminary basic results

1. (Strong Markov Property) If T is a stopping time then

$$Pr\{X_{T+1} = z \mid X_T = x, T = n\} = Pr\{X_1 = z \mid X_0 = x\}.$$

Here T is a random variable with the property that we can determine the truth of $T = n$, if we know the sequence $X_0 = x_0, X_1 = x_1, \dots, X_{k-1} = x_{k-1}, X_n = x$. Informally, the truth of $T = n$ depends only on the past and the immediate present and not on the future. The Markov chains we deal with are time homogeneous, so we could choose any (fixed) time n as our zero time. The strong Markov property allows us to use the random variable T as our zero time. This essentially means that we can take a time n as our starting time under the condition that certain things have taken place by that time.

A consequence: We have

$$\rho_{yy}^{(k)} = Pr_y\{T_y + n < \infty \mid T_y^{k-1} = n\} \times \rho_{yy}^{(k-1)} = \rho_{yy} \times \rho_{yy}^{(k-1)}.$$

By induction it follows that $\rho_{yy}^{(k)} = \rho_{yy}^k$.

[Here are the details to prove $\rho_{yy}^{(2)} = \rho_{yy}^2$. We have

$$\rho_{yy}^{(2)} \equiv Pr_y\{T_y^2 < \infty\}$$

\equiv Probability that $X_0 = X_m = X_n = y$ with $X_i \neq y, 0 < i < m$ and $X_i \neq y, m < i < n$, for all values of m, n with $m \leq n$.

Let $X_0 = y$. Let T'_y be the first m when $X_m = y$, for $m \geq 1$, and let T''_y be the first k when $X_{k+T'_y} = y$, for $k \geq 1$.

Observe that T'_y, T''_y are random variables and that $T_y^2 = T'_y + T''_y$.

So $Pr_y\{T_y^2 < \infty\} = \sum_k \sum_m Pr_y\{T'_y = m, T''_y = k\} = \sum_k \sum_m Pr\{T''_y = k \mid T'_y = m\} \times Pr_y\{T'_y = m\}$.

Now by strong Markov property $Pr\{T''_y = k \mid T'_y = m\} = Pr_y\{T'_y = k\}$. Thus we have

$$\begin{aligned} \rho_{yy}^{(2)} &= Pr_y\{T_y^2 < \infty\} = \sum_k \sum_m Pr_y\{T'_y = k\} \times Pr_y\{T'_y = m\} \\ &= \sum_k Pr_y\{T'_y = k\} \times \sum_m Pr_y\{T'_y = m\} = \rho_{yy} \times \rho_{yy}. \end{aligned}$$

Another consequence:

$$Pr_x\{T_y^k < \infty\} = Pr\{T_y^{k-1} < \infty \mid X_T = y\} \times Pr_x\{T_y < \infty\} = \rho_{yy}^{k-1} \rho_{xy}.$$

2. (Expectation of infinite sums of random variables) Let $X = \sum_{i=0}^{\infty} X_i$, where X_i are nonnegative random variables. Then

- if $\sum_{i=0}^n E(X_i)$ converges as $n \rightarrow \infty$, then $E(X) = \sum_{i=0}^{\infty} E(X_i)$,
- if $\sum_{i=0}^n E(X_i)$ becomes (positively) unbounded as $n \rightarrow \infty$, then $E(X) = \infty$.

Observe that this result has to be proved separately and is not immediate from 'finite' linearity of expectation of random variables.

Proof: Let $Y_k \equiv \sum_0^k X_i$. Then, by linearity of expectation, $E(Y_k) = \sum_0^k E(X_i)$. Now suppose $\sum_0^k E(X_i)$ converges as $k \rightarrow \infty$. It follows that $E(Y_k)$ converges too. Further since X_i are nonnegative it follows that $Y_k \leq Y_{k+1}, \forall k$, and $Y_k \leq X$. Now Lebesgue's monotone convergence theorem states that if $Y_k \leq Y_{k+1} \forall k$ and $\lim_{k \rightarrow \infty} Y_k = X$, then $\lim_{k \rightarrow \infty} E(Y_k) = E(X)$, provided $\lim_{k \rightarrow \infty} E(Y_k)$ exists. The conditions of this theorem are clearly satisfied in the present case of Y_k and X . This completes the proof of the first part of the statement.

Next suppose $E(Y_k) = \sum_0^k E(X_i)$ is (positively) unbounded as $k \rightarrow \infty$. Since $Y_k \leq X$, we must have $E(Y_k) \leq E(X)$ making $E(X)$ also (positively) unbounded.

4.3 Intermediate results

1. If X is a discrete random variable taking values $1, 2, \dots$, then $X = \sum_{n=1}^{\infty} 1_{X \geq n}$, and therefore

$$E(X) = \sum_{n=1}^{\infty} Pr\{X \geq n\}.$$

2. (Recall that the random variable $N_y \equiv$ number of visits to y , and that $Pr_x\{T_y^k < \infty\} = Pr_x\{N(y) \geq k\}$.)

$$Pr_x\{N(y) \geq k\} = \rho_{yy}^{k-1} \rho_{xy}.$$

3. •

$$E_x(N_y) = \rho_{xy} (\sum_{k=0}^{\infty} \rho_{yy}^k) = \rho_{xy} / (1 - \rho_{yy}),$$

where $E_x(N_y)$ is the expectation of N_y , when the starting state is x .

- A state y is recurrent (i.e., $\rho_{yy} = 1$) iff $E_y(N_y) = \infty$.

4.

$$E_x(N_y) = \sum_{n=1}^{\infty} p^n(x, y),$$

where $p^n(x, y)$ refers to the probability of reaching from x to y in n steps.

Proof is by recognizing $N_y = \sum_{n=1}^{\infty} 1_{X_n=y}$, and taking expectation on both sides.

5. If y is recurrent and $\rho_{yx} > 0$, then x is recurrent.

Proof: Since $\rho_{yx} > 0$, and y is recurrent $\rho_{yy} \neq 0$. So we must have for some $j, k, p^j(x, y) > 0, p^k(y, x) > 0$. Now

$$E_x(N_x) = \sum_{n=1}^{\infty} p^n(x, x) \geq p^j(x, y) (\sum_{n=1}^{\infty} p^n(y, y)) p^k(y, x) = \infty.$$

4.4 Final Results

1. • A finite, irreducible Markov chain has atleast one recurrent state.

Proof is by noting that

$$\sum_y E_x(N_y) = \sum_y \sum_n p^n(x, y) = \sum_n (\sum_y p^n(x, y)) = \sum_n (1) = \infty.$$

Therefore atleast one of the y must be such that $E_x(N_y) = \infty$.

- A finite, irreducible Markov chain has all its states recurrent.

Proof: At least one of the states y is recurrent and for each x , $\rho(x, y), \rho(y, x)$ are nonzero. By a result proved earlier this is sufficient for x to be recurrent.

2. The states of any finite Markov chain can be partitioned into sets T, R_1, \dots, R_k , where the states in T are transient, the remaining states are recurrent and the restrictions of the original Markov chain to each of the R_i is irreducible.

Proof: The set T is composed of states x for which for some y , we have $\rho(x, y) > 0$, and $\rho(y, x) = 0$, i.e., we can reach y from x but not return. The remaining states form equivalence classes of states which can be reached from each other. These are the R_i . Once you enter an R_i you cannot leave. The Markov chain restricted to the nodes in any of the R_i is irreducible since one can reach any node in it from any other and therefore the states are recurrent.

Remark: One can build an infinity of stationary distributions if there are atleast two R_i . Let μ^1 be the stationary distribution where $\mu(x) = 0, x \in T, \mu(z) = 0, z \in R_2$, and on R_i, μ^1 agrees with the unique stationary distribution μ_1 of the restriction of the Markov chain on R_1 . Define μ^2 interchanging R_1, R_2 in the definition of μ^1 . Every convex combination of μ^1, μ^2 would yield a stationary distribution of the original Markov chain.

3. Every finite irreducible Markov chain has a unique stationary distribution.

The algebraic proof for this fact is simple and already given.

When the irreducible Markov chain is aperiodic, i.e., the gcd of the period of return starting from any node (state) is 1, we can show that starting with any x , the limit as $n \rightarrow \infty$, of $p^n(x, y)$ tends to $\pi(y)$, the value of the stationary distribution at y and therefore $\pi(\cdot)$ is unique.

An FDTM has a finite number of states on which it can rest at discrete times $1, 2, \dots$. Usually the state at time n is denoted by X_n , and the possible set of states could be denoted by lower case symbols $x, y \dots$ etc.

If the chain is at state i at time n , it moves to state j at time $n+1$, with probability $p(i, j)$. This could be stated alternatively as ‘the conditional probability of the chain being at state j at time $n+1$, given that it is at state i at time n is $p(i, j)$.’ The key idea is that the probability of the chain being at state j at time $n+1$, does not depend on what happened before time n . Formally, we have the ‘Markov property’,

$$Pr\{X_{n+1} = j \mid X_n = i\} = Pr\{X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0\}.$$

We have further assumed ‘temporal homogeneity’, i.e, that $p(i, j)$ does not depend upon n .

In particular this means

$$Pr\{X_{n+1} = j \mid X_n = i\} = Pr\{X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, \dots, X_k = i_k\}, \quad (1)$$

because this latter expression, when we shift the time origin to k , is the same as

$$Pr\{X_{n+1-k} = j \mid X_{n-k} = i\} = Pr\{X_{n+1-k} = j \mid X_{n-k} = i, X_{n-k-1} = i'_{n-k-1}, \dots, X_0 = i'_0\},$$

where $i'_r \equiv i_{r+k}$.

Just to be clear on how we should go about proving formal versions of informal statements consider the following. Suppose what is given is not the immediately previous state but say a state k steps before. Surely, what happened before time $n + 1 - k$ is unimportant? Let us prove this for $k = 2$. Formally, let us prove

$$Pr\{X_{n+1} = j \mid X_{n-1} = i_{n-1}\} = Pr\{X_{n+1} = j \mid X_{n-1} = i_{n-1}, \dots, X_0 = i_0\}.$$

We will use the law of total probability. The LHS

$$Pr\{X_{n+1} = j \mid X_{n-1} = i_{n-1}\} = \sum_{\text{all states } s} Pr\{X_{n+1} = j \mid X_n = s, X_{n-1} = i_{n-1}\} \times Pr\{X_n = s \mid X_{n-1} = i_{n-1}\}.$$

The RHS

$$Pr\{X_{n+1} = j \mid X_{n-1} = i_{n-1}, \dots, X_0 = i_0\}$$

$$= \sum_{\text{all states } s} Pr\{X_{n+1} = j \mid X_n = s, X_{n-1} = i_{n-1}, \dots, X_0 = i_0\} \times Pr\{X_n = s \mid X_{n-1} = i_{n-1}, \dots, X_0 = i_0\}.$$

Now we use the basic Markov property on each of the two terms being multiplied and using (1) write the above expression as

$$\sum_{\text{all states } s} Pr\{X_{n+1} = j \mid X_n = s, X_{n-1} = i_{n-1}\} \times Pr\{X_n = s \mid X_{n-1} = i_{n-1}\}.$$

Since this is the same as the LHS, the proof is complete.

Next let us prove the following variation of the Markov property

$$Pr\{X_{n+1} = j \mid X_n = i\} = Pr\{X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, \dots, X_2 = i_2, X_0 = i_0\}.$$

We have skipped $X_1 = i_1$ in the right side conditional probability.

Note that the RHS is equal to

$$\sum_{\text{all states } s} Pr\{X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, \dots, X_2 = i_2, X_1 = s, X_0 = i_0\} \\ \times Pr\{X_1 = s \mid X_n = i, \dots, X_2 = i_2, X_0 = i_0\}.$$

But this is the same as

$$\sum_{\text{all states } s} Pr\{X_{n+1} = j \mid X_n = i\} \times Pr\{X_1 = s \mid X_n = i, \dots, X_2 = i_2, X_0 = i_0\}$$

This is equal to

$$Pr\{X_{n+1} = j \mid X_n = i\} \times \sum_{\text{all states } s} Pr\{X_1 = s \mid X_n = i, \dots, X_2 = i_2, X_0 = i_0\}$$

But

$$\sum_{\text{all states } s} Pr\{X_1 = s \mid X_n = i, \dots, X_2 = i_2, X_0 = i_0\} = 1.$$

So the RHS=LHS.

The same trick can be used even when there are multiple gaps.

Temporal homogeneity implies we can start from $X_r = i_0$ and repeat the above arguments, replacing X_n by X_{n+r} , etc.

(Natural question at this stage: surely we could also use models where the chain state at $n+1$ depends on states at time $n, n-1, \dots, n-k$?

Answer: They are not needed! This model is sufficiently versatile.)

Note here that everything we know about an FDTM is captured by the set of $p(i, j)s$. So we could represent the FDTM

- pictorially through a graph with vertices as states with, for each ordered pair (i, j) , an edge directed from i to j with weight $p(i, j)$. Of course if $p(i, j) = 0$, we omit the edge from i to j .
- in terms of a ‘transition matrix’ with rows and columns named by the states with the $(i, j)^{th}$ entry being $p(i, j)$.

Since, from a given state i at time n , the chain *must* move to *some* state (including i) at time $n+1$, it follows that that the outgoing edges (including selfloops, if any) must have weights adding up to 1, and, in the transition matrix, the rows should sum up to 1.

STEP 1 in the study of FDTM: Convert word problems in the book to FDTMs specifying the transition matrix and drawing the FDTM graph.

Some fundamental notions associated with an FDTM are listed below.

1. Steady state distribution on the states: this is a probability distribution $\pi(\cdot)$, such that if we pick initial states i with probability $\pi(i)$, the next states j will occur with probability $\pi(j)$.

Note here that if we pick initial states according to some probability distribution $\pi'(\cdot)$, the probability that the next state is j is given by

$$\sum_i Pr(i) \times Pr(j \mid i) = \sum_i \pi'(i) \times p(i, j).$$

Let $(\pi')^T$ denote the row vector whose i^{th} entry is $\pi'(i)$, and let P denote the transition matrix of FDTM (with $p(i, j)$ as the (i, j) entry).

So if the initial probability distribution is $\pi'(\cdot)$, the next state distribution is $(\pi')^T P$.

In the case of the stationary distribution $\pi(\cdot)$, we have $(\pi)^T = (\pi)^T P$.

2. Transient and recurrent states: if you start from a ‘recurrent’ state, you will return to it, with probability 1. If this probability is < 1 , then the state is said to be ‘transient’. We can prove that the following ‘event’ has probability 1 : if we start from a transient state we will return to it only a finite number of times.
3. The expected time of return starting from a given recurrent state. This is self explanatory.
4. $\rho_{xy} \equiv Pr\{we\ can\ go\ from\ x\ to\ y\ in\ finite\ time\}$.
5. Absorbing state: A state is said to be absorbing if, when we start from it, we remain at it with probability 1, i.e., i is an absorbing state iff $p(i, i) = 1$.
6. Irreducible Markov Chain: The graph of the FDTM has the property that given any ordered pair of states (i, j) , there is a directed path from i to j in the graph of the FDTM.

STEP 2 in the study of FDTM: Internalize the above definitions.

Through the lectures we will attempt to answer the following questions rigorously:

1. How to identify recurrent and transient states by looking only at the graph?
 Answer: Check if we can go from the given state i to some state j through a directed path such that there is no return directed path. If this is true the state is transient, otherwise recurrent.
2. Does the FDTM have a stationary distribution? If it exists, is it unique?
 Answer: FDTMs always have a stationary distribution. It is unique if the FDTM is irreducible.
3. Starting from a given recurrent state i , what is the expected time of return to it?
 Answer: If the FDTM is recurrent and has the stationary distribution $\pi(\cdot)$, the expected time of return to the recurrent state i is $1/\pi(i)$.
4. Starting from a transient state what is the expected time to reach a specified absorbing state?
 Answer: We will describe a method of computing this quantity.
5. Suppose we start from any state and keep running the markov chain according to its transition matrix. Will we encounter the different states with frequency in proportion with their $\pi(\cdot)$ value?
 Answer: Yes, if the FDTM is irreducible.
6. Suppose we start the Markov chain according to a probability distribution $\pi^0(\cdot)$ and let it run forever. Let $\pi^n(\cdot)$ denote the probability distribution after n steps. Will $\pi^n(\cdot)$ converge to the stationary distribution in the limit as $n \rightarrow \infty$?
 Answer: Yes if the FDTM is irreducible and aperiodic. (Periodic means all states recur after a certain fixed period > 1 .)

To answer the above questions we will need to understand the following:

1. What are the characteristic features of a transition matrix (or equivalently the graph) of an FDTM?

Answer: The matrix should have nonnegative entries and the rows should add up to 1. The outgoing edges from any node in the graph should have the sum of their weights equal to 1.

Let us call these respectively Markov transition matrix and Markov chain graph.

2. What kind of a matrix is P^k ? What is its significance?

Answer: P^k is also a Markov transition matrix because its entries are nonnegative and the row sum is 1.

Its significance is the following:

Its $(i, j)^{th}$ entry gives the probability of reaching j at the time $n + k$ starting with i at time n , i.e., the probability of reaching j from i in k steps.

The above requires an understanding of the famous ‘Chapman-Kolmogorov equation’. Below, we give a sketch of the proof.

Let us denote by $p^s(i, j)$ the probability of reaching state j starting from state i in s steps. ‘Chapman-Kolmogorov equation’ states

$$p^{m+n}(i, j) = \sum_k p^m(i, k) \times p^n(k, j).$$

The proof is straight forward. To reach from i to j in $m + n$ steps we must reach some intermediate state in m steps. So the probability of reaching from i to j in $m + n$ steps through the intermediate state k is the product $p^m(i, k) \times p^n(k, j)$. Now the events ‘reaching k ’ in m steps and then reaching j in n steps are disjoint for different k s and together make up the event ‘reaching from i to j in $m + n$ steps. So

$$p^{m+n}(i, j) = \sum_k p^m(i, k) \times p^n(k, j).$$

But this is exactly how the matrix P^{m+n} is computed in terms of the matrices P^m, P^n .

(To formalize the sketch we must use conditional probabilities.

Example: $p^s(i, j) = Pr\{X_s = j \mid X_0 = i\}$.

Probability of reaching from i to j in $m + n$ steps passing through k after m steps

$$= Pr\{X_{m+n} = j, X_m = k \mid X_0 = i\}$$

$$= Pr\{X_{m+n} = j \mid X_m = k, X_0 = i\} \times Pr\{X_m = k \mid X_0 = i\}$$

$$= Pr\{X_{m+n} = j \mid X_m = k\} \times Pr\{X_m = k \mid X_0 = i\} = p^n(k, j) \times p^m(i, k).$$

3. A random variable T that takes values $0, 1, 2, \dots$ is called a *stopping time* or *Markov time* if whether $T = k$ or not depends only on the states X_0, \dots, X_k , that the markov chain takes at times $0, 1, \dots, k$. Are the following stopping times?
 - (a) $T = n$, if starting at x at time 0 we reach y at time n .
 - (b) $T = n$, if starting at x at time 0 we reach y at time n and $n \leq 100$.
 - (c) $T = n$, if we are at y for the last time.

Answer: The first two random variables are stopping times. The last is not, since, to know that we are at y for the last time is impossible knowing only the past history.

‘Stopping time’ is one of the most slippery and useful notions that we encounter while studying FDTMs. For us, its power lies in the fact that we can, because of the ‘strong Markov property,’ treat X_T as X_0 , even though T is a random variable.

STEP 3 in the study of FDTM: Know how to prove the above results rigorously. But do not lose the ‘common touch’- understand them in the commonsensical intuitive way too.

Special conditions make the computation of stationary probability easy for some Markov chains. Some of these are discussed below.

1. **Doubly stochastic chains** For these chains the transition matrix has its columns also adding up to 1. Now, if we sum all the row vectors, we get the vector $\mathbf{1}^T = (1, 1, \dots, 1)$, i.e., $\mathbf{1}^T P = \mathbf{1}^T$. So scaling $\mathbf{1}^T$ by the reciprocal of the number of states will satisfy the same equation while being a probability distribution. Thus, for such chains, the uniform probability distribution is a stationary distribution.

2. **Duality** Let Markov chain X_n on states S have transition probability $p(i, j)$ and (unique) stationary probability $\pi(i), i \in S$.

Define a new Markov chain Y_n on S with transition probability $p'(i, j) \equiv \pi(j) \times p(j, i) / \pi(i)$. Let us call this the *dual* Markov chain to X_n . It can be verified that $\sum_i p'(i, j) = 1$.

It can be seen that, whenever there is an edge from i to j in the original Markov chain with probability $p(i, j)$, in the dual there is an edge from j to i with probability $p'(j, i)$. Let $\pi'(\cdot)$ denote the stationary distribution of the dual Markov chain. We can show that $\pi'(\cdot) = \pi(\cdot)$.

3. **Reversibility** When the Markov chain is its own dual, we have the *detailed balance condition*

$$p(i, j) = \pi(j) \times p(j, i) / \pi(i).$$

We call such a Markov chain reversible.

The algorithm given below verifies whether a given chain is reversible even where the distribution $\pi(\cdot)$ is not available, and if reversible, computes $\pi(\cdot)$.

The algorithm is simple. Let us suppose that from any node in the graph we can go to any other node through a directed path (i.e., always going along the direction of the arrow of the edge). We can show then that all entries of $\pi(\cdot)$ are positive even if the Markov chain is not reversible.

Firstly if $p(i, j) > 0$ we must have $p(j, i) > 0$ for the detailed balance condition

$p(i, j) = \pi(j) \times p(j, i) / \pi(i)$, to hold. So we should check if every edge has, in parallel, an edge going in the opposite direction. If this is not true, we can declare the chain to be not reversible.

Let us start from some node say 0 and assign it the value $\pi(0) = 1$.

(We will be scaling the $\pi(\cdot)$ values later appropriately, if our algorithm is able to terminate properly.)

If node 1 has an edge with value $p(0, 1)$ coming into it, for the detailed balance condition to be satisfied, we need $\pi(1) = \pi(0) \times p(0, 1)/p(1, 0)$. So $\pi(1)$ is fixed.

We repeat this process:

Starting from the set of nodes V for which the $\pi(\cdot)$ value is currently fixed, check if any out going edges are there to other nodes. Suppose there is an edge from node $j \in V$ to a node k outside. Fix the value $\pi(k)$ as $\pi(k) = \pi(j) \times p(j, k)/p(k, j)$.

We will stop when there are no outgoing edges from V . This would also mean that we have assigned a $\pi(\cdot)$ value to each node in the graph. Observe that by this time, within a scaling factor, $\pi(\cdot)$ is unique. The scaling factor arises because we could have assigned any positive value for $\pi(0)$.

However, we do not know whether for *every* edge (i, j) , the detailed balance condition is satisfied. So we check this now. If for some edge there is failure, we declare the Markov chain is not reversible.

If there is no failure, the chain is reversible and we scale the $\pi(\cdot)$ values so that they add up to 1. This is the desired stationary probability distribution.

An immediate consequence of the above discussion is that if

- every edge has a parallel edge in the opposite direction
- there are no directed loops other than the parallel edges

the Markov chain is reversible. This is because by the time all nodes have been assigned $\pi(\cdot)$ values by our algorithm, the detailed balance conditions would have also been verified since there are no other edges.

Example: The Ehrenfest chain graph is a simple straight line, if we replace parallel edges with single edges. So there are no loops except the parallel edges and the chain is reversible.

The algorithm is illustrated with an example below.

A test for Reversibility of Markov Chains

Assume that M is an irreducible aperiodic markov chain. Thus it has a unique stationary distribution, with a positive probability at every vertex.

Step 1 Build the R-graph of the Markov chain.

R-graph: Vertices are states of the Markov Chain.If p_{ij} is not equal to zero, put a directed edge $e(i, j)$ from i to j with weight $w_{ij} = \frac{p_{ij}}{p_{ji}}$ in the case of FDTM . If the weight of any edge is infinite STOP. The MC is not reversible since we have

$$\pi_i p_{ij} \neq \pi_j p_{ji} \tag{2}$$

for any stationary distribution π on the states.

Example:

Let the transition matrix in the case of FDTM be the matrix A below

$$A = \begin{bmatrix} a_{11} & a_{12} & 0 & a_{14} \\ a_{21} & a_{22} & a_{23} & 0 \\ 0 & a_{32} & a_{33} & a_{34} \\ a_{41} & 0 & a_{43} & a_{44} \end{bmatrix}$$

The R-graph is shown in Figure 3

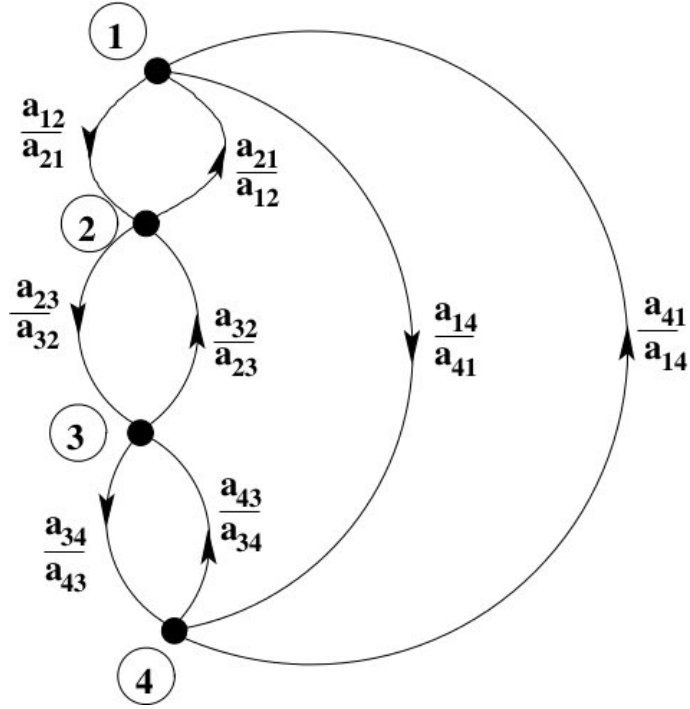


Figure 3:

Step 2 a Start from any node n . Assign it value 1. Suppose a set N_A of nodes have been assigned values. If N_A is not the full set of nodes pick any edge from $i \in N_A$ to $j \notin N_A$. Assign $\pi_j = w_{ij}\pi_i$. If no edges leave these nodes to the (non null) complement declare the Markov Chain to be not irreducible.

b If N_A is the full set of nodes GO TO STEP 3.

Step 3 For every edge $e(i, j)$ verify that $\pi_i w_{ij} = \pi_j$. If this is violated declare chain is not reversible. If not, Scale π_i by k so that

$$\sum k\pi_i = 1 \tag{3}$$

Output π_i as the stationary distribution and STOP.

Justification : If step 3 is satisfied we have the detailed balance equation satisfied for every pair $\{i, j\}$ for which $p_{ij} \neq 0$ and hence the MC is reversible.

4. New Markov Chains from old

- (a) *Merging*: Suppose a subset S_k of k states all have the same probability under the stationary distribution $\pi(\cdot)$. We build a new Markov chain by fusing all these states into a single ‘super-state’. All the edges leaving the original states have the probabilities associated multiplied by $1/k$, edges entering have the same probability as before, edges going between nodes in S_k become self loops, with probability multiplied by $1/k$. If parallel edges (of the same direction) result they are replaced by a single edge with the new probability equal to the sum of those of the original parallel edges. If several self loops result at a node, they can be replaced by a single selfloop at the same node with the edge probabilities added. For this Markov chain, we can show that the new stationary distribution $\pi'(\cdot)$ can be obtained by putting $\pi'(S_k) = k \times \pi(i), i \in S_k$, and leaving all other probabilities unchanged.

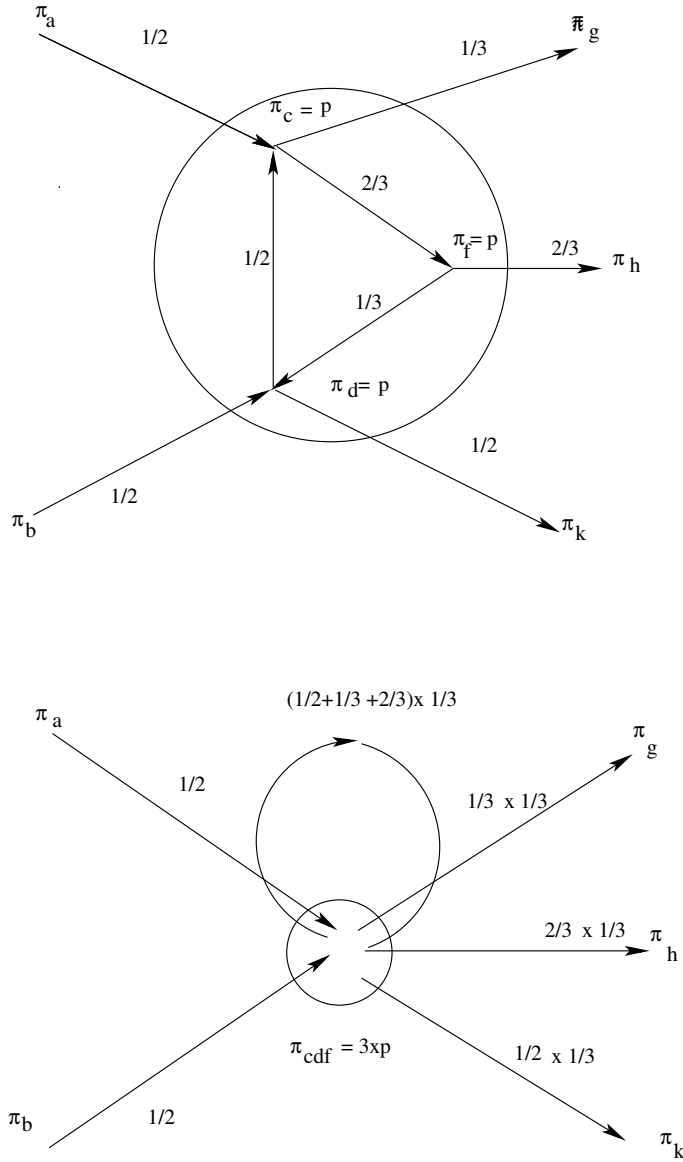


Figure 4: The merging process

The proof of these statements is best done by visualizing the above process of merging in terms

of the transition matrix.

Let P be the transition matrix. Let S_2 be the set of k states to be merged and S_1 be the complement. The k states in S_2 have the same $\pi(\cdot)$ value. Let $s \in S_2$. For the original Markov chain we have $\pi^T P = \pi^T$. After partitioning the rows and columns according to S_1, S_2 ,

$$(\pi_1^T, \pi_2^T) \left(\begin{array}{c|c} P_{11} & P_{12} \\ \hline P_{21} & P_{22} \end{array} \right) = (\pi_1^T, \pi_2^T). \quad (4)$$

Now this is the same as

$$(\pi_1^T, \pi(s) \times k) \left(\begin{array}{c|c} P_{11} & P_{12} \\ \hline \hat{p}_{21} & \hat{p}_{22} \end{array} \right) = (\pi_1^T, \pi_2^T), \quad (5)$$

where $(\hat{p}_{21}, \hat{p}_{22})$ is $1/k \times$ *sum of rows of* (P_{21}, P_{22}) . Now in the above equation, if we sum the second set of columns of the matrix

$$\left(\begin{array}{c|c} P_{11} & P_{12} \\ \hline \hat{p}_{21} & \hat{p}_{22} \end{array} \right), \quad (6)$$

on the LHS, yielding

$$\left(\begin{array}{c|c} P_{11} & \tilde{P}_{12} \\ \hline \hat{p}_{21} & \tilde{p}_{22} \end{array} \right), \quad (7)$$

the equation would remain correct if we sum the second set of columns of the row vector on the right side too. But summing the second set of columns of (π_1^T, π_2^T) , yields $(\pi_1^T, \pi(s) \times k)$.

We thus have

$$(\pi_1^T, \pi(s) \times k) \left(\begin{array}{c|c} P_{11} & \tilde{P}_{12} \\ \hline \hat{p}_{21} & \tilde{p}_{22} \end{array} \right) = (\pi_1^T, \pi(s) \times k). \quad (8)$$

To verify that the matrix in (7), is the transition matrix of the new merged Markov chain, we just note that the edges entering correspond to entries in \tilde{P}_{12} , edges leaving correspond to entries in \hat{p}_{21} , and self loops correspond to entries in \tilde{p}_{22} . It follows therefore that $\pi'(\cdot) = (\pi_1^T, \pi(s) \times k)$, is the stationary distribution for the merged Markov chain.

- (b) *splitting*: This operation could be regarded as one way of reversing the operation of ‘merging’. In the old Markov chain graph, we take any node and split it into k nodes. Every incoming edge should be duplicated k times, with the edge probability $1/k$ times the previous one, and feed into each of the split nodes. Every outgoing edge must be duplicated k times and feed out of each of the split nodes, with the same edge probability. Every self loop must be duplicated k times and attached to each of the resulting split nodes with the same probability as before. Let the old stationary distribution be $\pi(\cdot)$, and the new one be $\pi'(\cdot)$. Let node s in the old Markov chain graph split into the identical k nodes s_i which form set S_k . Then we can show $\pi'(s_i) = 1/k \times \pi(s)$. The probabilities associated with the nodes which have not been split remain as before.

Observe that splitting followed by merging will return to the original Markov chain, but merging followed by splitting might not.

(c) *Metropolis-Hastings method* This is a powerful method of building a reversible Markov chain which has a desired stationary distribution $\pi(\cdot)$, on a given set of states. Let X_n be a Markov chain with edge probability $q(i, j)$, with the additional condition that $q(i, j) \neq 0$ implies $q(j, i) \neq 0$. Let $r(i, j)$ denote $\min [\pi(j)q(j, i)/\pi(i)q(i, j), 1]$.

Let Y_n be the Markov chain on the same graph but with transition probability

$$p(i, j) = q(i, j)r(i, j).$$

Let us verify that Y_n satisfies the detailed balance condition with respect to the distribution $\pi(\cdot)$. Suppose $\pi(j)q(j, i) \leq \pi(i)q(i, j)$. We then have

$$\pi(i)p(i, j) = \pi(i)q(i, j)r(i, j) = \pi(i)q(i, j) \times \pi(j)q(j, i)/\pi(i)q(i, j) = \pi(j)q(j, i).$$

$$\pi(j)p(j, i) = \pi(j)q(j, i)r(j, i) = \pi(j)q(j, i) \times 1 = \pi(j)q(j, i),$$

verifying the detailed balance condition for Y_n .

STEP 4 in the study of FDTM: DO LOTS OF PROBLEMS.