

CONTINUOUS-TIME MARKOV CHAINS

by

Ward Whitt

Department of Industrial Engineering and Operations Research
Columbia University
New York, NY 10027-6699
Email: ww2040@columbia.edu
URL: www.columbia.edu/~ww2040

December 6, 2014

©Ward Whitt

Contents

1	Introduction	1
2	Transition Probabilities and Finite-Dimensional Distributions	2
3	Modelling	4
3.1	A DTMC with Exponential Transition Times	6
3.2	Transition Rates and ODE's	7
3.3	Competing Clocks with Exponential Timers	11
3.4	Uniformization: A DTMC with Poisson Transitions	14
4	Birth-and-Death Processes	16
5	Stationary and Limiting Probabilities for CTMC's	23
6	Reverse-Time CTMC's and Reversibility	31
7	Open Queueing Networks (OQN's)	36
7.1	The Model	36
7.2	The Traffic Rate Equations and the Stability Condition	36
7.3	The Limiting Product-Form Distribution	37
7.4	Extensions: More Servers or Different Service Scheduling Rules	39
7.5	Steady State in Discrete and Continuous Time	40
8	Closed Queueing Networks (CQN's)	41
8.1	Why Are Closed Models Interesting?	41
8.2	A Normalized Product-Form Distribution	42
8.3	Computing the Normalization Constant: The Convolution Algorithm	44
9	Stochastic Loss Models	45
9.1	The Erlang Loss Model	45
9.2	Stochastic Loss Networks	47
9.3	Insensitivity in the Erlang Loss Model	48
10	Regularity and Irregularity in Infinite-State CTMC's	50
10.1	Motivating Examples	51
10.2	Instantaneous Transitions, Explosions and the Minimal Construction	52
10.3	Conditions for Regularity and Recurrence	53
11	More on Reversible CTMC's and Birth-and-Death Processes	54
11.1	Spectral Representation in Reversible Markov Chains	54
11.2	Fitting BD Processes to Data	56
11.3	Comparing BD processes	57
11.4	First-Passage Times in BD Processes	59
12	Some Next Steps	61

1. Introduction

We now turn to **continuous-time Markov chains (CTMC's)**, which are a natural sequel to the study of discrete-time Markov chains (DTMC's), the Poisson process and the exponential distribution, because CTMC's combine DTMC's with the Poisson process and the exponential distribution. Most properties of CTMC's follow directly from results about DTMC's, the Poisson process and the exponential distribution. .

Like DTMC's, CTMC's are Markov processes that have a discrete state space, which we can take to be the positive integers. Just as with DTMC's, we will initially (in §§1-5) focus on the special case of a **finite state space**, but the theory and methods extend to infinite discrete state spaces, provided we impose additional regularity conditions; see §10. We will usually assume that the state space is the set $\{0, 1, 2, \dots, n\}$ containing the first $n + 1$ nonnegative integers for some positive integer n , but any finite set can be so labelled. Just as with DTMC's, a finite state space allows us to apply square (finite) matrices and elementary linear algebra. The main difference is that we now consider **continuous time**. We consider a stochastic process $\{X(t) : t \geq 0\}$, where time t is understood to be any nonnegative real number. The random variable $X(t)$ is the state occupied by the CTMC at time t .

As we will explain in §3, a CTMC can be viewed as a DTMC with altered transition times. Instead of unit times between successive transitions, the times between successive transitions are allowed to be independent exponential random variables with means that depend only on the state from which the transition is being made. Alternatively, as we explain in §3.4, a CTMC can be viewed as a DTMC (a different DTMC) in which the transition times occur according to a Poisson process. In fact, we already have considered a CTMC with just this property (but infinite state space), because the Poisson process itself is a CTMC. For that CTMC, the associated DTMC starts in state 0 and has only unit upward transitions, moving from state i to state $i + 1$ with probability 1 for all i . A CTMC generalizes a Poisson process by allowing other transitions. For a Poisson process, $X(t)$ goes to infinity as $t \rightarrow \infty$. We will be interested in CTMC's that have proper limiting distributions as $t \rightarrow \infty$.

Organization. Here is how the chapter is organized: We start in §2 by discussing transition probabilities and the way they can be used to specify the finite-dimensional distributions, which in turn specify the probability law of the CTMC. Then in §3 we describe four different ways to construct a CTMC model, giving concrete examples. In §4 we discuss the special case of a birth-and-death process, in which the only possible transitions are up one or down one to a neighboring state. The number of customers in a queue (waiting line) can often be modeled as a birth-and-death process. The special structure of a birth-and-death process makes the limiting probabilities especially easier to compute. Afterwards, in §5 we indicate how to calculate the limiting probabilities for a general irreducible CTMC. There are different ways, with the one that is most convenient usually depending on the modeling approach.

The second part is more advanced, focusing on reversibility and stochastic networks. We start in §6 by introducing reverse-time CTMC's and reversibility. We apply those notions to a CTMC consisting of several queues in series. in §§7 and 8 we present the basic theory of open and closed queueing networks, respectively. in §9 we discuss loss models, starting with the classical Erlang loss model and then continuing to multi-class multi-facility generalizations: stochastic loss networks.

We conclude with a brief treatment of some more advanced topics. In §10 we discuss the regularity conditions needed for infinite-state CTMC's. Finally, we discuss four special topics for reversible CTMC's and birth-and-death processes: (i) rates of convergence to steady state characterized via the spectral representation of the transition function, (ii) fitting birth-and-

death models to data, (iii) stochastic comparisons for birth-and-death processes and (iv) ways to compute first-passage-time distributions in birth-and-death processes. Much more material is available in the references.

2. Transition Probabilities and Finite-Dimensional Distributions

Just as with discrete time, a continuous-time stochastic process is a **Markov process** if the conditional probability of a future event given the present state and additional information about past states depends only on the present state. A **CTMC** is a continuous-time Markov process with a discrete state space, which can be taken to be a subset of the nonnegative integers. That is, a stochastic process $\{X(t) : t \geq 0\}$ (with an integer state space) is a CTMC if

$$P(X(s+t) = j | X(s) = i, X(r) = i_r, r \in A_s \subseteq [0, s]) = P(X(s+t) = j | X(s) = i) \quad (2.1)$$

for all states i and j and for all times $s > 0$ and $t > 0$. On the left in (2.1), we are conditioning on the values of $X(r)$ for all times r in a subset of “past” times A_s in addition to the value at the “present” time s . In general, A_s could be an arbitrary subset of $[0, s) \equiv \{r : 0 \leq r < s\}$, but to have the conditional probability in (2.1) well defined by elementary methods, we assume that A_s is a finite subset.

The conditional probabilities $P(X(s+t) = j | X(s) = i)$ are called the **transition probabilities**. We will consider the special case of **stationary transition probabilities** (sometimes referred to as homogeneous transition probabilities), occurring when

$$P(X(s+t) = j | X(s) = i) = P(X(t) = j | X(0) = i) \equiv P_{i,j}(t) \quad (2.2)$$

for all states i and j and for all times $s > 0$ and $t > 0$; the independence of s characterizes the stationarity. We assume stationary transition probabilities unless stipulated otherwise.

Thus a key concept for CTMC’s is the notion of transition probabilities. However, the transition probabilities of CTMC’s are not so easy to work with. As a consequence, we usually do not directly use transition probabilities when we construct and analyze CTMC models. First, when we construct a CTMC model, we invariably do not directly define the transition probabilities (although their structure will be implied by what we do define). Second, after constructing a CTMC model, we usually do not calculate the transition probabilities. Instead, we usually calculate the associated **limiting probabilities**, denoted by α_j :

$$\alpha_j \equiv \lim_{t \rightarrow \infty} P_{i,j}(t) \equiv \lim_{t \rightarrow \infty} P(X(t) = j | X(0) = i), \quad (2.3)$$

because they are much easier to calculate, and because they usually serve as excellent approximations for the exact transition probabilities $P_{i,j}(t)$ when t is large. (We use the notation α for the limiting probability vector of the CTMC, instead of π , because we reserve π for the limiting probability vector for an associated DTMC; see §3.1 and Theorem 5.2.)

Consistent with what we have written in (2.3), under regularity conditions, the limiting probabilities α_j will not depend on the initial state. Indeed, that will be true provided the CTMC is **irreducible**, which means (just as in discrete time) that it is possible with some positive probability to get from any state to any other state at some finite time, which may involve multiple transitions. (Just as in discrete time, for irreducibility, we do not require that we reach these other states in a single transition.) We assume irreducible CTMC’s unless stipulated otherwise.

This chapter is largely about constructing CTMC models and calculating the limiting probability vector $\alpha \equiv (\alpha_0, \alpha_1, \dots, \alpha_n)$. As with DTMC’s, we will also want to apply the

limiting probability vector α to answer a variety of related questions of interest. But, to repeat, neither constructing the CTMC model nor calculating the limiting probability vector α will directly involve the transition probabilities. Nevertheless, the transition probabilities are very important for understanding CTMC's.

Just as in discrete time, the evolution of the transition probabilities over time is described by the Chapman-Kolmogorov equations, but they take a different form in continuous time. In formula (2.4) below, we consider a sum over all possible states at some intermediate time. In doing so, we simply write a sum over integers. When we do that, we understand the sum to be over all possible states.

Lemma 2.1. (Chapman-Kolmogorov equations) For all $s \geq 0$ and $t \geq 0$,

$$P_{i,j}(s+t) = \sum_k P_{i,k}(s)P_{k,j}(t) . \quad (2.4)$$

Proof. We can compute $P_{i,j}(s+t)$ by considering all possible places the chain could be at time s . We then condition and uncondition, invoking the Markov property to simplify the conditioning; i.e.,

$$\begin{aligned} P_{i,j}(s+t) &= P(X(s+t) = j | X(0) = i) \\ &= \sum_k P(X(s+t) = j, X(s) = k | X(0) = i) \\ &= \sum_k P(X(s) = k | X(0) = i) P(X(s+t) = j | X(s) = k, X(0) = i) \\ &= \sum_k P(X(s) = k | X(0) = i) P(X(s+t) = j | X(s) = k) \quad (\text{Markov property}) \\ &= \sum_k P_{i,k}(s) P_{k,j}(t) \quad (\text{stationary transition probabilities}) . \quad \blacksquare \end{aligned}$$

Using matrix notation, we write $P(t)$ for the square matrix of transition probabilities ($P_{i,j}(t)$), and call it the **transition function**. In matrix notation, the Chapman-Kolmogorov equations reduce to a simple relation among the transition functions involving **matrix multiplication**:

$$P(s+t) = P(s)P(t) \quad (2.5)$$

for all $s \geq 0$ and $t \geq 0$. It is important to recognize that (2.5) means (2.4). From the perspective of abstract algebra, equation (2.5) says that the transition function has a semi-group property, where the single operation is matrix multiplication. (It is not a group because an inverse is missing.)

A CTMC is well specified if we specify: (1) its initial probability distribution - $P(X(0) = i)$ for all states i - and (2) its transition probabilities - $P_{i,j}(t)$ for all states i and j and positive times t . First, we can use these two elements to compute the distribution of $X(t)$ for each t , namely,

$$P(X(t) = j) = \sum_i P(X(0) = i) P_{i,j}(t) . \quad (2.6)$$

However, in general, we want to do more. We want to know about the joint distributions in order to capture the dependence structure. Recall that the **probability law of a stochastic process** is understood to be the set of all its finite-dimensional distributions. A **finite-dimensional distribution** is

$$P(X(t_1) = j_1, X(t_2) = j_2, \dots, X(t_k) = j_k) \quad (2.7)$$

for states j_i and times t_i satisfying $0 \leq t_1 < t_2 < \dots < t_k$. The probability law is specified by all these finite-dimensional distributions, considering all positive integers k , and all sets of k states and k ordered times. It is important that we can express any finite-dimensional distribution in terms of the initial distribution and the transition probabilities. For example, assuming that $t_1 > 0$, we have

$$\begin{aligned} P(X(t_1) = j_1, X(t_2) = j_2, \dots, X(t_k) = j_k) \\ = \sum_{j_0} P(X(0) = j_0) P_{j_0, j_1}(t_1) P_{j_1, j_2}(t_2 - t_1) \times \dots \times P_{j_{k-1}, j_k}(t_k - t_{k-1}) . \end{aligned} \quad (2.8)$$

In summary, equation (2.8) shows that we succeed in specifying the full probability law of the DTMC, as well as all the marginal distributions via (2.6), by specifying the initial probability distribution - $P(X(0) = i)$ for all i - and the transition probabilities $P_{i,j}(t)$ for all t , i and j or, equivalently, the transition function $P(t)$. However, when we construct CTMC models, as we do next, we do not directly specify the transition probabilities. We will see that, at least in principle, the transition probabilities can be constructed from what we do specify, but we usually do not carry out that step.

3. Modelling

We now turn to modelling: constructing a CTMC model. We saw that a DTMC model is specified by simply specifying its one-step transition matrix P and the initial probability distribution. Unfortunately, the situation is more complicated in continuous time.

In this section we will describe **four different approaches** to constructing a CTMC model. With each approach, we will need to specify the initial distribution, so we are focusing on specifying the model beyond the initial distribution. The four approaches are equivalent: You get to the same result from each and you can get to each from any of the others. Even though these four approaches are redundant, they are useful because they together give a different more comprehensive view of a CTMC. We see different things from different perspectives, much like the Indian fable about the blind men and the elephant, recaptured in the poem by John Godfrey Saxe (1816-1887):

The Blind Men and the Elephant

It was six men of Indostan
To learning much inclined,
Who went to see the Elephant
(Though all of them were blind),
That each by observation
Might satisfy his mind.

The First approached the Elephant,
And happening to fall
Against his broad and sturdy side,
At once began to bawl:
"God bless me! but the Elephant
Is very like a wall!"

The Second, feeling of the tusk,
Cried, "Ho! what have we here
So very round and smooth and sharp?
To me tis mighty clear
This wonder of an Elephant
Is very like a spear!"

The Third approached the animal,
And happening to take
The squirming trunk within his hands,
Thus boldly up and spake:
"I see," quoth he, "the Elephant
Is very like a snake!"

The Fourth reached out an eager hand,
And felt about the knee.
"What most this wondrous beast is like
Is mighty plain," quoth he;
" Tis clear enough the Elephant
Is very like a tree!"

The Fifth, who chanced to touch the ear,
Said: "Een the blindest man
Can tell what this resembles most;
Deny the fact who can
This marvel of an Elephant
Is very like a fan!"

The Sixth no sooner had begun
About the beast to grope,
Than, seizing on the swinging tail
That fell within his scope,
"I see," quoth he, "the Elephant
Is very like a rope!"

And so these men of Indostan
Disputed loud and long,
Each in his own opinion
Exceeding stiff and strong,
Though each was partly in the right,
And all were in the wrong!

For some applications, one modelling approach may be more natural than the others. Or one modelling approach may be more convenient for analyzing the model.

3.1. A DTMC with Exponential Transition Times

In order to construct a CTMC model, it is natural to build on our knowledge of DTMC's. So we first consider a way to exploit DTMC's in our construction of the CTMC model. To do so in the strongest way, we start with a DTMC having a transition matrix P , and then modify the way the transitions occur. Instead of having each transition take unit time, now we assume that each transition takes a random time. In particular, we assume that the time required to make a transition from state i has an exponential distribution with rate ν_i , and thus mean $1/\nu_i$, independent of the history before reaching state i .

This modelling approach is convenient for simulating the CTMC; we can recursively generate successive transitions. This modelling approach also avoids technical complications that arise in the conventional transition-rate approach, to be introduced in the next subsection. This modelling approach is also appealing because many applications are naturally expressed in this way.

Example 3.1. (Pooh Bear and the Three Honey Trees) A bear of little brain named Pooh is fond of honey. Bees producing honey are located in three trees: tree A , tree B and tree C . Tending to be somewhat forgetful, Pooh goes back and forth among these three honey trees randomly (in a Markovian manner) as follows: From A , Pooh goes next to B or C with probability $1/2$ each; from B , Pooh goes next to A with probability $3/4$, and to C with probability $1/4$; from C , Pooh always goes next to A . Pooh stays a random time at each tree. (Assume that the travel times can be ignored.) Pooh stays at each tree an exponential length of time, with the mean being 5 hours at tree A or B , but with mean 4 hours at tree C . Construct a CTMC enabling you to find the limiting proportion of time that Pooh spends at each honey tree.

Note that this problem is formulated directly in terms of the DTMC, describing the random motion at successive transitions, so it is natural to use this initial modelling approach. Here the transition matrix for the DTMC is

$$P = \begin{array}{c} A \\ B \\ C \end{array} \left(\begin{array}{ccc} 0 & 1/2 & 1/2 \\ 3/4 & 0 & 1/4 \\ 1 & 0 & 0 \end{array} \right) .$$

In the displayed transition matrix P , we have only labelled the rows. The columns are assumed to be labeled in the same order. As specified above, the exponential times spent at the three trees have means $1/\nu_A = 1/\nu_B = 5$ hours and $1/\nu_C = 4$ hours. ■

Given that we have already studied DTMC's, it is natural to wonder how the steady-state probability vector of the CTMC is related to the steady-state probability vector of the DTMC. For the DTMC with transition matrix P (looking at the transition epochs of the CTMC), the steady state probability vector is π , the unique probability vector satisfying the equation

$$\pi = \pi P . \tag{3.1}$$

It is significant that, **in general, the steady-state probability vector of the CTMC is not π** . Thus, we use different notation, referring to the steady-state probability vector of the CTMC as α . Fortunately, though, the two steady-state probability vectors turn out to be intimately related. In particular, in §5 we will see that

$$\alpha_j = \frac{(\pi_j/\nu_j)}{\sum_k (\pi_k/\nu_k)} . \tag{3.2}$$

In §5 we will see how to justify formula (3.2) above and relate it to other ways to calculate the limiting probabilities for this CTMC. We will then be able to answer the question about the long-run proportion of time that Pooh spends at each tree.

Finally, when we study renewal theory and semi-Markov processes, we will see that the same steady-state probability vector α in (3.2) also holds if the random holding times in each state are not exponentially distributed, provided that they are independent and identically distributed (i.i.d.) with those same means $1/\nu_j$. Indeed, this first modelling approach corresponds to treating the CTMC as a special case of a **semi-Markov process (SMP)**. An SMP is a DTMC with independent random transition times, but it allows the distributions of the intervals between transitions to be non-exponential.

With this initial modelling approach, it is natural to assume, as was the case in Example 3.1, that there are no one-step transitions in the DTMC from any state immediately back to itself, but it is not necessary to make that assumption. We get a CTMC from a DTMC and exponential transition times without making that assumption.

However, to help achieve a simple relation between the first two modelling approaches, we make that assumption here: **We assume that there are no one-step transitions from any state to itself in the DTMC**; i.e., we assume that $P_{i,i} = 0$ for all i . However, we emphasize that this assumption is not critical, as we will explain after we introduce the third modelling approach. Indeed, we will want to allow transitions from a state immediately to itself in the fourth - uniformization - modelling approach. That is a crucial part of that modelling approach.

3.2. Transition Rates and ODE's

A second modelling approach is based on representing the transition probabilities as the solution of a system of ordinary differential equations, which allows us to apply well-established modelling techniques from the theory of differential equations in a deterministic setting; e.g., see Simmons (1991). With this second modelling approach, we directly specify transition rates.

We proceed with that idea in mind, but without assuming knowledge about differential equations. We focus on the transition probabilities of the CTMC, even though they have not yet been specified. With the transition probabilities in mind, we assume that there are well-defined derivatives (from above or from the right) of the transition probabilities at 0. We assume these derivatives exist, and call them transition rates.

But first we must define **zero-time transition probabilities**, which we do in the obvious way: We let $P(0) = I$, where I is the identity matrix; i.e., we set $P_{i,i}(0) = 1$ for all i and we set $P_{i,j}(0) = 0$ whenever $i \neq j$. We are just assuming that you cannot go anywhere in zero time.

We then let the transition rate from state i to state j be defined in terms of the **derivatives**:

$$Q_{i,j} \equiv P'_{i,j}(0+) \equiv \left. \frac{dP_{i,j}(t)}{dt} \right|_{t=0+} . \quad (3.3)$$

In (3.3) $0+$ appears to denote the right derivative at 0 because $P_{i,j}(t)$ is not defined for $t < 0$.

This approach is used in most treatments of CTMC's, but without mentioning derivatives or right-derivatives. Instead, it is common to assume that

$$P_{i,j}(h) = Q_{i,j}h + o(h) \quad \text{as } h \downarrow 0 \quad \text{if } j \neq i \quad (3.4)$$

and

$$P_{i,i}(h) - 1 = Q_{i,i}h + o(h) \quad \text{as } h \downarrow 0 , \quad (3.5)$$

where $o(h)$ is understood to be a quantity which is asymptotically negligible as $h \downarrow 0$ after dividing by h . (Formally, $f(h) = o(h)$ as $h \downarrow 0$ if $f(h)/h \rightarrow 0$ as $h \downarrow 0$.)

For a finite state space, which we have assumed, and for infinite state spaces under extra regularity conditions, we will have

$$-Q_{i,i} \equiv \sum_{j,j \neq i} Q_{i,j} \quad (3.6)$$

because the transition probabilities $P_{i,j}(t)$ sum over j to 1. Moreover, we have

$$-Q_{i,i} = \nu_i \quad \text{for all } i, \quad (3.7)$$

because we have assumed that $P_{i,i} = 0$ in the first modelling approach.

In other words, these two assumptions mean that

$$\lim_{h \downarrow 0} \frac{P_{i,j}(h) - P_{i,j}(0)}{h} = Q_{i,j} \quad \text{for all } i \text{ and } j, \quad (3.8)$$

which is just what is meant by (3.3).

In summary, we first assumed that transition probabilities are well defined, at least for zero time and small positive time intervals, and then assume that they are differentiable from the right at 0. We remark that it is possible to weaken that assumption, and only assume that the transition probabilities are continuous at 0: $P(h) \rightarrow P(0) \equiv I$ as $h \downarrow 0$. Then it is possible to prove that the derivatives exist; see §§II.1 and II.2 of Chung (1967).

Having defined the transition rates in terms of the assumed behavior of the transition probabilities in a very short (asymptotically negligible) interval of time, we can specify the CTMC model by specifying these transition rates; i.e., we specify the transition-rate matrix Q , having elements $Q_{i,j}$. (But we do not first fully define the transition probabilities themselves!) Thus, just as we specify a DTMC model via a matrix P , we can specify a CTMC model via the transition-rate matrix Q .

When specifying the transition-rate matrix Q , it suffices to specify the off-diagonal elements $Q_{i,j}$ for $i \neq j$, because the diagonal elements $Q_{i,i}$ are always defined by (3.6). The off-diagonal elements are always nonnegative, whereas the diagonal elements are always negative. Each row sum of Q is zero.

Even though this modelling approach for CTMC's is similar to what we did for DTMC's, it is more complicated, because the rate matrix Q is harder to interpret than the one-step transition matrix P . (The discussion above is intended to help interpretation.) In fact, this approach to CTMC modelling is perhaps best related to modelling with ordinary differential equations, as mentioned at the beginning of this subsection.

To construct the transition probabilities $P_{i,j}(t)$ from the transition rates $Q_{i,j} \equiv P'_{i,j}(0+)$, we apply the Chapman-Kolmogorov equations in Lemma 2.1 in order to show that the transition probabilities satisfy **two systems of ordinary differential equations (ODE's)** generated by the transition rates. In matrix notation, these will be simple first-order linear ODE's.

Theorem 3.1. (Kolmogorov forward and backward ODE's) *The transition probabilities satisfy both the Kolmogorov **forward** differential equations*

$$P'_{i,j}(t) = \sum_k P_{i,k}(t) Q_{k,j} \quad \text{for all } i \text{ and } j, \quad (3.9)$$

which in matrix notation is the matrix ODE

$$P'(t) = P(t)Q, \quad (3.10)$$

and the Kolmogorov **backward** differential equations

$$P'_{i,j}(t) = \sum_k Q_{i,k} P_{k,j}(t) \quad \text{for all } i \text{ and } j, \quad (3.11)$$

which in matrix notation is the matrix ODE

$$P'(t) = QP(t), \quad (3.12)$$

Proof. We start with the forward equation, using matrix notation. We apply the Chapman-Kolmogorov equations in Lemma 2.1 to write

$$P(t+h) = P(t)P(h),$$

and then do an asymptotic analysis as $h \downarrow 0$. (This is tantamount to doing a careful asymptotic analysis of what happens in a small interval after time t .) We subtract $P(t)$ from both sides and divide by h , to get

$$\frac{P(t+h) - P(t)}{h} = P(t) \frac{P(h) - I}{h},$$

where I is the identity matrix. Recalling that $I = P(0)$, we can let $h \downarrow 0$ to get the desired result (3.10). To get the backward equation (3.12), we start with

$$P(t+h) = P(h)P(t)$$

and reason in the same way. (This is tantamount to doing a careful asymptotic analysis of what happens in a small interval after time 0, and then applying $P(t)$ thereafter.) ■

To help remember which ODE is forward and which is backwards, note that $P(t)Q$ appearing on the righthand side of the forward ODE is in alphabetic order, whereas $QP(t)$ appearing on the righthand side of the backward ODE is in reverse (backward) alphabetic order.

With a finite state space, both ODE's are always well defined. With an infinite state space, there can be technical problems, because there could be infinitely many transitions in finite time; see §10. With an infinite state space, the forward ODE can be more problematic, because it presumes the process got to time t before doing the asymptotic analysis. Here we assume a finite state space, so we do not encounter those pathologies. Under regularity conditions, those pathologies will not occur with infinite state spaces either.

To obtain the transition function $P(t)$ from the transition-rate matrix Q , we can solve one of these ODE's. In preparation, we review the simple one-dimensional story. Suppose that we have an ODE $f'(t) = cf(t)$, where f is understood to be a differentiable real-valued function f with known initial value $f(0)$. If we divide both sides by $f(t)$, we get $f'(t)/f(t) = c$. Since $f'(t)/f(t)$ is the derivative of $\log f(t)$, we can integrate to get

$$\log f(t) - \log f(0) = ct \quad \text{or} \quad f(t) = f(0)e^{ct}, \quad t \geq 0.$$

Thus we see that f must be an exponential function.

Closely paralleling the real-valued case, the matrix ODE's in (3.10) and (3.12) have an exponential solution, but now a matrix-exponential solution. (Since $P(0) = I$, the initial condition plays no role, just as above when $f(0) = 1$.) In particular, as a consequence of Theorem 3.1, we have the following corollary.

Theorem 3.2. (matrix exponential representation) *The transition function can be expressed as a matrix-exponential function of the rate matrix Q , i.e.,*

$$P(t) = e^{Qt} \equiv \sum_{n=0}^{\infty} \frac{Q^n t^n}{n!} \quad (3.13)$$

This matrix exponential is the unique solution to the two ODE's with initial condition $P(0) = I$.

Proof. If we verify or assume that we can interchange summation and differentiation in (3.13), we can check that the displayed matrix exponential satisfies the two ODE's:

$$\frac{d}{dt} \sum_{n=0}^{\infty} \frac{Q^n t^n}{n!} = \sum_{n=0}^{\infty} \frac{d}{dt} \frac{Q^n t^n}{n!} = \sum_{n=0}^{\infty} \frac{n Q^n t^{n-1}}{n!} = Q \sum_{n=0}^{\infty} \frac{Q^n t^n}{n!} = Q e^{Qt} .$$

We give a full demonstration at the end of §3.4. ■

However, in general the transition function $P(t)$ is not elementary to compute via (3.13); see Moler and Van Loan (2003). Indeed, one of the ways to evaluate the matrix-exponential function displayed in (3.13) is to numerically solve one of the ODE's as expressed in (3.10) or (3.12).

We now illustrate this second modelling approach with an example.

Example 3.2. (Copier Breakdown and Repair) Consider two copier machines that are maintained by a single repairman. Machine i functions for an exponentially distributed amount of time with mean $1/\gamma_i$, and thus rate γ_i , before it breaks down. The repair times for copier i are exponential with mean $1/\beta_i$, and thus rate β_i , but the repairman can only work on one machine at a time. Assume that the machines are repaired in the order in which they fail. Suppose that we wish to construct a CTMC model of this system, with the goal of finding the long-run proportions of time that each copier is working and the repairman is busy. How can we proceed?

An initial question is: What should be the state space? Can we use 4 states, letting the states correspond to the subsets of failed copiers? Unfortunately, the answer is “no,” because in order to have the Markov property we need to know which copier failed first when both copiers are down. However, we can use 5 states with the states being: 0 for no copiers failed, 1 for copier 1 is failed (and copier 2 is working), 2 for copier 2 is failed (and copier 1 is working), (1, 2) for both copiers down (failed) with copier 1 having failed first and being repaired, and (2, 1) for both copiers down with copier 2 having failed first and being repaired. (Of course, these states could be relabelled 0, 1, 2, 3 and 4, but we do not do that.)

From the problem specification, it is natural to work with transition rates, where these transition rates are obtained directly from the originally-specified failure rates and repair rates (the rates of the exponential random variables). In Figure 1 we display a **rate diagram** showing the possible transitions with these 5 states together with the appropriate rates. It can be helpful to construct such rate diagrams as part of the modelling process.

From Figure 1, we see that there are 8 possible transitions. The 8 possible transitions should clearly have transition rates

$$Q_{0,1} = \gamma_1, Q_{0,2} = \gamma_2, Q_{1,0} = \beta_1, Q_{1,(1,2)} = \gamma_2, Q_{2,0} = \beta_2, Q_{2,(2,1)} = \gamma_1, Q_{(1,2),2} = \beta_1, Q_{(2,1),1} = \beta_2 .$$

Rate Diagram

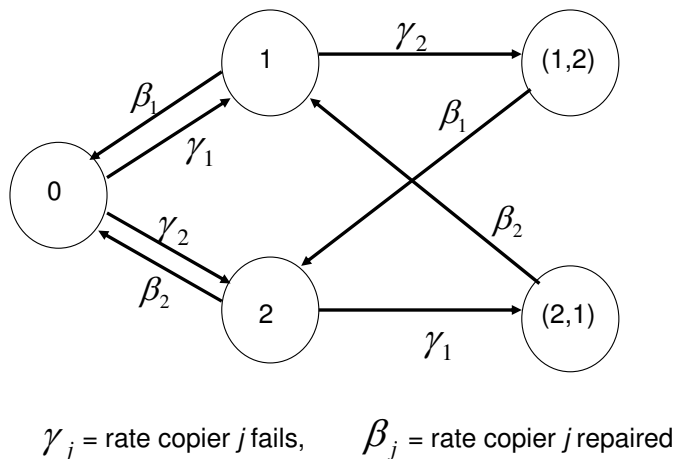


Figure 1: A rate diagram showing the transition rates among the 5 states in Example 3.2, involving copier breakdown and repair.

In other words, the rate matrix should be

$$Q = \begin{matrix} & \begin{matrix} 0 \\ 1 \\ 2 \\ (1,2) \\ (2,1) \end{matrix} \end{matrix} \begin{pmatrix} -(\gamma_1 + \gamma_2) & \gamma_1 & \gamma_2 & 0 & 0 \\ \beta_1 & -(\gamma_2 + \beta_1) & 0 & \gamma_2 & 0 \\ \beta_2 & 0 & -(\gamma_1 + \beta_2) & 0 & \gamma_1 \\ 0 & 0 & \beta_1 & -\beta_1 & 0 \\ 0 & \beta_2 & 0 & 0 & -\beta_2 \end{pmatrix}.$$

In §5, we will compute the limiting probability distribution of this CTMC and answer the questions posed above. ■

3.3. Competing Clocks with Exponential Timers

We now present a third modelling approach, which is an appealing constructive alternative to the second modelling approach based on rates of unknown transition functions. This third modelling approach is even more natural for Example 3.2. This third approach also helps link the first two modelling approaches.

With this third modelling approach, movement from state to state is determined by “competing” clocks with timers that go off at random, exponentially-distributed, times. For each state i , there is a clock associated with each state j the process could possibly move to in a single transition from state i . Let C_i be the set of states that the CTMC can possibly move to from state i in a single transition. Equivalently, C_i is the **set of active clocks** in state i . (We here assume that the process does not move from state i immediately back to state i .)

Each time the CTMC moves to state i , clocks with **timers are set or reset, if necessary** to go off at random times $T_{i,j}$ for each $j \in C_i$. Each clock has an exponential timer; i.e., the random time $T_{i,j}$ is given an exponential distribution with (positive finite) rate $Q_{i,j}$ and thus mean $1/Q_{i,j}$ (depending on i and j). Moreover, we assume that these newly set times $T_{i,j}$ are mutually independent and independent of the history of the CTMC prior to that transition time. By the lack-of-memory property of the exponential distribution, resetting running timers is equivalent (leaves the probability law of the stochastic process unchanged) to not resetting the times, and letting the timers continue to run.

Example 3.3. (Copier Breakdown and Repair Revisited) At this point we should reconsider Example 3.2 and observe that it is even more natural to define the CTMC through the proposed clocks with random timers. The random times triggering transitions are the exponential times to failure and times to repair specified in the original problem formulation. However, there is a difference: In the actual system, those random times do not get reset at each transition epoch. But, because of the lack-of-memory property of the exponential distribution, a timer that is still going can be reset at any time, including one of these random transition times, without changing the distribution of the remaining time. Thus the clocks with random timers does produce a valid representation of the desired CTMC model. ■

We now discuss the **implications** of the model specification in terms of exponential timers. As a consequence of these probabilistic assumptions, with probability 1, no two timers ever go off at the same time. (Since the exponential distribution is a continuous distribution, the probability of any single possible value is 0.) Let T_i be the time that the first timer goes off in state i and let N_i be the index j of the random timer that the first goes off; i.e.,

$$T_i \equiv \min_j \{T_{i,j}\} \quad (3.14)$$

and

$$N_i \equiv j \quad \text{such that} \quad T_{i,j} = T_i . \quad (3.15)$$

(The index j yielding the minimum is often called the *argmin*.) We then let the process move from state i next to state N_i after an elapsed time of T_i , and we repeat the process, starting from the new state N_i .

To understand the implications of these exponential clocks, recall basic properties of the exponential distribution. Recall that **the minimum of several independent exponential random variables** is again an exponential random variable with a rate equal to the sum of the rates. Hence, T_i has an exponential distribution, i.e.,

$$P(T_i \leq t) = 1 - e^{-\nu_i t}, \quad t \geq 0, \quad (3.16)$$

where

$$\nu_i \equiv -Q_{i,i} = \sum_{j,j \neq i} Q_{i,j}, \quad (3.17)$$

as in (3.6) and (3.7). (Again we use the assumption that $P_{i,i} = 0$ in the first modelling approach.)

Next, recall that, when considering several independent exponential random variables, each exponential random variable is **the exponential random variable yielding the minimum** with a probability proportional to its rate, so that

$$P(N_i = j) = \frac{Q_{i,j}}{\sum_{k,k \neq i} Q_{i,k}} = \frac{Q_{i,j}}{\nu_i} \quad \text{for} \quad j \neq i . \quad (3.18)$$

Finally, as discussed before in relation to the exponential distribution, the random variables T_i and N_i are **independent random variables**:

$$P(T_i \leq t, N_i = j) = P(T_i \leq t)P(N_i = j) = (1 - e^{-\nu_i t}) \left(\frac{Q_{i,j}}{\nu_i} \right) \quad \text{for all } t \text{ and } j .$$

After each transition, new timers are set, with the distribution of $T_{i,j}$ being the same at each transition to state i . So new timer values are set only at transition epochs. However, by the lack-of-memory property of the exponential distribution, the distribution of the remaining times $T_{i,j}$ and the associated random variables T_i and N_i would be the same any time we looked at the process in state i .

The analysis we have just done translates this clock formulation directly into a DTMC with exponential transition times, as in our first modelling approach in §3.1: The one-step transition matrix P of the DTMC is

$$P_{i,j} = P(N_i = j) = \frac{Q_{i,j}}{\sum_{k,k \neq i} Q_{i,k}} = \frac{Q_{i,j}}{\nu_i} \quad \text{for } j \neq i , \quad (3.19)$$

with $P_{i,i} = 0$ for all i , as specified in (3.18), while the rate ν_i of the exponential holding time in state i is specified in (3.17).

Moreover, it is easy to see how to define transition rates as required for the second modelling approach. We just let $Q_{i,j}$ be the rate of the exponential timer $T_{i,j}$. We have chosen the notation to make these associations obvious. Moreover, we can use the exponential timers to prove that the transition probabilities of the CTMC are well defined and do indeed have derivatives at the origin.

The construction here makes it clear how to relate the first two modelling approaches. Given the rate matrix Q , we define the one-step transition matrix \hat{P} of the DTMC by (3.19) and the rate $\hat{\nu}_i$ of the exponential transition time in state i by (3.17). That procedure gives us an underlying DTMC \hat{P} with $\hat{P}_{i,i} = 0$ for all i .

These equations also tell us how to go the other way: Given (P, ν) , we let

$$Q_{i,j} = \nu_i P_{i,j} \quad \text{for } j \neq i \quad \text{and} \quad Q_{i,i} = - \sum_{j,j \neq i} Q_{i,j} = \nu_i \quad \text{for all } i . \quad (3.20)$$

From this analysis, we see that the CTMC is uniquely specified by the rate matrix Q ; i.e., two different Q matrices produce two different CTMC's (two different probability laws, i.e., two different f.d.d.'s). That property also holds for the first modelling approach, provided that we assume that $P_{i,i} = 0$ for all i . Otherwise, the same CTMC can be represented by different pairs (P, ν) . There is only one if we require, as we have done, that there be no transitions from any state immediately back to itself.

We can also use this third modelling approach to show that the probability law of the CTMC is unaltered if there are initially one-step transitions from any state to itself. If we are initially given one-step transitions from any state to itself, we can start by removing them, but without altering the probability law of the original CTMC. If we remove a DTMC transition from state i to itself, we must compensate by increasing the transition probabilities to other states and increasing the mean holding time in state i . To do so, we first replace initial transition matrix P with transition matrix \hat{P} , where $\hat{P}_{i,i} = 0$ for all i . To do so without altering the CTMC, we must let the new transition probability be the old conditional probability given that there is no transition from state i to itself; i.e., we let

$$\hat{P}_{i,j} = \frac{P_{i,j}}{1 - P_{i,i}} \quad \text{for all } i \text{ and } j . \quad (3.21)$$

We never divide by zero, because $P_{i,i} < 1$ (assuming that the chain has more than two states and is irreducible). Since we have eliminated DTMC transitions from state i to itself, we must make the mean transition time larger to compensate. In particular, we replace $1/\nu_i$ by $1/\hat{\nu}_i$, where

$$1/\hat{\nu}_i = \frac{(1/\nu_i)}{1 - P_{i,i}} \quad \text{or} \quad \hat{\nu}_i = \nu_i(1 - P_{i,i}) . \quad (3.22)$$

Theorem 3.3. (removing transitions from a state back to itself) *The probability law of the CTMC is unaltered by removing one-step transitions from each state to itself, according to (3.21) and (3.22).*

Proof. The tricky part is recognizing what needs to be shown. Since (1) the transition rates determine the transition probabilities, as shown in §3.2, (2) the transition probabilities determine the finite-dimensional distributions and (3) the finite-dimensional distributions are regarded as the probability law of the CTMC, as shown in §2, it suffices to show that we have the right transition rates. So that is what we show.

Applying (3.20), we see that the transition rates of the new CTMC (denoted by a hat) are

$$\hat{Q}_{i,j} \equiv \hat{\nu}_i \hat{P}_{i,j} = \nu_i(1 - P_{i,i}) \frac{P_{i,j}}{(1 - P_{i,i})} = \nu_i P_{i,j} , \quad (3.23)$$

just as in (3.20). ■

In closing, we remark that this third modelling approach with independent clocks corresponds to treating the CTMC as a special case of a **generalized semi-Markov process (GSMP)**; e.g., see Glynn (1989). For general GSMP's, the clocks can run at different speeds and the timers can have nonexponential distributions.

3.4. Uniformization: A DTMC with Poisson Transitions

Our final modelling approach is not really a direct modelling approach, but rather an intermediate modelling approach, starting from the first modelling approach involving a DTMC with exponential transition times, that facilitates further analysis. Indeed, this modelling approach can be regarded as a special case of the first modelling approach. But it provides a different view of a CTMC.

In our first modelling approach, involving a DTMC with exponential transition times, the means of those transition times $1/\nu_i$ could vary from state to state. However, if it happened that these means were all the same, then we could represent the CTMC directly as a DTMC with transitions governed by an independent Poisson process, because in a Poisson process the times between transitions are IID exponential random variables.

Specifically, if the mean transition time is $1/\nu_0$ for all states, then we can generate all transitions from a Poisson process with rate ν_0 . Let $\{Y_n : n \geq 0\}$ be the DTMC with one-step transition matrix P and let $\{N(t) : t \geq 0\}$ be an independent Poisson process with rate ν_0 . Under that condition, the CTMC $\{X(t) : t \geq 0\}$ can be constructed as a **random time change** of the DTMC $\{Y_n : n \geq 0\}$ by the Poisson process $\{N(t) : t \geq 0\}$, i.e.,

$$X(t) = Y_{N(t)}, \quad t \geq 0 . \quad (3.24)$$

As a consequence,

$$P_{i,j}(t) \equiv P(X(t) = j | X(0) = i) = \sum_{k=0}^{\infty} P_{i,j}^k P(N(t) = k) = \sum_{k=0}^{\infty} P_{i,j}^k \frac{e^{-\nu_0 t} (\nu_0 t)^k}{k!} . \quad (3.25)$$

This situation may appear to be very special, but actually any finite-state CTMC can be represented in this way. We can achieve this representation by using the technique of **uniformization**, which means making the rates uniform or constant.

We make the rates uniform without changing the probability law of the CTMC by introducing one-step transitions from some states to themselves, which we can regard as **fictitious transitions**, because the process never actually moves. We can generate **potential transitions** from a Poisson process with rate λ , where λ is chosen so that

$$\nu_i \equiv -Q_{i,i} = \sum_{j,j \neq i} Q_{i,j} \leq \lambda \quad \text{for all } i, \quad (3.26)$$

as in (3.17).

When the CTMC is in state i , each of these potential transitions is a real transition (to another state) with probability ν_i/λ , while the potential transition is a fictitious transition (a transition from state i back to state i , meaning that we remain in state i at that time) with probability $1 - (\nu_i/\lambda)$, independently of past events. In other words, in each state i , we perform **independent thinning of the Poisson process** having rate λ , creating real transitions in state i according to a Poisson process having rate ν_i , just as in the original model.

The uniformization construction requires that we change the transition matrix of the embedded DTMC. The new one-step transition matrix allows transitions from a state to itself. In particular, the new one-step transition matrix \tilde{P} is constructed from the CTMC transition rate matrix Q and λ satisfying (3.26) by letting

$$\tilde{P}_{i,j} = \frac{Q_{i,j}}{\lambda} \quad \text{for } j \neq i \quad (3.27)$$

and

$$\tilde{P}_{i,i} = 1 - \sum_{j,j \neq i} \tilde{P}_{i,j} = 1 - \frac{\nu_i}{\lambda} = 1 + \frac{Q_{i,i}}{\lambda} = 1 - \frac{\sum_{j,j \neq i} Q_{i,j}}{\lambda}. \quad (3.28)$$

In matrix notation,

$$\tilde{P} = I + \lambda^{-1}Q. \quad (3.29)$$

Note that we have done the construction to ensure that \tilde{P} is a bonafide Markov chain transition matrix; it is nonnegative with row sums 1.

Uniformization is useful because it allows us to apply properties of DTMC's to analyze CTMC's. For the general CTMC characterized by the rate matrix Q , we have transition probabilities $P_{i,j}(t)$ expressed via \tilde{P} in (3.27)-(3.29) and λ as

$$P_{i,j}(t) \equiv P(X(t) = j | X(0) = i) = \sum_{k=0}^{\infty} \tilde{P}_{i,j}^k P(N(t) = k) = \sum_{k=0}^{\infty} \tilde{P}_{i,j}^k \frac{e^{-\lambda t} (\lambda t)^k}{k!}, \quad (3.30)$$

where \tilde{P} is the DTMC transition matrix constructed in (3.27)-(3.29). We also have representation (3.24) provided that the DTMC $\{Y_n : n \geq 0\}$ is governed by the one-step transition matrix \tilde{P} and the Poisson process $\{N(t) : t \geq 0\}$ has rate λ in (3.26).

But how do we know that equations (3.27) and (3.30) are really correct?

Theorem 3.4. (validity of uniformization) *The CTMC constructed via (3.27) and (3.30) leaves the probability law of the CTMC unchanged.*

Proof. We can justify the construction by showing that the transition rates are the same. Starting from (3.30), we see that, for $i \neq j$,

$$\begin{aligned} P_{i,j}(h) &= \sum_k \tilde{P}_{i,j}^k \frac{e^{-\lambda h} (\lambda h)^k}{k!} \\ &= \lambda h e^{-\lambda h} \tilde{P}_{i,j}^1 + o(h) = \lambda h e^{-\lambda h} \frac{Q_{i,j}}{\lambda} + o(h) = Q_{i,j} h + o(h) \quad \text{as } h \downarrow 0, \end{aligned} \quad (3.31)$$

consistent with (3.4), while

$$\begin{aligned} P_{i,i}(h) - 1 &= \sum_k \tilde{P}_{i,i}^k \frac{e^{-\lambda h} (\lambda h)^k}{k!} - 1 \\ &= \tilde{P}_{i,i}^0 e^{-\lambda h} + \lambda h e^{-\lambda h} \tilde{P}_{i,i}^1 + o(h) - 1 \\ &= (1 - \lambda h + o(h)) + (\lambda h + o(h)) \left(1 + \frac{Q_{i,i}}{\lambda}\right) + o(h) - 1 \\ &= Q_{i,i} h + o(h) \quad \text{as } h \downarrow 0, \end{aligned} \quad (3.32)$$

consistent with (3.5). ■

We now give a full proof of Theorem 3.2, showing that the transition function $P(t)$ can be expressed as the matrix exponential e^{Qt} .

Proof of Theorem 3.2. (matrix-exponential representation) Apply (3.27) to see that $\tilde{P} = \lambda^{-1}Q + I$. Then substitute for \tilde{P} in (3.30) to get

$$\begin{aligned} P(t) &= \sum_{k=0}^{\infty} \tilde{P}^k \frac{e^{-\lambda t} (\lambda t)^k}{k!} = \sum_{k=0}^{\infty} (\lambda^{-1}Q + I)^k \frac{e^{-\lambda t} (\lambda t)^k}{k!} = e^{-\lambda t} \sum_{k=0}^{\infty} \frac{(Q + \lambda I)^k t^k}{k!} \\ &= e^{-\lambda t} e^{(Q + \lambda I)t} = e^{-\lambda t} e^{Qt} e^{\lambda t} = e^{Qt} \equiv \sum_{k=0}^{\infty} \frac{Q^k t^k}{k!}. \quad \blacksquare \end{aligned}$$

In §5 we will show how uniformization can be applied to quickly determine existence, uniqueness and the form of the limiting distribution of a CTMC. Now we consider a special class of CTMC's that both often arise and are easy to analyze.

4. Birth-and-Death Processes

Many CTMC's have transitions that only go to neighboring states, i.e., either up one or down one; they are called birth-and-death processes. Motivated by population models, a transition up one is called a **birth**, while a transition down one is called a **death**. The birth rate in state i is denoted by λ_i , while the death rate in state i is denoted by μ_i . The rate diagram for a birth-and-death process (with state space $\{0, 1, \dots, n\}$) takes the simple linear form shown in Figure 2.

Thus, for a birth-and-death process, the CTMC transition rates take the special form

$$Q_{i,i+1} = \lambda_i, \quad Q_{i,i-1} = \mu_i \quad \text{and} \quad Q_{i,j} = 0 \quad \text{if } j \notin \{i-1, i, i+1\}, \quad 1 \leq i \leq n-1, \quad (4.1)$$

with

$$Q_{0,1} = \lambda_0, \quad Q_{0,j} = 0 \quad \text{if } j \notin \{0, 1\}, \quad Q_{n,n-1} = \mu_n \quad \text{and} \quad Q_{n,j} = 0 \quad \text{if } j \notin \{n-1, n\}. \quad (4.2)$$

Rate Diagram for a Birth-and-Death Process

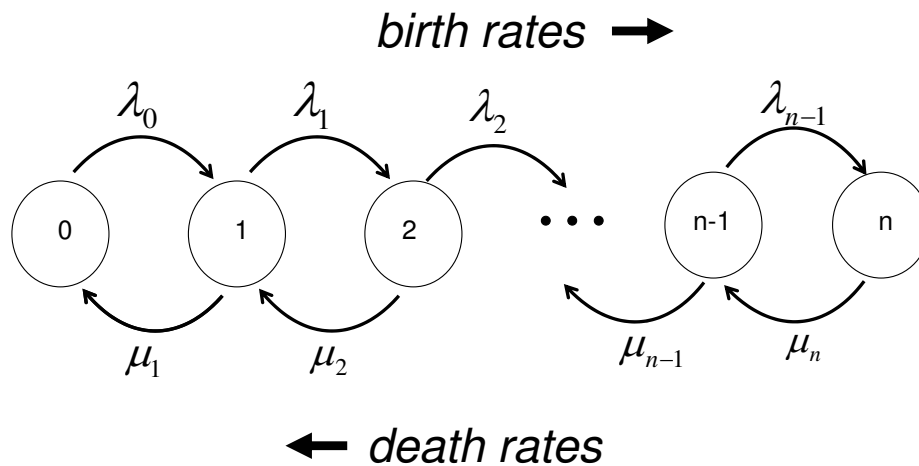


Figure 2: A rate diagram showing the transition rates for a birth-and-death process.

As before, the row sums of Q are zero.

A further special case is a **pure-birth process**, which only has transitions up one (equivalently, all death rates are 0). We have already encountered a special pure-birth process (on the nonnegative integers) - the Poisson process - which has constant birth rate, i.e., $\lambda_i = \lambda$ for all i . Similarly, a pure-death process has only transitions down one. For a finite state space, a pure-death process is equivalent to a pure-birth process, because we can just relabel the states.

The special structure of a birth-and-death process makes it easier to calculate the limiting probabilities. First, we observe that the global-balance equations (flow into state j equals flow out of state j), captured by the equation $\alpha Q = 0$, can be replaced by more elementary detailed-balance equations.

Theorem 4.1. (detailed-balance equations) *For a birth-and-death process, the limiting probability vector α is the unique solution to the **detailed-balance equations***

$$\alpha_j \lambda_j = \alpha_{j+1} \mu_{j+1} \quad \text{for all } j, \quad 0 \leq j \leq n-1, \quad (4.3)$$

with $\alpha e = 1$.

Proof. We give two different proofs: First, just as for a general CTMC, we can apply a **rate-conservation principle**, but now in a more special form. Because the birth-and-death process can move only to neighboring states, we can deduce that the steady-state rate of transitions up from state j to state $j+1$, $\alpha_j \lambda_j$, must equal the steady-state rate of transitions down from state $j+1$ to state j , $\alpha_{j+1} \mu_{j+1}$. That is, there is rate conservation between any two neighboring states. That yields the detailed-balance equations in (4.3). The rate-conservation

principle itself follows from a simple observation. In the time interval $[0, t]$, the number of transitions up from state j to state $j + 1$ can differ by at most by one from the number of transitions down from state $j + 1$ to state j .

From the perspective of a general CTMC, we already have established that α is the unique solution to $\alpha Q = 0$ with $\alpha \mathbf{e} = 1$. For a birth-and-death process, the j^{th} equation in the system $\alpha Q = 0$ is

$$(\alpha Q)_j = \alpha_{j-1}\lambda_{j-1} - \alpha_j(\lambda_j + \mu_j) + \alpha_{j+1}\mu_{j+1} = 0 \quad \text{for } 1 \leq j \leq n-1, \quad (4.4)$$

with

$$(\alpha Q)_0 = -\alpha_0\lambda_0 + \alpha_1\mu_1 = 0 \quad \text{and} \quad (\alpha Q)_n = \alpha_{n-1}\lambda_{n-1} - \alpha_n\mu_n = 0. \quad (4.5)$$

From (4.3)-(4.5), we see that the sum of the first j equations of the form (4.4) and (4.5) from $\alpha Q = 0$ coincides with the j^{th} detailed-balance equation in (4.3), while the difference between the j^{th} and $(j-1)^{\text{st}}$ detailed-balance equations coincides with the j^{th} equation from (4.4) and (4.5). Hence the two characterizations are equivalent. ■

In fact, it is not necessary to solve a system of linear equations each time we want to calculate the limiting probability vector α , because we can analytically solve the detailed-balance equations to produce an explicit formula.

Theorem 4.2. (limiting probabilities) *For a birth-and-death process with state space $\{0, 1, \dots, n\}$,*

$$\alpha_j = \frac{r_j}{\sum_{k=0}^n r_k} \quad 0 \leq j \leq n, \quad (4.6)$$

where

$$r_0 = 1 \quad \text{and} \quad r_j = \frac{\lambda_0 \times \lambda_1 \times \dots \times \lambda_{j-1}}{\mu_1 \times \mu_2 \times \dots \times \mu_j}. \quad (4.7)$$

Proof and Discussion. By virtue of Theorem 4.1, it suffices to solve the detailed-balance equations in (4.3). We can do that **recursively**:

$$\alpha_j = \frac{\lambda_{j-1}}{\mu_j} \alpha_{j-1} \quad \text{for all } j \geq 1,$$

which implies that

$$\alpha_j = r_j \alpha_0 \quad \text{for all } j \geq 1,$$

for r_j in (4.7). We obtain the final form (4.6) when we require that $\alpha \mathbf{e} = 1$. ■

Example 4.1. (a small barbershop) Consider a small barbershop, where there are only two barbers, each with his own barber chair. Suppose that there is only room for at most 5 customers, with 2 in service and 3 waiting. Assume that potential customers arrive according to a Poisson process at rate $\lambda = 6$ per hour. Customers arriving when the system is full are **blocked and lost**, leaving without receiving service and without affecting future arrivals. Assume that the duration of each haircut is an independent exponential random variable with a mean of $\mu^{-1} = 15$ minutes. Customers are served in a first-come first-served manner by the first available barber.

We can ask a variety of **questions**: (a) What is the long-run proportion of time there are two customers in service plus two customers waiting? (b) What is the (long-run) proportion of time each barber is busy? We might then go on to ask how this long-run behavior would change if we changed the number of barbers or the number of waiting spaces.

We start by **constructing the model**. Let $Q(t)$ denote the number of customers in the system at time t . Then the stochastic process $\{Q(t) : t \geq 0\}$ is a birth-and-death process with six states: 0, 1, 2, 3, 4, 5. Indeed, this is a standard **queueing model**, commonly referred to as **the M/M/2/3 queue**. The first M means a Poisson arrival process (M for Markov), the second M means IID exponential service times (M again for Markov), the 2 is for 2 servers, and the 3 is for 3 additional waiting spaces.) It is common to use λ to denote the arrival rate and μ the service rate of each server.

We can represent the CTMC in terms of competing exponential timers, as in §3.3. The possible triggering events are an arrival (birth), causing the state to go up 1, or a departure (death), causing the state to go down 1. It is of course important that these are independent exponential random variables.

The blocking alters the arrival process. The blocking means that no arrivals can enter in state 5. By making the state space $\{0, 1, \dots, 5\}$, we have accounted for the blocking. Since the interarrival times of a Poisson process have an exponential distribution, there are active clocks with exponential timers corresponding to the event of a new arrival in states 0-4. The arrival clock in state i has mean $1/\lambda_i = 1/\lambda$, where $\lambda = 6$ per hour is the arrival rate of the Poisson process. Hence the birth rates are $\lambda_i = 6$, $0 \leq i \leq 4$. We have $\lambda_5 = 0$, because there are no arrivals when the system is full.

Since the service times are independent exponential random variables, the active clocks corresponding to departures also can be represented as exponential random variables. (Recall that the minimum of independent exponential variables is again exponential with a rate equal to the sum of the rates.) There are active clocks with exponential timers corresponding to the event of a new departure in states 1-5. The departure clock in state i has mean $1/\mu_i$, where μ_i is the death rate to be determined. Since the mean service time is $1/\mu = 15$ minutes, the service rate for each barber is $\mu = 1/15$ per minute or $\mu = 4$ per hour. However, we must remember that the service rate applies to each server separately. Since we are measuring time in hours, the death rates are $\mu_1 = \mu = 4$, $\mu_2 = \mu_3 = \mu_4 = \mu_5 = 2\mu = 8$. We have $\mu_0 = 0$ since there can be no departures from an empty system.

Given the birth rates and death rates just defined, we can draw the rate diagram for the six states 0, 1, ..., 5, as in Figure 2. The associated rate matrix is now

$$Q = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} -6 & 6 & 0 & 0 & 0 & 0 \\ 4 & -10 & 6 & 0 & 0 & 0 \\ 0 & 8 & -14 & 6 & 0 & 0 \\ 0 & 0 & 8 & -14 & 6 & 0 \\ 0 & 0 & 0 & 8 & -14 & 6 \\ 0 & 0 & 0 & 0 & 8 & -8 \end{pmatrix} \end{matrix}$$

We now can apply Theorem 4.2 to calculate the **limiting probabilities**. From (4.6) and (4.7),

$$\alpha_j = \frac{r_j}{r_0 + r_1 + \dots + r_5}, \quad 0 \leq j \leq 5,$$

where $r_0 = 1$ and

$$r_j = \frac{\lambda_0 \times \lambda_1 \times \dots \times \lambda_{j-1}}{\mu_1 \times \mu_2 \times \dots \times \mu_j}, \quad 1 \leq j \leq 5.$$

Here

$$\begin{aligned} r_0 &= 1 = \frac{512}{512} \\ r_1 &= \frac{6}{4} = \frac{3}{2} = \frac{768}{512} \end{aligned}$$

$$\begin{aligned}
r_2 &= \frac{6 \times 6}{4 \times 8} = \frac{36}{32} = \frac{9}{8} = \frac{576}{512} \\
r_3 &= \frac{6 \times 6 \times 6}{4 \times 8 \times 8} = \frac{27}{32} = \frac{432}{512} \\
r_4 &= \frac{27 \times 6}{32 \times 8} = \frac{81}{128} = \frac{324}{512} \\
r_5 &= \frac{81 \times 6}{128 \times 8} = \frac{243}{512}
\end{aligned}$$

Hence,

$$\begin{aligned}
\alpha_0 &= \frac{512}{2855}, & \alpha_1 &= \frac{768}{2855}, & \alpha_2 &= \frac{576}{2855}, & \alpha_3 &= \frac{432}{2855}, \\
\alpha_4 &= \frac{324}{2855} & \text{and} & & \alpha_5 &= \frac{243}{2855}.
\end{aligned}$$

This particular calculation is admittedly a bit tedious, but it is much better than solving the system of equations based on $\alpha Q = 0$, which would be required for a general CTMC.

Given the steady-state probability vector $\alpha \equiv (\alpha_0, \dots, \alpha_5)$, you can then **answer the questions posed**: (a) The long-run proportion of time there are two customers in service plus two customers waiting is α_4 . (b) The (long-run) proportion of time each barber is busy is $(1/2)\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 + \alpha_5$. (The two barbers are each busy half of the time that only one barber is busy.) Finally, we could see how these answers would change if we changed the number of barbers or the number of waiting spaces. We would just perform a similar analysis with the alternative model(s). ■

Example 4.2. (customers that balk and abandon) Consider the same barbershop with two barbers and three waiting spaces, as specified above, but in addition suppose that customers may elect to balk or abandon. In particular, suppose that an arriving customer finding both barbers busy, but an available waiting space, will elect to stay, independently of all past events, with probability $2/3$; otherwise, the arrival will **balk**, i.e., refuse to join and instead immediately leave, without affecting future arrivals. Moreover, suppose that each arriving customer who is not blocked and who elects to wait is only willing to wait a certain time before starting service; otherwise the customer will **abandon**, i.e., leave without receiving service and without affecting future arrivals. Let the amount of patience of successive customers, i.e., these successive times to abandon, be IID exponential random variables with mean $\theta^{-1} = 10$ minutes.

Again we can ask a variety of **questions**: (a) What is the rate of customer abandonment? (b) What is the long-run proportion of potential arrivals that enter and then abandon? (c) What proportion of potential customers enter upon arrival (i.e., neither balk nor are blocked)? (d) What proportion of potential customers are served?

Even though it may not be entirely evident at first, the stochastic process representing the number of customers in the system over time is again a birth-and-death process. Again we can represent the CTMC in terms of competing exponential timers, as in §3.3. The possible triggering events are an arrival (birth), causing the state to go up 1, or a departure (death), causing the state to go down 1, where the departure may be due to service completion or abandonment. As noted before, the blocking means that no arrivals can enter in state 5. The balking alters the arrival process further. The balking corresponds to performing an independent thinning of the external Poisson arrival process in states 2-4. In those states, the actual arrivals form a Poisson process with arrival rate $\lambda \times (2/3) = 6 \times (2/3) = 4$ per hour.

Since the interarrival times of a Poisson process have an exponential distribution, there are active clocks with exponential timers corresponding to the event of a new arrival in states 0-4. The arrival clock in state i has mean $1/\lambda_i$, where λ_i is the birth rate to be determined. The birth rates in these states are: $\lambda_0 = \lambda_1 = \lambda = 6$ per hour and $\lambda_2 = \lambda_3 = \lambda_4 = 6 \times (2/3) = 4$ per hour. (The reduction is due to the balking. We have $\lambda_5 = 0$, because there are no arrivals when the system is full.)

Since the service times and times to abandon are independent exponential random variables, the active clocks corresponding to departures also can be represented as exponential random variables. As before, there are active clocks with exponential timers corresponding to the event of a new departure in states 1-5. The departure clock in state i has mean $1/\mu_i$, where μ_i is the death rate to be determined. As before, the service rate for each barber is $\mu = 4$ per hour. Since the mean time to abandon is $1/\theta = 10$ minutes for each customer, the individual abandonment rate is $\theta = 1/10$ per minute or 6 per hour. However, we must remember that the service rate applies to each server separately, while the abandonment rate applies to each waiting customer separately. Thus the death rates are $\mu_1 = \mu = 4$, $\mu_2 = 2\mu = 8$, $\mu_3 = 2\mu + \theta = 8 + 6 = 14$, $\mu_4 = 2\mu + 2\theta = 8 + 12 = 20$, $\mu_5 = 2\mu + 3\theta = 8 + 18 = 26$. ($\mu_0 = 0$.)

Given the new birth rates and death rates just defined, we can draw the new rate diagram for the six states $0, 1, \dots, 5$, as in Figure 2. The new rate matrix is

$$Q = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} -6 & 6 & 0 & 0 & 0 & 0 \\ 4 & -10 & 6 & 0 & 0 & 0 \\ 0 & 8 & -12 & 4 & 0 & 0 \\ 0 & 0 & 14 & -18 & 4 & 0 \\ 0 & 0 & 0 & 20 & -24 & 4 \\ 0 & 0 & 0 & 0 & 26 & -26 \end{pmatrix} \end{matrix}$$

We now can apply Theorem 4.2 to calculate the **limiting probabilities**. From (4.6),

$$\alpha_i = \frac{r_i}{r_0 + r_1 + \dots + r_5}, \quad 0 \leq i \leq 5,$$

where $r_0 = 1$ and

$$r_i = \frac{\lambda_0 \times \lambda_1 \times \dots \times \lambda_{i-1}}{\mu_1 \times \mu_2 \times \dots \times \mu_i}, \quad 1 \leq i \leq 5.$$

Here

$$\begin{aligned} r_0 &= 1 = \frac{3640}{3640} \\ r_1 &= \frac{6}{4} = \frac{3}{2} = \frac{5460}{3640} \\ r_2 &= \frac{6 \times 6}{4 \times 8} = \frac{36}{32} = \frac{9}{8} = \frac{4095}{3640} \\ r_3 &= \frac{6 \times 6 \times 4}{4 \times 8 \times 14} = \frac{36}{112} = \frac{9}{28} = \frac{1170}{3640} \\ r_4 &= \frac{9 \times 4}{28 \times 20} = \frac{9}{140} = \frac{234}{3640} \\ r_5 &= \frac{9 \times 4}{140 \times 26} = \frac{18}{1820} = \frac{9}{910} = \frac{36}{3640} \end{aligned}$$

Hence,

$$\begin{aligned} \alpha_0 &= \frac{3640}{14635}, & \alpha_1 &= \frac{5460}{14635}, & \alpha_2 &= \frac{4095}{14635}, & \alpha_3 &= \frac{1170}{14635}, \\ \alpha_4 &= \frac{234}{14635} & \text{and} & & \alpha_5 &= \frac{36}{14635}. \end{aligned}$$

Given the steady-state probability vector $\alpha \equiv (\alpha_0, \dots, \alpha_5)$, you can then **answer the other questions**: (a) The rate of customer abandonments is $\theta\alpha_3 + 2\theta\alpha_4 + 3\theta\alpha_5 = 6\alpha_3 + 12\alpha_4 + 18\alpha_5$. (b) The long-run proportion of potential customers that enter and abandon is the rate customers abandon, just determined, divided by the arrival rate, i.e.,

$$\frac{(\theta\alpha_3 + 2\theta\alpha_4 + 3\theta\alpha_5)}{\lambda} = \frac{(6\alpha_3 + 12\alpha_4 + 18\alpha_5)}{6} = \alpha_3 + 2\alpha_4 + 3\alpha_5 .$$

Questions (c) and (d) are more tricky, because they ask about the proportion of customers having a specified experience, instead of the long-run proportion of time. However, it turns out that these notions agree in this problem, because the arrival process of potential customers is a Poisson process. There is a principle called **Poisson Arrivals See Time Averages (PASTA)** that implies that the proportion of customers that see some state upon arrival coincides with the proportion of time the process spends in that state, provided that the arrival process is Poisson (and other regularity conditions hold, which they do here; e.g., see §5.16 of Wolff (1989), Melamed and Whitt (1990) or Stidham and El Taha (1999)). Hence, consistent with intuition, the long-run proportion of all potential customers that are blocked coincides with the long-run proportion of time that the system is full, which is α_5 . (But that property would not remain true if we made the arrival process non-Poisson.) Similarly, the long-run proportion of customers that balk is $1/3$ times the long-run proportion of time that the system is in one of the states 2, 3 or 4, which is $(1/3) \times (\alpha_2 + \alpha_3 + \alpha_4)$. (c) Hence the long-run proportion of potential customers enter upon arrival (i.e., neither balk nor are blocked) is $1 - (1/3) \times (\alpha_2 + \alpha_3 + \alpha_4) - \alpha_5$.

(d) We can find the long-run proportion of potential customers served in two different ways:

Method 1. The long-run proportion of customers served is 1 minus the sum of the proportions that balk, abandon and are blocked. We can thus apply the answers to previous questions. The answer is

$$1 - \frac{\alpha_2 + \alpha_3 + \alpha_4}{3} - \alpha_5 - (\alpha_3 + 2\alpha_4 + 3\alpha_5) ,$$

Rewriting, we get

$$1 - \frac{\alpha_2}{3} - \frac{4\alpha_3}{3} - \frac{7\alpha_4}{3} - 4\alpha_5 = \frac{11,020}{14,635} = 0.753 .$$

Method 2. The long-run proportion of customers served can be represented as the overall service completion rate divided by the external arrival rate, counting all potential arrivals. The denominator - the arrival rate - is $\lambda = 6$ per hour. The service completion rate is

$$(\alpha_1 \times 4) + (\alpha_2 + \alpha_3 + \alpha_4 + \alpha_5) \times 8 ,$$

because the service rate is 4 in state 1, while the service rate is $2 \times 4 = 8$ in states 2 – 5. Hence, the long-run proportion of customers served is

$$\frac{2\alpha_1}{3} + \frac{4(\alpha_2 + \alpha_3 + \alpha_4 + \alpha_5)}{3} = \frac{11,020}{14,635} .$$

Even though the two formulas are different, they give the same answer.

Finally, just as before, we could see how these answers would change if we changed the number of barbers or the number of waiting spaces. We would just perform a similar analysis with the alternative model(s). ■

In many applications it is natural to use birth-and-death processes with **infinite state spaces**. As in other mathematical settings, we primarily introduce infinity because it is more convenient. With birth-and-death processes, an infinite state space often simplifies the form of the limiting probability distribution. We illustrate by giving a classic queueing example.

Example 4.3. (the M/M/1 Queue) One of the most elementary queueing models is the $M/M/1$ queue, which has a single server and unlimited waiting room. As with the $M/M/s/r$ model considered in Example 4.1 (with $S = 2$ and $r = 3$), customers arrive in a Poisson process with rate λ and the service times are IID exponential random variables with mean $1/\mu$. The number of customers in the system at time t as a function of t , say $Q(t)$, is then a birth-and-death process. However, since there is unlimited waiting room, the state space is infinite.

With an infinite state space, we must guard against pathologies; see §10. In order to have a proper stationary distribution, it is necessary to require that the arrival rate λ be less than the maximum possible rate out, μ . Equivalently, we require that the traffic intensity $\rho \equiv \lambda/\mu$ be strictly less than 1.

When we apply the extension of Theorem 4.2 to infinite state spaces, under the assumption that $\rho < 1$, we get

$$\alpha_j = \frac{r_j}{\sum_{k=0}^{\infty} r_k}, \quad \text{where } r_j = \rho^j, \quad j \geq 0, \quad (4.8)$$

which implies that α is the **geometric distribution**; i.e.,

$$\lim_{t \rightarrow \infty} P(Q(t) = j | Q(0) = i) = \alpha_j = (1 - \rho)\rho^j, \quad j \geq 0, \quad (4.9)$$

which has mean $\rho/(1 - \rho)$.

If instead we considered the $M/M/1/r$ model (the $M/M/1$ model with a finite waiting room), which has 1 server and r extra waiting spaces, then the birth-and-death process has $r+2$ states, from 0 to $r+1$. The limiting distribution then becomes the **truncated geometric distribution**:

$$\alpha_j = \frac{(1 - \rho)\rho^j}{(1 - \rho^{r+2})}, \quad 0 \leq j \leq r + 1. \quad (4.10)$$

The geometric distribution in (4.9) is more appealing than the truncated geometric distribution in (4.10) because of its cleaner form. However, the finite-waiting room model applies without constraint on ρ ; a proper limiting distribution exists for $\rho \geq 1$ as well as for $\rho < 1$. ■

5. Stationary and Limiting Probabilities for CTMC's

Just as with DTMC's, the CTMC model specifies how the process moves locally. Just as with DTMC's, we use the CTMC model to go from the assumed local behavior to deduce global behavior. That is, we use the CTMC model to calculate its limiting probability distribution, as defined in (2.3). We then use that limiting probability distribution to answer questions about what happens in the long run. In this section we show how to compute limiting probabilities. The examples will illustrate how to apply the limiting distribution to answer other questions about what happens in the long run.

But first we want to establish a firm foundation. We will demonstrate existence and uniqueness of a limiting distribution, which justifies talking about “the” limiting distribution of an (irreducible) CTMC. We also want to show that the limiting distribution of a CTMC coincides with the (unique) stationary distribution of the CTMC. A probability vector β is a **stationary distribution** for a CTMC $\{X(t) : t \geq 0\}$ if $P(X(t) = j) = \beta_j$ for all t and j

whenever $P(X(0) = j) = \beta_j$ for all j . In general the two notions - limiting distribution and stationary distribution - are distinct, but for CTMC's there is a unique probability vector with both properties.

Example 5.1. (distinction between the concepts) Before establishing positive results for CTMC's, we show that in general the two notions are distinct: there are stationary distributions that are not limiting distributions; and there are limiting distributions that are not stationary distributions.

(a) Recall that a periodic irreducible finite-state DTMC has a unique stationary probability vector, which is not a limiting probability vector; the transitions probabilities $P_{i,j}^k$ alternate as k increases, assuming a positive value at most every d steps, where d is the period of the chain. (A CTMC cannot be periodic.)

(b) To go the other way, consider a stochastic process $\{X(t) : t \geq 0\}$ with continuous state space consisting of the unit interval $[0, 1]$. Suppose that the stochastic process moves deterministically except for its initial value $X(0)$, which is a random variable taking values in $[0, 1]$. After that initial random start, let the process move deterministically on the unit interval $[0, 1]$ according to the following rules: From state 0, let the process instantaneously jump to state 1. Otherwise, let the process move according to the ODE

$$X'(t) \equiv \frac{dX(t)}{dt} = -X(t), t \geq 0 .$$

Then $\{X(t) : t \geq 0\}$ is a Markov process with a unique limiting distribution. In particular,

$$\lim_{t \rightarrow \infty} X(t) = 0 \quad \text{with probability } 1 ,$$

so that the limiting distribution is unit probability mass on 0. However, that limit distribution is not a stationary distribution. Indeed, $P(X(t) = 0) = 0$ for all $t > 0$ and all distributions of $X(0)$. If $P(X(0) = 0) = 1$, then $P(X(t) = e^{-t})$ for all $t = 1$. Even though this Markov process has a unique limiting probability distribution, there is no stationary probability vector for this Markov process. ■

But the story is very nice for irreducible finite-state CTMC's: Then there always exists a unique stationary probability vector, which also is a limiting probability vector. The situation is somewhat cleaner for CTMC's than for DTMC's, because we cannot have periodic CTMC's. That is implied by the following result.

Lemma 5.1. (positive transition probabilities) *For an irreducible CTMC, $P_{i,j}(t) > 0$ for all i, j and $t > 0$.*

Proof. The argument going forward in time is easy: By Lemma 2.1, if $P_{i,j}(s) > 0$, then

$$P_{i,j}(s+t) = \sum_k P_{i,k}(s)P_{k,j}(t) \geq P_{i,j}(s)P_{j,j}(t) \geq P_{i,j}(s)e^{Q_{j,j}t} > 0 \quad \text{for all } t > 0 ,$$

because $P_{j,j}(t)$ is bounded below by the probability of no transition at all from state j in time t , which is $e^{Q_{j,j}t}$. (Recall that $Q_{j,j} < 0$.) More generally, we apply representation (3.30). Since the CTMC is irreducible, $P_{i,j}(t) > 0$ for some t . By representation (3.30), we thus have $\tilde{P}_{i,j}^k > 0$ for some k , implying that the embedded DTMC with transition matrix \tilde{P} is irreducible. From here on, we argue by contradiction: Suppose that $P_{i,j}(t) = 0$ for some $t > 0$. Then, by representation (3.30), $\tilde{P}_{i,j}^k = 0$ for all k , which would imply that \tilde{P} is reducible. Since that is a contradiction, we must have $P_{i,j}(t) > 0$ for all $t > 0$, as claimed. ■

Theorem 5.1. (existence and uniqueness) For an irreducible finite-state CTMC, there exists a unique **limiting probability vector** α ; i.e., there exists a unique probability vector α such that

$$\lim_{t \rightarrow \infty} P_{i,j}(t) = \alpha_j \quad \text{for all } i \quad \text{and } j . \quad (5.1)$$

Moreover, that limiting probability vector α is the unique **stationary probability vector**, i.e., if

$$P(X(0) = j) = \alpha_j \quad \text{for all } j ,$$

then

$$P(X(t) = j) = \alpha_j \quad \text{for all } j \quad \text{and } t > 0 . \quad (5.2)$$

Proof. We will apply established results for DTMC's in the setting of the fourth modelling approach in §3.4; i.e., we will apply uniformization. To do so, we apply representation (3.30). From that representation and Lemma 5.1, it follows immediately that the CTMC is irreducible if and only if the embedded Markov chain with transition matrix \tilde{P} is irreducible. Assuming that the CTMC is indeed irreducible, the same is true for that embedded DTMC. By making λ in (3.26) larger if necessary, we can have $\tilde{P}_{i,i} > 0$ for all i , so that the embedded DTMC with transition matrix \tilde{P} can be taken to be aperiodic as well.

Given that the DTMC with transition matrix \tilde{P} is irreducible and aperiodic, we know that the embedded DTMC has a unique stationary distribution $\tilde{\pi}$ satisfying

$$\tilde{\pi} = \tilde{\pi}\tilde{P} \quad \text{and} \quad \tilde{\pi}\mathbf{e} = 1 ,$$

with the additional property that

$$\tilde{P}_{i,j}^k \rightarrow \tilde{\pi}_j \quad \text{as } k \rightarrow \infty$$

for all i and j . From representation (3.30), it thus follows that $\tilde{\pi}$ is also the limiting distribution for the CTMC; i.e., we have

$$\alpha_j = \tilde{\pi}_j \quad \text{for all } j .$$

Here is a **detailed mathematical argument**: For any $\epsilon > 0$ given, first choose k_0 such that $|\tilde{P}_{i,j}^k - \tilde{\pi}_j| < \epsilon/2$ for all $k \geq k_0$. Then choose t_0 such that $P(N(t) < k_0) < \epsilon/4$ for all $t \geq t_0$. As a consequence, for $t > t_0$,

$$\begin{aligned} |P_{i,j}(t) - \tilde{\pi}_j| &= |P(X(t) = j | X(0) = i, N(t) < k_0) - \tilde{\pi}_j| P(N(t) < k_0) \\ &\quad + |P(X(t) = j | X(0) = i, N(t) \geq k_0) - \tilde{\pi}_j| P(N(t) \geq k_0) \\ &\leq 2P(N(t) < k_0) + P(N(t) \geq k_0) \frac{\epsilon}{2} \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} \leq \epsilon . \end{aligned} \quad (5.3)$$

Moreover, there can be no other stationary distribution, because any stationary distribution of the CTMC has to coincide with the limiting distribution of the DTMC, again by (3.30). ■

We now turn to calculation. We give three different ways to calculate the limiting distribution, based on the different modelling frameworks. (We do not give a separate treatment for the competing clocks with exponential timers. We treat that case via the transition rates.) To sum row vectors in matrix notation, we right-multiply by a column vector of 1's. Let \mathbf{e} denote such a column vector of 1's.

Theorem 5.2. (calculation)

(a) Given a CTMC characterized as a DTMC with one-step transition matrix \tilde{P} and transitions according to a Poisson process with rate λ , as in §3.4,

$$\alpha_j = \tilde{\pi}_j \quad \text{for all } j, \quad (5.4)$$

where $\tilde{\pi}$ is the unique solution to

$$\tilde{\pi} = \tilde{\pi}\tilde{P} \quad \text{and} \quad \tilde{\pi}\mathbf{e} = 1, \quad (5.5)$$

with \tilde{P} given in (3.27) or, equivalently,

$$\sum_i \tilde{\pi}_i \tilde{P}_{i,j} = \tilde{\pi}_j \quad \text{for all } j \quad \text{and} \quad \sum_j \tilde{\pi}_j = 1. \quad (5.6)$$

(b) Given a CTMC characterized in terms of a DTMC with one-step transition matrix P and exponential transition times with means $1/\nu_i$, as in §3.1,

$$\alpha_j = \frac{(\pi_j/\nu_j)}{\sum_k (\pi_k/\nu_k)}, \quad (5.7)$$

where π is the unique solution to

$$\pi = \pi P \quad \text{and} \quad \pi\mathbf{e} = 1. \quad (5.8)$$

(c) Given a CTMC characterized by its transition-rate matrix Q , as in §3.2, α is the unique solution to

$$\alpha Q = 0 \quad \text{and} \quad \alpha\mathbf{e} = 1 \quad (5.9)$$

or, equivalently,

$$\sum_i \alpha_i Q_{i,j} = 0 \quad \text{for all } j \quad \text{and} \quad \sum_i \alpha_i = 1. \quad (5.10)$$

(d) Given a CTMC characterized by its transition function $P(t)$, perhaps as constructed in §3.4, α is the unique solution to

$$\alpha P(t) = \alpha \quad \text{for any } t > 0 \quad \text{and} \quad \alpha\mathbf{e} = 1 \quad (5.11)$$

or, equivalently,

$$\sum_i \alpha_i P_{i,j}(t) = \alpha_j \quad \text{for all } j \quad \text{and} \quad \sum_i \alpha_i = 1. \quad (5.12)$$

Proof and Discussion. (a) Exploiting Uniformization. In our proof of Theorem 5.1 above, we have already shown that α coincides with $\tilde{\pi}$.

(b) Starting with the embedded DTMC. Since Theorem 5.1 establishes the existence of a unique stationary probability distribution, it suffices to show that the distribution displayed in (5.7) is that stationary distribution. Equivalently, it suffices to show that $\tilde{\pi} = \alpha$ for α in (5.7), where $\tilde{\pi}$ is the unique solution to

$$\tilde{\pi} = \tilde{\pi}\tilde{P} \quad \text{and} \quad \tilde{\pi}\mathbf{e} = 1.$$

To see that is the case, observe that $\alpha_j = c\pi_j/\nu_j$ for α defined in (5.7). To show that $\alpha\tilde{P} = \alpha$, observe that

$$\begin{aligned}
(\alpha\tilde{P})_j &= \sum_i \alpha_i \tilde{P}_{i,j} = c \sum_i \frac{\pi_i}{\nu_i} \tilde{P}_{i,j} \\
&= c \left(\sum_{i,i \neq j} \frac{\pi_i}{\nu_i} \frac{Q_{i,j}}{\lambda} + \frac{\pi_j}{\nu_j} \tilde{P}_{j,j} \right) \\
&= c \left(\sum_{i,i \neq j} \frac{\pi_i}{\nu_i} \left(\frac{\nu_i}{\lambda} P_{i,j} \right) + \frac{\pi_j}{\nu_j} \left(1 - \frac{\sum_{i,i \neq j} \nu_j P_{j,i}}{\lambda} \right) \right) \\
&= c \left(\sum_{i,i \neq j} \frac{\pi_i P_{i,j}}{\lambda} + \frac{\pi_j}{\nu_j} - \frac{\pi_j \sum_{i,i \neq j} P_{j,i}}{\lambda} \right) \\
&= c \left(\sum_i \frac{\pi_i P_{i,j}}{\lambda} - \frac{\pi_j P_{j,j}}{\lambda} + \frac{\pi_j}{\nu_j} - \frac{\pi_j (1 - P_{j,j})}{\lambda} \right) \\
&= c \left(\frac{\pi_j}{\lambda} - \frac{\pi_j P_{j,j}}{\lambda} + \frac{\pi_j}{\nu_j} - \frac{\pi_j}{\lambda} + \frac{\pi_j P_{j,j}}{\lambda} \right) \\
&= c \frac{\pi_j}{\nu_j} = \alpha_j .
\end{aligned} \tag{5.13}$$

From (5.7), we see that $\alpha\mathbf{e} = 1$, where \mathbf{e} is again a column vector of 1's. That completes the proof.

We now give a **separate direct informal argument** (which can be made rigorous) to show that α has the claimed form. Let $Z_{i,j}$ be the time spend in state i during the j^{th} visit to state i and let $N_i(n)$ be the number of visits to state i among the first n transitions. Then the actual proportion of time spent in state i during the first n transitions, say $T_i(n)$, is

$$T_i(n) = \frac{\sum_{j=1}^{N_i(n)} Z_{i,j}}{\sum_k \sum_{j=1}^{N_k(n)} Z_{k,j}} , \tag{5.14}$$

However, by properties of DTMC's, $n^{-1}N_i(n) \rightarrow \pi_i$ with probability 1 as $n \rightarrow \infty$. Moreover, by the law of large numbers,

$$\frac{1}{n} \sum_{j=1}^{N_i(n)} Z_{i,j} = \left(\frac{N_i(n)}{n} \right) \left(\frac{\sum_{j=1}^{N_i(n)} Z_{i,j}}{N_i(n)} \right) \rightarrow \pi_i E[Z_{i,j}] = \pi_i/\nu_i \quad \text{as } n \rightarrow \infty . \tag{5.15}$$

Thus, combining (5.14) and (5.15), we obtain

$$T_i(n) \rightarrow \frac{(\pi_i/\nu_i)}{\sum_k (\pi_k/\nu_k)} \quad \text{as } n \rightarrow \infty , \tag{5.16}$$

supporting (5.7). For a full proof, we need to show that this same limit holds at arbitrary times t as $t \rightarrow \infty$. It is intuitively clear that holds, but we do not prove that directly.

(c) Starting with the transition rates. We give several different arguments, from different perspectives, to show that α is characterized as the unique solution of $\alpha Q = 0$ with $\alpha\mathbf{e} = 1$.

We first apply a **rate-conservation principle**: In steady state, the rate of transitions into state j has to equal the rate of transitions out of state j , for each state j . The steady-state rate of transitions into state j is

$$\sum_{i, i \neq j} \alpha_i Q_{i,j}$$

for the limiting probability vector α to be determined, while the steady-state rate of transitions out of state j is

$$\sum_{i, i \neq j} \alpha_j Q_{j,i} = -\alpha_j Q_{j,j} .$$

Setting these two steady-state rates equal yields

$$\sum_i \alpha_i Q_{i,j} = 0 ,$$

which, when applied to all j , is equivalent to $\alpha Q = 0$ in matrix notation.

Alternatively, we can **start from the ODE's**. From Theorem 5.1, we know that $P_{i,j}(t) \rightarrow \alpha_j$ as $t \rightarrow \infty$ for all i and j . Thus the right side of the backwards ODE $P'(t) = QP(t)$ converges, which implies that

$$P'_{i,j}(t) = \sum_k Q_{i,k} P_{k,j}(t) \rightarrow \sum_k Q_{i,k} \alpha_j \quad \text{as } t \rightarrow \infty .$$

However, since $\sum_k Q_{i,k} = 0$ for all i ,

$$P'_{i,j}(t) \rightarrow 0 \quad \text{as } t \rightarrow \infty \quad \text{for all } i \quad \text{and } j .$$

When we apply these established limits for $P(t)$ and $P'(t)$ in the **forward ODE**, $P'(t) = P(t)Q$, we immediately obtain the desired $0 = \alpha Q$, where $\alpha \mathbf{e} = 1$.

We can instead work with the **DTMC transition matrix** \tilde{P} . From (3.20) and (3.27), we see that

$$Q = \lambda(\tilde{P} - I) . \tag{5.17}$$

Multiply on the left by α in (5.17) to get

$$\alpha Q = \lambda(\alpha \tilde{P} - \alpha) ,$$

which implies that $\alpha Q = 0$ if and only if $\alpha \tilde{P} = \alpha$.

(d) Starting with the transition function $P(t)$. This final characterization is similar to part (a). Apply the explicit expression for $P(t)$ in (3.30) with the expression $\tilde{P} = \lambda^{-1}Q + I$ to deduce that $\alpha Q = 0$ if and only if $\alpha P(t) = \alpha$. ■

To illustrate, we now return once more to Examples 3.1 and 3.2.

Example 5.2. (Pooh Bear and the Three Honey Trees Revisited) In Example 3.1 the CTMC was naturally formulated as a DTMC with exponential transition times, as in the first modelling approach in §3.1. We exhibited the DTMC transition matrix P and the mean transition times $1/\nu_i$ before. Thus it is natural to apply Theorem 5.2 (b) in order to calculate the limiting probabilities. From that perspective, the limiting probabilities are

$$\alpha_j = \frac{\pi_j(1/\nu_j)}{\sum_k \pi_k(1/\nu_k)} ,$$

where the limiting probability vector π of the discrete-time Markov chain with transition matrix P is obtained by solving $\pi = \pi P$ with $\pi \mathbf{e} = 1$, yielding

$$\pi = \left(\frac{8}{17}, \frac{4}{17}, \frac{5}{17} \right).$$

Then final steady-state distribution, accounting for the random holding times, is

$$\alpha = \left(\frac{1}{2}, \frac{1}{4}, \frac{1}{4} \right).$$

We were then asked to find the limiting proportion of time that Pooh spends at each of the three trees. Those limiting proportions coincide with the limiting probabilities. That can be demonstrated by applying the renewal-reward theorem from renewal theory. ■

Example 5.3. (Copier Maintenance Revisited Again) Let us return to Example 3.2 and consider the question posed there: What is the long-run proportion of time that each copier is working and what is the long-run proportion of time that the repairman is busy? To have concrete numbers, suppose that the failure rates are $\gamma_1 = 1$ per month and $\gamma_2 = 3$ per month; and suppose the repair rates are $\beta_1 = 2$ per month and $\beta_2 = 4$ per month.

We first substitute the specified numbers for the rates γ_i and β_i in the rate matrix Q in (3.2), obtaining

$$Q = \begin{matrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ (1,2) & & & & & \\ (2,1) & & & & & \end{matrix} \begin{pmatrix} 0 & -4 & 1 & 3 & 0 & 0 \\ 1 & 2 & -5 & 0 & 3 & 0 \\ 2 & 4 & 0 & -5 & 0 & 1 \\ 0 & 0 & 2 & -2 & 0 & \\ 0 & 4 & 0 & 0 & -4 & \end{pmatrix}.$$

Then we solve the system of linear equations $\alpha Q = 0$ with $\alpha \mathbf{e} = 1$, which is easy to do with a computer and is not too hard by hand. Just as with DTMC's, one of the equations in $\alpha Q = 0$ is redundant, so that with the extra added equation $\alpha \mathbf{e} = 1$, there is a unique solution. Performing the calculation, we see that the limiting probability vector is

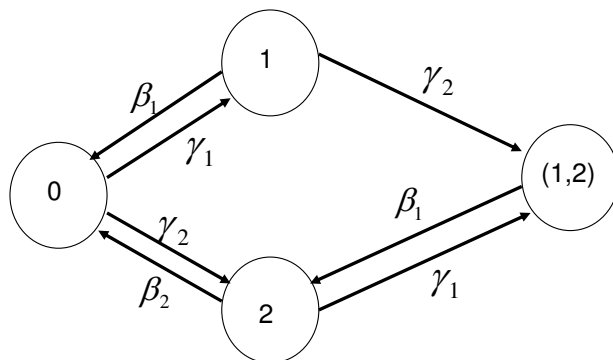
$$\alpha \equiv (\alpha_0, \alpha_1, \alpha_2, \alpha_{(1,2)}, \alpha_{(2,1)}) = \left(\frac{44}{129}, \frac{16}{129}, \frac{36}{129}, \frac{24}{129}, \frac{9}{129} \right).$$

Thus, the long-run proportion of time that copier 1 is working is $\alpha_0 + \alpha_2 = 80/129 \approx 0.62$, while the long-run proportion of time that copier 2 is working is $\alpha_0 + \alpha_1 = 60/129 \approx 0.47$. The long-run proportion of time that the repairman is busy is $\alpha_1 + \alpha_2 + \alpha_{(1,2)} + \alpha_{(2,1)} = 1 - \alpha_0 = 85/129 \approx 0.659$,

Now let us consider **an alternative repair strategy**: Suppose that copier 1 is more important than copier 2, so that it is more important to have it working. Toward that end, suppose the repairman always work on copier 1 when both copiers are down. In particular, now suppose that the repairman stops working on copier 2 when it is down if copier 1 also subsequently fails, and immediately shifts his attention to copier 1, returning to work on copier 2 after copier 1 has been repaired. How do the long-run proportions change?

With this alternative repair strategy, we can revise the state space. Now it does suffice to use 4 states, letting the state correspond to the set of failed copiers, because now we know what the repairman will do when both copiers are down; he will always work on copier 1. Thus it suffices to use the single state $(1, 2)$ to indicate that both machines have failed. There now

Revised Rate Diagram



$\gamma_j =$ rate copier j fails, $\beta_j =$ rate copier j repaired

Figure 3: A revised rate diagram showing the transition rates among the 4 states in Example 3.3, where the repairman always works on copier 1 first when both have failed.

is only one possible transition from state $(1, 2)$: $Q_{(1,2),2} = \mu_1$. We display the revised rate diagram in Figure 3 below.

The associated rate matrix is now

$$Q = \begin{matrix} 0 \\ 1 \\ 2 \\ (1,2) \end{matrix} \begin{pmatrix} -(\gamma_1 + \gamma_2) & \gamma_1 & \gamma_2 & 0 \\ \beta_1 & -(\gamma_2 + \beta_1) & 0 & \gamma_2 \\ \beta_2 & 0 & -(\gamma_1 + \beta_2) & \gamma_1 \\ 0 & 0 & \beta_1 & -\beta_1 \end{pmatrix}$$

or, with the numbers assigned to the parameters,

$$Q = \begin{matrix} 0 \\ 1 \\ 2 \\ (1,2) \end{matrix} \begin{pmatrix} -4 & 1 & 3 & 0 \\ 2 & -5 & 0 & 3 \\ 4 & 0 & -5 & 1 \\ 0 & 0 & 2 & -2 \end{pmatrix}.$$

Just as before, we obtain the limiting probabilities by solving $\alpha Q = 0$ with $\alpha \mathbf{e} = 1$. Now we obtain

$$\alpha \equiv (\alpha_0, \alpha_1, \alpha_2, \alpha_{(1,2)}) = \left(\frac{20}{57}, \frac{4}{57}, \frac{18}{57}, \frac{15}{57} \right).$$

Thus, the long-run proportion of time that copier 1 is working is $\alpha_0 + \alpha_2 = 38/57 = 2/3 \approx 0.67$, while the long-run proportion of time that copier 2 is working is $\alpha_0 + \alpha_1 = 24/57 \approx 0.42$. The new strategy has increased the long-run proportion of time copier 1 is working from 0.62 to 0.67, at the expense of decreasing the long-run proportion of time copier 2 is working from

0.47 to 0.42. The long-run proportion of time the repairman is busy is $1 - \alpha_0 = 37/57 \approx 0.649$, which is very slightly less than before.

We conclude by making some **further commentary**. We might think that the revised strategy is wasteful, because the repairman quits working on copier 2 when copier 1 fails after copier 2 previously failed. By shifting to work on copier 1, we might think that the repairman is being inefficient, “wasting” his expended effort working on copier 2, making it more likely that both copiers will remain failed. In practice, under other assumptions, that might indeed be true, but here because of the lack-of-memory property of the exponential distribution, the expended work on copier 2 has no influence on the remaining required repair times. From a pure efficiency perspective, it might be advantageous to give one of the two copiers priority at this point, but not because of the expended work on copier 2. On the other hand, we might prefer the original strategy from a “fairness” perspective. In any case, the CTMC model lets us analyze the consequences of alternative strategies. As always, the relevance of the conclusions depends on the validity of the model assumptions. But even when the model assumptions are not completely realistic or not strongly verified, the analysis can provide insight. ■

6. Reverse-Time CTMC’s and Reversibility

Just as for DTMC’s, an important concept for CTMC’s is reversibility. A stochastic process $\{X(t) : -\infty < t < \infty\}$ is said to be **reversible** if it has the same probability law as $\{X(-t) : -\infty < t < \infty\}$. Thus a CTMC is reversible if we get the same CTMC if we run time backwards.

Just as for DTMC’s, we can start by constructing the **reverse-time CTMC** associated with a given CTMC with transition-rate matrix Q and transition function $P(t)$. We obtain the reverse-time Markov chain from the original (forward) CTMC by reversing time. The reverse-time transition probabilities describe where the process came from instead of where it is going. That is, we let

$$\overleftarrow{P}_{i,j}(t) \equiv P(X(s) = j | X(s+t) = i) . \quad (6.1)$$

We can then then apply basic properties of conditional probabilities to express these reverse-time transition probabilities in terms of given forward-time transition probabilities; i.e.,

$$\overleftarrow{P}_{i,j}(t) \equiv P(X(s) = j | X(s+t) = i) = \frac{P(X(s) = j, X(s+t) = i)}{P(X(s+t) = i)} = \frac{P(X(s) = j)P_{j,i}(t)}{P(X(s+t) = i)} . \quad (6.2)$$

Unfortunately, however, when we do this, we see that in general we do not obtain bonafide Markov transition probabilities; if we sum over j in (6.2), the transition probabilities do not necessarily sum to 1, as required. We need to assume more. What we assume in addition is that **the Markov chain is in equilibrium**. That is, we assume that the given (forward) CTMC with transition function $P(t)$ is irreducible with initial distribution equal to its stationary distribution α .

Since many reversible processes have infinite-state space, **we henceforth allow an infinite state space**. We assume that the CTMC is irreducible, regular, recurrent and positive recurrent; see §10. Thus the equation $\alpha Q = 0$ has a unique solution such that α is a probability distribution.

In this setting we consider the system in equilibrium by letting $P(X(t) = j) = \alpha_j$ for all t . Under this extra equilibrium condition, equation (6.2) can be expressed as

$$\overleftarrow{P}_{i,j}(t) \equiv P(X(s) = j | X(s+t) = i) = \frac{P(X(s) = j, X(s+t) = i)}{P(X(s+t) = i)} = \frac{\alpha_j P_{j,i}(t)}{\alpha_i} . \quad (6.3)$$

Theorem 6.1. (reverse-time CTMC) *If an positive recurrent irreducible CTMC is put in equilibrium by letting its initial distribution be its stationary distribution α , then the transition probabilities in (6.3) are bonafide transition probabilities, with the same stationary probability distribution α , and the stochastic process satisfies the reverse-time Markov property*

$$\begin{aligned} P(X(s) = j | X(s+t) = i, X(u) = i_u, u \in A_u \subseteq (s+t, \infty)) &= P(X(s) = j | X(s+t) = i) \\ &\equiv \overleftarrow{P}_{i,j}(t). \end{aligned} \quad (6.4)$$

The reverse-time Markov property in (6.4) says that the conditional probability of a “past” state j at time s , given the “present” state i at time $s+t$ plus the states at “future” times u in the set A_u depends only on the present state i at time $s+t$.

Proof. First it is clear that the alleged transition probabilities in (6.3) are nonnegative. Since $\alpha P(t) = \alpha$ for each t , by Theorem 5.2 (d), the row sums of \overleftarrow{P} now do indeed sum to 1, as required. To see that the Markov property in (6.4) does indeed hold, apply properties of conditional probabilities to rewrite (6.4) as

$$\begin{aligned} &P(X(s) = j | X(s+t) = i, X(u) = i_u, u \in A_u \subseteq (s+t, \infty)) \\ &= \frac{P(X(s) = j, X(s+t) = i, X(u) = i_u, u \in A_u \subseteq (s+t, \infty))}{P(X(s+t) = i, X(u) = i_u, u \in A_u \subseteq (s+t, \infty))} \\ &= \frac{P(X(s) = j)P_{j,i}(t)P(X(u) = i_u, u \in A_u \subseteq (s+t, \infty) | X(s+t) = i)}{P(X(s+t) = i)P(X(u) = i_u, u \in A_u \subseteq (s+t, \infty) | X(s+t) = i)} \\ &= \frac{P(X(s) = j)P_{j,i}(t)}{P(X(s+t) = i)} = \frac{\alpha_j P_{j,i}(t)}{\alpha_i}, \end{aligned}$$

where in the last line we have first canceled the common term $P(X(u) = i_u, u \in A_u \subseteq (s+t, \infty) | X(s+t) = i)$ from the numerator and the denominator and then exploited the equilibrium property. ■

From (6.3), we see right away (by looking at the derivative at 0) that the reverse-time CTMC associated with a CTMC having transition-rate matrix Q itself has transition-rate matrix \overleftarrow{Q} , where

$$\overleftarrow{Q}_{i,j} = \frac{\alpha_j Q_{j,i}}{\alpha_i}. \quad (6.5)$$

Note that any irreducible CTMC in equilibrium (initialized with its stationary distribution α) has an associated reverse-time CTMC with transition-rate matrix \overleftarrow{Q} . But that does not make the CTMC reversible. That is a stronger property: A CTMC is said to be time-reversible or just **reversible** if the reverse-time CTMC coincides with the original CTMC in equilibrium, i.e., if $\overleftarrow{Q} = Q$ or, equivalently, if

$$\alpha_i Q_{i,j} = \alpha_j Q_{j,i} \quad \text{for all } i \text{ and } j. \quad (6.6)$$

For CTMC’s, reversibility as defined at the beginning of this section is equivalent to equations (6.6). From a rate-conservation perspective, reversibility holds for a CTMC if and only if the steady-state rate of transitions from state i to state j equals the steady-state rate of transitions from state j to state i for all states i and j . Thus reversibility is characterized by the **detailed-balance equations** in (6.6), which are generalizations of the detailed-balance equations for birth-and-death processes in (4.3). We summarize these properties in the following theorem.

Theorem 6.2. (reversibility of CTMC's) *A CTMC with transition rate matrix Q is reversible with stationary probability vector α if and only if the detailed balance equations in (6.6) hold.*

Proof. We have seen that the detailed balance equations in (6.6) imply reversibility, given that α is the stationary vector. Summing those equations on either i or j gives $\alpha Q = 0$, implying that α must in fact be the stationary probability vector. ■

As a consequence, we immediately have the following result.

Theorem 6.3. (reversibility of birth-and-death processes) *All birth-and-death processes are reversible.*

By the statement in Theorem 6.3, we mean that the birth-and-death process is reversible *provided that it has a stationary probability vector α and is initialized by that stationary probability vector α .*

We now observe that reversibility is inherited by truncation. We say that a CTMC with state space S and rate matrix Q is *truncated* to the subset $A \subset S$ if we disallow all transitions out of the subset A ; i.e., if we set $Q_{i,j} = 0$ for $i \in A$ and $j \in S - A$. We obtain a new CTMC with state space A by letting $Q_{i,j}^{(A)} = Q_{i,j}$ for all i and j in A with $i \neq j$ and $Q_{i,i}^{(A)} = -\sum_{j \in A} Q_{i,j}$ for all $i \in A$.

Theorem 6.4. (truncation) *If a reversible CTMC with rate matrix Q and stationary probability vector α is truncated to a subset A , yielding the rate matrix $Q^{(A)}$ defined above, and remains irreducible, then the truncated CTMC with the rate matrix $Q^{(A)}$ is also reversible and has stationary probability vector*

$$\alpha_j^{(A)} = \frac{\alpha_j}{\sum_{k \in A} \alpha_k}, \quad \text{for all } j \in A. \quad (6.7)$$

Proof. It is elementary that the truncated CTMC with the probability vector $\alpha^{(A)}$ in (6.7) satisfies the detailed balance equations in (6.6) if the original CTMC does, which holds by Theorem 6.2. ■

We have seen an instance of Theorem 6.3 when we looked at the $M/M/1/r$ queueing model in (4.10), following Example 4.3.

We now apply reversibility to get something new and interesting. For that, we consider the $M/M/s$ queue with s servers, unlimited waiting room, Poisson arrival process with arrival rate λ and IID exponential service times with mean $1/\mu$. The number in system over time is a birth-and-death process with constant birth rate λ and death rates $\mu_j = \min\{j, s\}\mu$. Since there is an infinite state space, we require that $\rho \equiv \lambda/s\mu < 1$ in order to have a proper limiting distribution.

Theorem 6.5. (departures from an M/M/s queue) *The departure process from an M/M/s queue in equilibrium (with $\rho \equiv \lambda/s\mu < 1$) is a Poisson process with departure rate equal to the arrival rate λ .*

Proof. Having $\rho \equiv \lambda/s\mu < 1$ ensures that a proper limiting probability vector α exists. Put the system in equilibrium by letting the initial distribution be that limiting distribution. By Theorem 6.3, the CTMC is reversible. Thus, in equilibrium, the process counting the number of customers in the system at any time in reverse time has the same probability law as

the original process. However, reversing time changes departures (jumps down) into arrivals (jumps up) and vice versa. So the departures must form a Poisson process with rate λ . ■

Theorem 6.5 is quite remarkable. It takes some effort to even directly show that the time between successive departures in an $M/M/1$ queue in equilibrium has an exponential distribution with mean $1/\lambda$. That is an instructive exercise.

We can do more: We can establish an even more surprising result. Let $Q(t)$ be the number in system at time t , either waiting or being served, and let $D(t)$ count the number of departures in the interval $[0, t]$. We can show that the distribution of the queue length at time t is independent of the departures prior to time t , for the $M/M/s$ queue in equilibrium!

Theorem 6.6. (more about departures from an $M/M/s$ queue) *For an $M/M/s$ queue in equilibrium (with $\rho \equiv \lambda/s\mu < 1$), the number in system at time t , $Q(t)$, is independent of $\{D(s) : 0 \leq s \leq t\}$, the departure process before time t .*

Proof. Just as for Theorem 6.5, having $\rho \equiv \lambda/s\mu < 1$ ensures that a proper limiting probability vector α exists. As before, put the system in equilibrium by letting the initial distribution be the limiting distribution. By Theorem 6.3, the CTMC is reversible. Thus, in equilibrium, the process in reverse time has the same probability law as the original process. With any one ordering of time, the arrival process after time t is independent of the queue length at time t . by the independent-increments property of the Poisson process. Since arrivals and departures switch roles when we reverse time, we deduce the asserted conclusion as well. ■

We now go on to consider networks of queues. We first combine Theorems 6.5 and 6.6 with Example 4.3 to obtain the limiting distribution for the number of customers at each station for two or more single-server queues in series. In particular, consider an $M/M/1$ model with arrival rate λ and service rate μ_1 , where $\rho_1 \equiv \lambda/\mu_1 < 1$. Let all departures from this $M/M/1$ queue proceed next to a second single-server queue with unlimited waiting room and IID exponential service times with individual service rate μ_2 , where also $\rho_2 \equiv \lambda/\mu_2 < 1$. This model is often referred to as the $M/M/1 \rightarrow /M/1$ tandem queue. Let $Q_i(t)$ be the number of customers at station i at time t , either waiting or being served.

Theorem 6.7. (the limiting probabilities for then $M/M/1 \rightarrow /M/1$ tandem queue) *For the $M/M/1 \rightarrow /M/1$ tandem queue in equilibrium (with $\rho_i \equiv \lambda/s\mu_i < 1$ for each i), the departure processes from the two queues are Poisson processes with rate λ and*

$$\lim_{t \rightarrow \infty} P(Q_1(t) = j_1, Q_2(t) = j_2) = (1 - \rho_1)\rho_1^{j_1}(1 - \rho_2)\rho_2^{j_2} \quad \text{for all } j_1 \text{ and } j_2. \quad (6.8)$$

Theorem 6.7 concludes that the limiting probabilities for the two queues in series are the same as for two independent $M/M/1$ queue, each with the given arrival rate λ . (However, the two stochastic processes $\{Q_1(t) : t \geq 0\}$ and $\{Q_2(t) : t \geq 0\}$. starting in equilibrium, are not independent. The result is for a single time t .) We say that the limiting distribution has **product form**. That product form means that the two marginal distributions are independent.

Proof. Suppose that we initialize the system with the alleged limiting probability distribution. Since it is the product of two geometric distributions, the two marginal distributions are independent. Hence we can first focus on the first queue. By Theorem 6.5, the departure process from this first queue is Poisson with rate λ . Hence, the second queue by itself is an $M/M/1$ queue in equilibrium. Hence each queue separately has the displayed geometric stationary distribution, which coincides with its limiting distribution. Now, considering the system in equilibrium, by Theorem 6.6, at any time t , the random number $Q_1(t)$ is independent of the departure process from the first queue up to time t . That implies that $Q_2(t)$ and

$Q_1(t)$ must be independent for each t , which implies the product-form limiting distribution in (6.8). ■

Note that the entire $M/M/1 \rightarrow /M/1$ tandem queue is not itself reversible; it is possible to go from state (i, j) to state $(i - 1, j + 1)$ (with a departure from queue 1), but it is not possible to go back. So the detailed-balance conditions in (6.6) cannot hold. We established Theorem 6.7 by exploiting the reversibility of only the first station by itself. However, there is an alternative way to prove Theorem 6.7 exploiting only reverse-time CTMC's, which has other applications, e.g., to treat networks of Markovian queues that are not acyclic (do not have flow in one direction only).

Theorem 6.8. (Kelly's Lemma: exploiting reverse-time chains without reversibility) *Let Q be the transition-rate matrix of an irreducible CTMC. If we can find numbers α_j and $\overleftarrow{Q}_{i,j}$ such that*

$$\alpha_i Q_{i,j} = \alpha_j \overleftarrow{Q}_{j,i} \quad \text{for all } i \neq j \quad (6.9)$$

and

$$-Q_{i,i} \equiv \sum_{j,j \neq i} Q_{i,j} = \sum_{j,j \neq i} \overleftarrow{Q}_{i,j} \equiv -\overleftarrow{Q}_{i,i} \quad \text{for all } i, \quad (6.10)$$

then $\overleftarrow{Q}_{i,j}$ are the transition rates for the reverse-time CTMC associated with Q and α is the limiting probability vector for both CTMC's.

Proof. Add over i with $i \neq j$ in (6.9) and apply (6.10) to get

$$\sum_{i,i \neq j} \alpha_i Q_{i,j} = \sum_{i,i \neq j} \alpha_j \overleftarrow{Q}_{j,i} = \alpha_j \sum_{i,i \neq j} Q_{j,i},$$

which implies that $\alpha Q = 0$. Hence α is the limiting distribution for the CTMC with transition-rate matrix Q . Consequently, $\overleftarrow{Q}_{i,j}$ are the transition rates associated with the reverse-time CTMC based on Q in equilibrium. That implies (by summing on j in (6.9)) that α is the limiting distribution for the reverse-time CTMC as well. ■

We now apply Theorem 6.8 to give an alternative proof of Theorem 6.7, which again has the advantage that it does not require directly solving the equation $\alpha Q = 0$.

Alternative proof of Theorem 6.7. As for many harder problems, the first step is to guess the form of the limiting distribution; i.e., we guess that α has the product form in (6.8). We then guess that the reverse-time CTMC should be itself a $M/M/1 \rightarrow /M/1$ tandem queue with arrival rate λ and the given service rates. Going forward from state (i, j) , we have three possible transitions: (1) to state $(i + 1, j)$ due to an arrival, (2) to state $(i - 1, j + 1)$ due to departure from queue 1 and (3) to $(i, j - 1)$ due to departure from queue 2. We have possible flows in the other direction for the reverse-time CTMC to state (i, j) in three possible ways: (1) from state $(i + 1, j)$ due to a departure from queue 1 (in original order), (2) from state $(i - 1, j + 1)$ due to departure from queue 2 and (3) to $(i, j - 1)$ due to an arrival from outside at queue 2. We thus have three equations to check in order to verify (6.9):

$$\begin{aligned} (1 - \rho_1)\rho_1^i(1 - \rho_2)\rho_2^j\lambda &= (1 - \rho_1)\rho_1^{i+1}(1 - \rho_2)\rho_2^j\mu_1 \\ (1 - \rho_1)\rho_1^i(1 - \rho_2)\rho_2^j\mu_1 &= (1 - \rho_1)\rho_1^{i-1}(1 - \rho_2)\rho_2^{j+1}\mu_2 \\ (1 - \rho_1)\rho_1^i(1 - \rho_2)\rho_2^j\mu_2 &= (1 - \rho_1)\rho_1^i(1 - \rho_2)\rho_2^{j-1}\lambda, \end{aligned} \quad (6.11)$$

which are easily seen to be satisfied. It is also easy to see that (6.10) holds. For a state (i, j) with $i > 0$ and $j > 0$, the total rate out to a new state is $\lambda + \mu_1 + \mu_2$, corresponding to the possibilities of an arrival or a service completion at one of the two queues. For the state $(0, j)$ with $j > 0$, the total rate out to a new state is $\lambda + \mu_2$ in both cases, excluding the possibility of a service completion from queue 1 in both cases. The case $(i, 0)$ is similar. For the state $(0, 0)$, the total rate out to new states is λ in both cases, corresponding to an arrival. ■

7. Open Queueing Networks (OQN's)

A minor variation of the argument we have just used to prove Theorem 6.7 for tandem queues by applying Theorem 6.8 (Kelly's lemma) applies to treat a general **Markovian open network of single-server queues (OQN)**, a special case of a Jackson QN. We give the highlights, but we only begin to introduce the rich theory of queueing networks; e.g., see Chen and Yao (2001), Kelly (1979), Sauer and Chandy (1981), Serfozo (1999), van Dijk (1993), Walrand (1988) and Whittle (1986).

7.1. The Model

Let there be m single-server queues, each with unlimited waiting room. Let there be an external Poisson arrival process at queue i with rate $\lambda_{e,i}$, where $\lambda_{e,i} \geq 0$ and $\lambda_{e,i} > 0$ for at least one i . Let the service times at queue i be exponential with mean $1/\mu_i$. Let the arrival processes be mutually independent. Let all the service times be mutually independent and independent of the arrival processes. Let there be **Markovian routing**, independent of the arrival and service processes; i.e., let each customer, immediately after completing service at queue i , go next to queue j with probability $P_{i,j}$, independently of all previous events. Let each customer depart from the network from queue i with probability $1 - \sum_{j=1}^m P_{i,j}$. If we include outside of the network as a single state $m + 1$, then the routing is characterized by a DTMC. In this routing DTMC we assume that $P_{m+1,m+1} = 1$, making the outside state absorbing. Moreover, we assume that all other states are transient states. That is, we assume that each arriving customer will eventually depart from the system.

Consider the vector valued process $Q(t) \equiv (Q_1(t), Q_2(t), \dots, Q_m(t))$, where $Q_i(t)$ is the number of customers at queue i , either waiting or being served, at time t . It is easy to see that the stochastic process $\{Q(t) : t \geq 0\}$ is a CTMC. The possible events are an arrival from outside or a service completion at one of the m queues. Those all governed by the specified rates. We will show that this CTMC also has a product-form limiting distribution, under regularity conditions. Given the possibility of feedback now, it is even more remarkable that the marginal distributions of the limiting probability distribution should be independent.

7.2. The Traffic Rate Equations and the Stability Condition

To characterize the limiting behavior, we first need to find the **total arrival rate** at each queue, i.e., the sum of the external and internal arrival rates. In order to find the total arrival rate to each queue, we need to solve a system of linear equations, **the traffic-rate equations**:

$$\lambda_j = \lambda_{e,j} + \sum_{i=1}^m \lambda_i P_{i,j} \quad \text{for } 1 \leq j \leq m, \quad (7.1)$$

or, equivalently, in matrix notation,

$$\Lambda = \Lambda_e + \Lambda P, \quad (7.2)$$

which implies that

$$\Lambda = \Lambda_e(I - P)^{-1} . \quad (7.3)$$

The inverse $(I - P)^{-1}$ is the fundamental matrix of the absorbing routing DTMC.

In order for the solution of the traffic-rate equations to be valid total arrival rates, we have to be sure that the net arrival rate is less than the maximum possible service rate at each queue. As before, let the **traffic intensity** at queue j be

$$\rho_j = \frac{\lambda_j}{\mu_j}, \quad 1 \leq j \leq m . \quad (7.4)$$

We assume that $\rho_j < 1$ for all j , $1 \leq j \leq m$.

7.3. The Limiting Product-Form Distribution

With those assumptions, the CTMC $\{Q(t) : t \geq 0\}$ has a product-form limiting distribution.

Theorem 7.1. (Markovian open network of single-server queues) *Consider the Markovian open network of single server queues defined above, where the inverse $(I - P)^{-1}$ is well defined for the routing matrix P and $\rho_i < 1$ for each i . Then the limiting distribution has product form, with geometric marginals of the $M/M/1$ queue; i.e.,*

$$\lim_{t \rightarrow \infty} P(Q(t) = (j_1, \dots, j_m)) = \alpha(j_1, \dots, j_m) = \prod_{k=1}^m (1 - \rho_k) \rho_k^{j_k} . \quad (7.5)$$

Proof. A direct proof is to guess that the solution is of product form, as in (7.5), and then simply verify that $\alpha Q = 0$. That verification step is simplified by applying Theorem 6.8. To do so, we need to define the candidate reverse-time CTMC with transition rates \overleftarrow{Q} . Just as with Theorem 6.7, we guess that the reverse-time CTMC itself corresponds to an open network of single-server queues, with the same service-time distributions at the queues, but we need to guess the appropriate external arrival rates $\overleftarrow{\lambda}_{e,i}$ and routing probabilities $\overleftarrow{P}_{i,j}$. The idea is to guess those quantities by seeing what is required to have the flow rates balance in equilibrium.

First, the flow rate through each queue should be the same in either direction, so we should have

$$\overleftarrow{\lambda}_i = \lambda_i \quad \text{for all } i . \quad (7.6)$$

Next, the reverse-time external arrival rate at queue i should be equal to the forward-time departure rate from queue i , i.e.,

$$\overleftarrow{\lambda}_{e,i} = \lambda_i \left(1 - \sum_{j=1}^m P_{i,j}\right) \quad \text{for all } i . \quad (7.7)$$

To complete the construction, note that the stationary flow from queue i to queue j in forward time should equal the stationary reverse flow in reverse time, i.e.,

$$\lambda_i P_{i,j} = \overleftarrow{\lambda}_j \overleftarrow{P}_{j,i} . \quad (7.8)$$

As a consequence, we have

$$\overleftarrow{P}_{j,i} = \frac{\lambda_i P_{i,j}}{\overleftarrow{\lambda}_j} . \quad (7.9)$$

Combining equations (7.6), (7.7) and (7.9), we have defined the reverse-time model elements $\overleftarrow{\lambda}_{e,i}$, $\overleftarrow{\lambda}_i$ and $\overleftarrow{P}_{j,i}$ for all i and j in terms of the corresponding forward-time modelling elements.

For the reverse-time queueing network, we should have an analog of the traffic-rate equations in (7.1) in reverse time, namely,

$$\overleftarrow{\lambda}_j = \overleftarrow{\lambda}_{e,j} + \sum_{i=1}^m \overleftarrow{\lambda}_i \overleftarrow{P}_{i,j} \quad \text{for } 1 \leq j \leq m, \quad (7.10)$$

or, equivalently, in matrix notation,

$$\overleftarrow{\Lambda} = \overleftarrow{\Lambda}_e + \overleftarrow{\Lambda} \overleftarrow{P}. \quad (7.11)$$

And, indeed, it is easy to check that these reverse-time traffic rate equations are valid, by applying the definitions above.

Just as for the tandem OQN in Theorem 6.7, we assume that the reverse-time OQN has independent Poisson arrival processes at each queue and Markovian routing. The external arrival rate at queue j is then $\overleftarrow{\lambda}_{e,j}$ satisfying (7.7). The routing probabilities are $\overleftarrow{P}_{j,i}$ as given in (7.9). It then remains to verify that these guesses yield the right answer; i.e., we need to verify equations (6.9) and (6.10) in this setting, remembering that the states now correspond to m -tuples (i_1, \dots, i_m) . That is straightforward, paralleling the proof of Theorem 6.7. As before, all transitions are triggered by arrivals and service completions (followed by a random routing). ■

From the reverse-time construction just completed, we also can deduce the following corollary.

Corollary 7.1. (Output Theorem: departure processes from an OQN) *Under the assumptions of Theorem 7.1, the m departure processes from the network from the individual queues are m independent Poisson processes, with the departure process from queue i having rate*

$$\delta_i \equiv \lambda_i \left(1 - \sum_{j=1}^m P_{i,j}\right). \quad (7.12)$$

Moreover, the total departure process is a Poisson process with rate

$$\delta \equiv \sum_{j=1}^m \delta_j = \sum_{j=1}^m \lambda_{e,j}. \quad (7.13)$$

Proof. By construction, the external arrival processes to the queues of the reverse-time OQN are independent Poisson processes, just as assumed for the original OQN. However, the arrival processes in the reverse-time OQN coincide with the forward-time departure processes from the network, with the specified rates. We only look at those departures that actually leave the entire network. We do not look at arrivals routed to some other queue.) To complete the proof, it is good to have a characterization of a Poisson process that does not depend on the ordering of time. With our usual “forward thinking,” we think of the arrivals in an interval $[t, t + u]$, being independent of the process before time t . However, instead it is useful to characterize the Poisson process as a process with single jumps that has stationary and independent increments. That characterization applies both forwards and backwards. Thus a reverse-time Poisson process necessarily also is a Poisson process. For (7.13), recall that the superposition of independent Poisson processes is Poisson with a rate equal to the sum of the rates. To see that the total rate in equals the total rate out, as we would expect, compare the sum over i of (7.12) to the sum over j of (7.1). ■

Remark 7.1. (*non-Poisson internal flows*) While the flows into and out of the network at the queues are independent Poisson processes, in general the internal flows are *not* Poisson processes, because of the feedback. To illustrate, consider a simple OQN consisting of one single-server queue, where each customer after completing service returns to the back of the queue for an additional service with probability p , independent of the history prior to that service completion. If the external arrival rate at that queue is λ_e , then the net arrival rate is $\lambda = \lambda_e/(1 - p)$. Assume that $\lambda < \mu$, where μ is the service rate at the queue, so that the number of customers in the system is a positive recurrent CTMC. We now want to see that the internal departure process cannot be Poisson.

To see that the departure (service completion) process cannot be Poisson, we observe that it cannot have independent increments. Consider two successive increments $D(t) - D(t - u)$ and $D(t + u) - D(t)$. As $D(t) - D(t - u)$ increases, the next increment $D(t + u) - D(t)$ is more likely to be big, because many of the departures in the interval $[t - u, t]$ we will return for an additional service. To make this effect evident, Suppose that $p = 1/2$, $\lambda_e = \epsilon$ and $\mu = 2$, so that $\lambda = 2\epsilon$ and $\rho = \epsilon$, so that the mean number in the system at any time is $\epsilon/(1 - \epsilon)$. Then let ϵ be small. Since the flow rates are very small, $P(D(t + u) - D(t) > 1)$ will be very small, but the conditional probability $P(D(t + u) - D(t) > 1 | D(t) - D(t - u) = M)$ can be much larger if M is large, because many of the departures in the interval $[t - u, t]$ we will return for an additional service.

7.4. Extensions: More Servers or Different Service Scheduling Rules

We have assumed that each queue has only a single server with the **first-come first-served (FCFS)** service discipline, but the theory extends to other models for the individual queues.

1. More Servers. First, Theorem 7.1 and Corollary 7.1 extend to Markovian OQN's with **multi-server queues**, where there may be different numbers of servers at each queue. Again we assume that the service-time distribution is exponential and the FCFS service-scheduling rule is used. Again there is a product-form limiting distribution, but now the marginal steady-state distributions at each queue have the limiting distribution of the $M/M/s$ queue, where s is the number of servers at that queue. Even more generally, the service-rate at each queue may be state dependent with a rate depending on the number of customers there. The marginal distribution is again given by the birth-and-death process steady-state. In all these extensions we need to assume that the net rate into each queue is strictly less than the maximum rate out, where the rate in is determined by the traffic rate equations. For the $M/M/s_i$ queue at node i , the revised traffic intensity is

$$\rho_i \equiv \frac{\lambda_i}{s_i \mu_i}.$$

We need $\rho_i < 1$ for all i .

2. Non-FCFS Service-Scheduling Rules. Second, Theorem 7.1 and Corollary 7.1 extend to Markovian OQN's with **single-server queues** but certain **special non-FCFS scheduling rules**. There are two alternative scheduling rules for single-server queues: **processor-sharing (PS)** and **last-come first-served with preemptive-resume (LCFS-PR)**. With these alternative non-FCFS service-scheduling rules, the steady-state has the **insensitivity** property: The service-time distribution need not be exponential. The successive service times can be i.i.d. with a general distribution (having finite mean). The product-form steady-state distribution depends on that general service-time distribution at that queue only through its mean; the steady state distribution is the same as if it were exponential.

In managing systems with queues, service scheduling can be very important. Thus, the extension to non-FCFS service-scheduling rules has proven to be very useful. The PS service

discipline has been extensively used. It serves as a good approximation of the **round robin (RR) scheduling rule**. With *RR* and *PS*, all customers requiring service can be said to be in service. Customers have **service requirements** determined upon arrival to the queue. The customer remains at the queue until all its required service has been provided. With *PS*, at each instant, the available service rate of the server is divided equally among all customers requiring service at that instant. Hence, in general, unlike FCFS, customers do not depart in the order they arrive.

The PS scheduling rule is a convenient mathematical approximation for the RR scheduling rule. With the RR rule, only one customer receives service at any time, but the customers take turns, in a round robin order, each receiving a small fixed quantum of service, until their requirement is met. The customer departs when the initial service requirement is met; otherwise the customer returns to the end of the line, to wait until it can receive the next quantum of service. Unfortunately, the RR rule, is not so easy to analyze, and does not have the insensitivity property, so that the PS model is used as an approximation for the RR rule.

The LCFS-PR scheduling rule is less frequently used, but it also leads to the product-form steady-state distribution with the insensitivity property. With LCFS-PR, all service rate is allocated to the last customer to arrive. Whenever a new arrival comes, that new arrival replaces the customer in service, while the customer that was in service goes to the head of the queue of waiting customers. However, the partial service provided to each customer is not lost. When a customer who had his service interrupted returns to the server, his service time resumes where it left off. The remaining requirement is the initial requirement reduced by all service provided so far. With LCFS-PR, immediately enters service upon arrival and might complete service before another customer arrives. However, customers might have their service interrupted multiple times. Given that the model is stable (the traffic intensity is less than 1 at each queue), all customers will eventually be served.

The product-form with the insensitivity property also holds for queues that are **infinite-server (IS) queues**; then each customer enters service immediately upon arrival. The product form is then the same as if the queue were an $M/M/\infty$ queue, which is covered by the first extension above. For these three alternative individual queue models – PS, LCFS-PR and IS, the service-time distribution need not be exponential. Nevertheless, the product-form steady-state distribution of the OQN is the same as if the queues were FCFS queues with exponential service distributions.

7.5. Steady State in Discrete and Continuous Time

Throughout our study of CTMC's, we focus on the limiting steady-state probability vector α , which describes the probability distribution of the model at an arbitrary time in steady state. It is important to be aware that this same distribution might not be apply at other special times. Indeed, we have already seen in §3.1 and §5 that the steady-state distribution in continuous time, α , is in general different from the steady-state at transition epochs, which we have denoted by π ; see (3.2) and (5.7).

Similarly, the steady-state distribution at arrival epochs and at departure epochs need not agree with the continuous-time steady-state distribution. However, there is a principle called **Poisson Arrivals See Time Averages (PASTA)** that implies that the proportion of customers that see some state upon arrival coincides with the proportion of time the process spends in that state, provided that the arrival process is Poisson (and other regularity conditions hold, which they do here; e.g., see §5.16 of Wolff (1989), Melamed and Whitt (1990) or Stidham and El Taha (1999). In particular, PASTA always holds for a Poisson process independent of the system, as would be seen by an outside observer, but it also holds for a

Poisson process that is part of the system under an additional lack of anticipation assumption.

Closely related to PASTA is the **arrival theorem**, which describes the steady-state seen be an arrival at a queue in a queueing network. The arrival theorem is valid for internal flows in the network even though they are not Poisson. In an OQN, the arrivals at a queue see the steady-state even though the arrival may come on an internal flow that is not Poisson; e.g., see Melamed and Whitt (1990). For associated closed queueing networks, having a fixed population of customers, which we consider next, arrivals see the network with that customer removed. That is the steady-state distribution for the network with the population reduced by 1.

The connections between stationary distributions for the same system viewed in discrete and continuous time has been studied for more general stochastic processes than Markov processes. This issue is at the heart of the theory of stationary point processes and stationary marked point processes.

8. Closed Queueing Networks (CQN's)

We obtain a closed queueing network (CQN) from an OQN if we eliminate all external arrival processes from the model and instead assume that there is a fixed number (population) of customers in the system, say N . The resulting model is a Gordon-Newell QN or a closed Jackson QN.

8.1. Why Are Closed Models Interesting?

From a modeling perspective, it is natural to wonder why there would be much interest in a CQN, because in most systems customers or jobs arrive from outside, move around and then eventually depart when all required service is complete. At first glance, an OQN seems much more natural than a CQN. However, CQN's have been used extensively, perhaps even more than OQN's, especially in models of computer systems.

When closed models are used, people think of a new job replacing an old one whenever the old job has completed its required service. Which model is appropriate depends on what you think can be directly specified and controlled and what needs to be calculated by analyzing a model. With some systems, it is natural to think of the **population** of jobs in the system as the **independent variable**, which can be directly specified and controlled, while the **throughput** (the departure rate of completed service from the network) is the **dependent variable**, which needs to be determined by analyzing a model. That perspective leads to a CQN.

In contrast, with an OQN, the throughput is the independent variable, which is directly specified as the total arrival rate (which equals the total departure rate), while the population is a dependent *random* variable, which is to be determined, e.g., via its steady-state distribution or the expected value of that steady-state distribution.

However, OQN's tend to be easier to analyze than CQN's, especially when more general non-Markovian models are contemplated. It thus can be useful to observe that the CQN perspective can be achieved with OQN's by using the **fixed population mean (FPM) method**, discussed in Whitt (1984) and §5.5 of Walrand (1988). The idea is to regard the expected steady-state system population as given, and then do a search over all external arrival rates to find that external arrival rate that produces the desired mean steady-state total number in the OQN. For small systems, the performance of the OQN with the FPM method tends to be very different from the performance of the CQN, because the steady-state number in the system tends to be quite variable in the OQN. Surprisingly, however, the performance tends to be much less in large systems, either with many queues or with a large population. In such

larger OQN's, the variance of the total system population in steady state tends to be relatively small compared to the mean.

8.2. A Normalized Product-Form Distribution

Markovian CQN's can be analyzed much like OQN's. They are still CTMC's and Kelly's lemma can still be applied, but there are some significant differences.

As for the OQN, let the CQN have m single-server queues with service rate μ_i at queue i . As before, let the routing matrix be denoted by $P \equiv (P_{i,j})$, but now we assume that it is the transition matrix of an irreducible DTMC instead of a transient DTMC.

Let the vector of queue lengths (number at each queue, including in service, if any) be $Q(t) \equiv (Q_1(t), \dots, Q_m(t))$ and let possible states be vectors $n \equiv (n_1, \dots, n_m)$, where $n_i \geq 0$ for all i and we impose the **CQN population constraint**

$$\sum_{i=1}^m n_i = N. \quad (8.1)$$

Each transition is a service completion combined with a routing of that customer to the end of another queue. In particular, the rate of a service completion at queue i that goes to queue j is $\mu_i P_{i,j}$ for all states $n \equiv (n_1, \dots, n_m)$ such that $n_i > 0$. To concisely describe transitions among states, let e_i be a vector with a 1 in the i^{th} place and 0's elsewhere. The CTMC $\{(Q_1(t), \dots, Q_m(t)) : t \geq 0\}$ is specified by its rate "matrix." For any state vector n with $n_i > 0$, transitions to state vector $n - e_i + e_j$ occur at rate $\mu_i P_{i,j}$, i.e., the transition rates are specified by

$$Q_{n, n - e_i + e_j} \equiv \mu_i P_{i,j} \quad \text{for all } n \text{ with } n_i > 0.$$

The rate matrix is specified by having

$$\sum_{n'} Q_{n, n'} = 0.$$

(This is just the vector analog of having all row sums in the rate matrix be 0.) In this setting, we see that the limiting steady-state probability vector is the familiar $\alpha \equiv \alpha_n \equiv \alpha_{n_1, \dots, n_m}$ such that

$$\alpha Q = 0.$$

What remains is to show that this steady-state probability vector has a simple modification of the product form we saw for OQN's. In particular, it is a normalization of the steady-state distribution of the OQN to enforce the population constraint in (8.1) above. To characterize this steady-state distribution, we need the CQN traffic rate equations. Since there are no external arrival processes, instead of (7.1) and (7.2), we have the **CQN traffic rate equations**

$$\lambda_j = \sum_{i=1}^m \lambda_i P_{i,j} \quad \text{for } 1 \leq j \leq m, \quad (8.2)$$

or equivalently, in matrix notation,

$$\Lambda = \Lambda P, \quad (8.3)$$

where $\Lambda \equiv (\lambda_1, \dots, \lambda_m)$ is a $1 \times m$ row vector and P is the usual $m \times m$ Markov transition matrix.

We immediately observe that the CQN traffic equations in (8.2) and (8.3) agrees with the fundamental DTMC equation $\pi = \pi P$, which has a unique solution as a probability

vector, because P is the transition matrix for an irreducible m -state DTMC. Hence, the unique stationary probability vector π is one such solution. However, here there is no reason why the flow rates should sum to 1. Any positive multiple $c\pi$ for $c > 0$ is another solution. From the DTMC theorem, those are the only solutions. This one-parameter family of solutions to the traffic rate equations does specify the **relative arrival rates**. The ratio of the flow rate through queues i and j necessarily is $\lambda_i/\lambda_j = c\pi_i/c\pi_j$, which is the same as π_i/π_j . So we have determined the relative traffic rates. The remaining degree of freedom is removed by the population constraint (8.1).

Theorem 8.1. (Markovian closed network of single-server queues) *Consider the CQN with single-server queues defined above, where $\Lambda \equiv (\lambda_1, \dots, \lambda_m)$ is a solution to the traffic rate equations in (8.2) or (8.3). Then the limiting distribution has the normalized product form,*

$$\alpha_n \equiv \alpha_{(n_1, n_2, \dots, n_m)} \equiv \frac{f(n)}{G(N)}, \quad (8.4)$$

where

$$f(n) \equiv f((n_1, n_2, \dots, n_m)) = \prod_{i=1}^m \rho_i^{n_i}, \quad (8.5)$$

with $\rho_i \equiv \lambda_i/\mu_i$, $1 \leq i \leq m$,

$$G(N) \equiv \sum_{n: n \in S(N)} f(n) = \sum_{n_1=0}^N \sum_{n_2=0}^{N-n_1} \dots \sum_{n_{m-1}=0}^{N-\sum_{i=1}^{m-1} n_i} f\left(n_1, \dots, n_{m-1}, N - \sum_{i=1}^{m-1} n_i\right) \quad (8.6)$$

and

$$S(N) \equiv \{n \equiv (n_1, \dots, n_m) : n_i \geq 0 \text{ for all } i \text{ and } \sum_{j=1}^m n_j = N\}. \quad (8.7)$$

First Proof. We can apply Theorem 6.8 (Kelly's lemma). We guess that the reverse-time CTMC is just another CQN of the same form with the same service rates and the same relative throughput rates $\overleftarrow{\lambda}_i = \lambda_i$ for all i . Then just as in (7.8) and (7.9), we should have $\lambda_i P_{i,j} = \overleftarrow{\lambda}_j \overleftarrow{P}_{j,i}$, so that $\overleftarrow{P}_{j,i} = \lambda_i P_{i,j} / \overleftarrow{\lambda}_j$. Then, paralleling the proofs of Theorems 6.7 and 7.1, we can check that the detailed balance equation holds, namely,

$$\alpha_n Q_{n, n-e_i+e_j} = \alpha_{n-e_i+e_j} \overleftarrow{Q}_{n-e_i+e_j, n}, \quad (8.8)$$

so that the proof is completed by invoking Theorem 6.8. ■

Second Proof. The difficulty with the first proof is that it is not so clear how to guess the form of the steady-state probability vector α . The second approach is to apply a related OQN to guess the form of the steady-state distribution for the associated CQN.

We construct the OQN by performing a cut on one arc with positive flow. Choose two states in the CQN where there is positive routing probability in one direction. Without loss of generality, this can be from queue 2 to queue 1, possibly after relabeling the queues. Assuming that $P_{2,1} > 0$, create an associated OQN by making two changes: (i) introduce a Poisson external arrival process at queue 1 with small rate $\epsilon > 0$ and (ii) let all flow that would have gone from queue 2 to queue 1 leave the network.

The resulting OQN has total external arrival rate ϵ and thus total departure rate ϵ . We now solve the OQN. For sufficiently small ϵ , we will have $\rho_i \equiv \lambda_i/\mu_i < 1$ at each queue in

the OQN. Under that condition, the OQN will have a unique stationary probability vector α satisfying $\alpha Q = 0$. In particular, we will have

$$f_{OQN}(n) \equiv f_{OQN}((n_1, \dots, n_m)) = \prod_{i=1}^m (1 - \rho_i) \rho_i^{n_i} \quad (8.9)$$

However this unique invariant measure for the OQN is also an invariant measure for the CQN.

Of course, (8.9) is not the same as $f(n)$ in (8.5) because there is the extra term

$$A \equiv \prod_{i=1}^m (1 - \rho_i), \quad (8.10)$$

in (8.9), but A in (8.10) is a constant, so *both* (8.9) and (8.5) are stationary measures for the CQN CTMC, by the same argument given above. The differing constant will be removed by the normalization. From this point, we can apply Theorem 6.8 (Kelly's lemma) to complete the proof. This reasoning also leads to formulas (8.4)-(8.7) above. ■

8.3. Computing the Normalization Constant: The Convolution Algorithm

One difficulty remains: We still need to compute the normalization constant $G(N)$ in (8.6), which is also called the partition function. This can be done by a recursive algorithm, called the convolution algorithm or the Buzen algorithm. Let the queues be labeled in order from 1 to m . Instead of the notation $G(N)$, now use new notation that counts the number of queues, with the given initial ordering. Let

$$G(N, m) = \sum_{\{\mathbf{n}: n_1 + n_2 + \dots + n_m = N\}} \prod_{i=1}^m \rho_i^{n_i}, \quad (8.11)$$

where, as before, $\rho_i \equiv \lambda_i / \mu_i$, $1 \leq i \leq m$, with $\lambda \equiv (\lambda_1, \dots, \lambda_m)$ a solution to (8.2).

Theorem 8.2. (*Buzen's convolution algorithm*) *For any $N > 1$ and $M > 1$, the normalization constant $G(N, m)$ defined in (8.11) can be solved by the two-dimensional recursion*

$$G(N, m) = g(N, m - 1) + \rho_m G(N - 1, m). \quad (8.12)$$

Proof. The first term on the right in (8.12) covers the case in which $n_m = 0$, while the second case on the right covers the complementary case in which $n_m > 0$. ■

An alternative way to compute the normalization constant is to numerically invert its generating function. This approach becomes more important in multi-chain multi-class CQN's where there are multiple population constraints, producing normalization constants of the form $G(N) = G((N_1, \dots, N_p))$; see Choudhury et al. (1995a). This alternative approach is appealing because the numerical inversion makes it possible to compute any single value $G(N)$ without computing the function for all vectors N' with $N' < N$. The generating function can also be exploited to perform asymptotic approximations.

The generating functions have a remarkably simple form. For the basic model with normalization constant $G(N)$ in (8.6), its generating function is defined as

$$\hat{G}(z) \equiv \sum_{N=0}^{\infty} G(N) z^N. \quad (8.13)$$

Proposition 1. (*generating function of the normalization constant*) For the CQN above with m single-server queues,

$$\hat{G}(\mathbf{z}) \equiv \sum_{N=0}^{\infty} G(N)z^N = \prod_{i=1}^m (1 - \rho_i z^i)^{-1}. \quad (8.14)$$

Proof. The main idea is to change the order of summation, which is carried out in the third line below. Substituting in $G(N)$ from (8.6) into (8.13), we obtain

$$\begin{aligned} \hat{G}(\mathbf{z}) &\equiv \sum_{N=0}^{\infty} G(N)z^N = \sum_{N=0}^{\infty} \left(\sum_{\mathbf{n} \in S(N)} f(\mathbf{n}) \right) z^N \\ &= \sum_{N=0}^{\infty} \left(\sum_{\mathbf{n} \in S(N)} \prod_{i=1}^m \rho_i^{n_i} \right) z^N = \sum_{N=0}^{\infty} \left(\sum_{\mathbf{n} \in S(N)} \prod_{i=1}^m \rho_i^{n_i} z^{n_i} \right) \\ &= \sum_{n_1=0}^{\infty} \cdots \sum_{n_m=0}^{\infty} \left(\prod_{i=1}^m \rho_i^{n_i} z^{n_i} \right) \\ &= \prod_{i=1}^m \left(\sum_{n_i=0}^{\infty} \cdots \sum_{n_q=0}^{\infty} \rho_i^{n_i} z^{n_i} \right) = \prod_{i=1}^m (1 - \rho_i z)^{-1} \end{aligned} \quad (8.15)$$

as claimed. ■

9. Stochastic Loss Models

Another important class of stochastic networks are the loss networks. These can be viewed as generalizations of the classical Erlang loss model. Reversibility also plays a role here. We illustrate how reversibility can be used to deduce the insensitivity property of the Erlang loss model.

9.1. The Erlang Loss Model

This is the Markovian $M/M/N/0$ queueing model with Poisson arrival process, N servers, exponential service times and no extra waiting space. The number of customers in the system (all in service) is a finite-state BD process. The birth rates are $\lambda_k \equiv \lambda$ and the death rates are $\mu_k = k\mu$ for an overall arrival rate of λ and an individual service rate of μ . The detailed balance equations give

$$\alpha_j \lambda = \alpha_{j+1} (j+1) \mu, \quad (9.1)$$

so that, by Theorem 4.2,

$$\alpha_j = \frac{a^j}{G j!}, \quad 0 \leq j \leq N, \quad (9.2)$$

where $a \equiv \lambda/\mu$ in (9.2) is the **offered load** (the expected number of busy servers in the associated infinite-server model) and the normalization constant $G \equiv G(N) \equiv G(N, a)$ is chosen so that the steady-state probabilities α_j sum to 1. Hence

$$G(N) = \sum_{j=0}^N \frac{a^j}{j!} \quad (9.3)$$

This system can be viewed as a **truncation** of the associated $M/M/\infty$ infinite-server system, where

$$G(N) = \sum_{j=0}^{\infty} \frac{a^j}{j!} = e^a; \quad (9.4)$$

then the steady-state distribution is Poisson with mean a . (The truncation can be viewed as a consequence of Theorem 6.4.)

The principle performance measure for the Erlang loss model is the steady state blocking probability experienced by an arrival, which coincides with the steady-state probability α_N by the PASTA property, which can be expressed directly in terms of the normalization constant:

$$B(N, a) = \alpha_N = 1 - \frac{G(N-1, a)}{G(N, a)}. \quad (9.5)$$

For large N and a , it is convenient to calculate the blocking probability $B(N, a)$ by a recursion. The standard recursion is obtained in the following proposition. The final expression involving the traffic intensity $\rho \equiv a/N$ tends to be convenient for numerical calculations because ρ tends to be near 1 in applications. It is instructive to write programs to compute $B(10^k, 10^k)$ for $k = 1, 2, 3$ and 4; compare obvious alternatives to the recursion below. For more on the Erlang blocking formula, see Jagerman (1974) and Whitt (2002).

Proposition 2. (*recursion for the blocking probability*) *The blocking probability in the $M/M/N/0$ model satisfies the recursion*

$$B(N, a) = \frac{1}{1 + [N/aB(N-1, a)]} = \frac{aB(N-1, a)}{N + aB(N-1, a)} = \frac{\rho B(N-1, a)}{1 + \rho B(N-1, a)}, \quad (9.6)$$

where $B(0, a) \equiv 1$, $a \equiv \lambda/\mu$ is the offered load and $\rho \equiv a/N$ is the traffic intensity.

Proof. It is convenient to work with the reciprocal, because it tends to be easier to analyze. Note that

$$R(N, a) \equiv \frac{1}{B(N, a)} = \frac{G(N)}{a^N/N!} = \frac{G(N-1) + (a^N/N!)}{a^N/N!} = \frac{NR(N-1, a)}{a} + 1,$$

from which the recursion (9.6) follows easily. ■

To understand the application of numerical inversion to calculate normalization constants (partition functions) by inverting their generating functions in more general loss models, we show how to calculate the generating function here.

Proposition 3. (*generating function of the normalization constant*) *For the $M/M/N/0$ normalization constant $G(N)$ in (9.3), the generating function is*

$$\hat{G}(z) \equiv \sum_{N=0}^{\infty} G(N)z^N = \frac{e^{az}}{1-z}. \quad (9.7)$$

where $a \equiv \lambda/\mu$ is the offered load.

Proof. We reason as for Proposition 1, changing the order of summation at the end of the first line:

$$\begin{aligned}
\hat{G}(z) &\equiv \sum_{N=0}^{\infty} G(N)z^N = \sum_{N=0}^{\infty} \left(\sum_{j=0}^N \frac{a^j}{j!} \right) z^N = \sum_{j=0}^{\infty} \sum_{N=j}^{\infty} \frac{a^j}{j!} z^N \\
&= \sum_{j=0}^{\infty} \frac{(az)^j}{j!} \sum_{N=j}^{\infty} z^{N-j} = \sum_{j=0}^{\infty} \frac{(az)^j}{j!} \sum_{N=0}^{\infty} z^N \\
&= \left(\frac{1}{1-z} \right) \sum_{j=0}^{\infty} \frac{(az)^j}{j!} = \frac{e^{az}}{1-z}. \quad \blacksquare
\end{aligned} \tag{9.8}$$

9.2. Stochastic Loss Networks

Stochastic loss networks are generalizations of the Erlang loss model model with multiple classes of customers and multiple facilities. Customers require and use subsets of resources from each facility and hold them all for the duration of their time in the system; see Kelly (1991).

As in Kelly (1991), we use telecommunications terminology in our description of the model. Consider a **network** with J **links** (the facilities) indexed by j , $1 \leq j \leq J$. Suppose that link j contains C_j **circuits** (the resources). Let there be R customer **classes** (often called routes, **classes** \equiv **routes**), indexed by r , $1 \leq r \leq R$. Customers of class r **arrive** according to a Poisson process with rate λ_r . If the customer can get all its required resources in the network, the class- r customer uses these resources (so that they are temporarily not available to other customers) and stays in the network for an exponential length of time (**service or holding time** with mean $1/\mu_r$). When a customer leaves, it releases all its resources. A class- r customer requires $A_{j,r}$ circuits (units of resource) on link j , $1 \leq j \leq J$. Thus customers may use more than one circuit on a link and require circuits on more than one link. If a customer can not be given all its required resources immediately upon arrival, then this customer is blocked without affecting the future arrivals.

Let $X(t) = (X_1(t), \dots, X_R(t))$ be the number of customers from each class in the system at time t . The R -dimensional stochastic process is a CTMC. It is easily seen to be a **reversible CTMC**. The local balance equation is

$$\alpha_n \lambda_r = \alpha_{n+e_r} \mu_r (n_r + 1), \tag{9.9}$$

where $n \equiv (n_1, n_2, \dots, n_R)$ and e_i is an R -dimensional vector with a 1 in the i^{th} place and 0 elsewhere. The system has a steady-state distribution that is a natural generalization of the Erlang loss model. Following Kelly (1991), we let $\nu_r \equiv \lambda_r / \mu_r$ be the offered load for class/route r .

Theorem 9.1. (*steady-state of a stochastic loss network*) *The stochastic loss network defined above is a reversible finite-state CTMC with steady-state distribution*

$$\alpha_n = G(C)^{-1} f(n), \quad n \in S(C), \tag{9.10}$$

where

$$f(n) = f(n_1, n_2, \dots, n_R) \equiv \prod_{r=1}^R \frac{\nu_r^{n_r}}{n_r!} \tag{9.11}$$

$\nu \equiv \lambda_r / \mu_r$,

$$S(C) \equiv S(C_1, C_2, \dots, C_J) \equiv \{n \in \mathbb{R}_+^R : An \leq C\} \tag{9.12}$$

and

$$G(C) \equiv G(C_1, C_2, \dots, C_J) \equiv \sum_{n \in S(C)} f(n). \quad (9.13)$$

The set $S(C)$ contains the possible states of the process, while $G(C)$ in (9.13) is the normalization constant or partition function. The blocking probability of an arbitrary class- r request is

$$B_r = 1 - \frac{G(C - Ae_r)}{G(C)}, \quad (9.14)$$

where e_r is a vector that is 1 in the r^{th} place and 0 otherwise.

Proof. The statement may seem complicated, but both it and the proof are elementary when viewed properly. First suppose that there are no capacity constraints. Then the R classes become independent and the number of route/class r evolves as the number of busy servers in an $M/M/\infty$ queueing model with arrival rate λ_r and service rate μ_r . Hence, the overall steady-state distribution is just the product of these R steady-state Poisson distributions. Moreover, each of these is a birth and death process and so reversible, In addition, the entire system is reversible. The final result is then obtained by truncation, i.e., by applying Theorem 6.4, just as for the Erlang loss model. ■

What is challenging here, just as with closed queueing networks, is to compute the normalization constant. One approach is to develop asymptotic approximations for large systems, while another approach is to develop efficient numerical algorithms. Just as for CQN's and the Erlang loss model, the normalization constant can be computed by numerically inverting its generating function. Again, the generating function has a surprisingly simple form. For numerical inversion of generating functions of this normalization constant, see Choudhury et al. (1995b).

Proposition 4. (*generating function*) For the stochastic loss model above, the generating function of the normalization constant in (9.13) is a generalization of (10.5), namely,

$$\begin{aligned} \hat{G}(z) &\equiv \hat{G}(z_1, \dots, z_J) \equiv \sum_C G(C) \prod_{j=1}^J z_j^{C_j} \\ &= \sum_{C_1=0}^{\infty} \cdots \sum_{C_J=0}^{\infty} G(C_1, \dots, C_J) z_1^{C_1} \times \cdots \times z_J^{C_J} \\ &= \frac{e^{\left(\sum_{r=1}^R \nu_r \prod_{j=1}^J z_j^{A_{j,r}}\right)}}{\prod_{j=1}^J (1 - z_j)}. \end{aligned} \quad (9.15)$$

(There is a typo in (2.12) on p. 1115 of Choudhury et al. (1995b); the second sum in the numerator there should be a product.)

The story about loss networks is much longer. See Kelly (1991) and references there for more interesting material including asymptotic approximations such as the Erlang fixed point theorem and phase transitions.

9.3. Insensitivity in the Erlang Loss Model

The $M/M/N/0$ Erlang loss model has the insensitivity property; i.e., the same steady state distribution still holds in the $M/GI/N/0$ model, having i.i.d. service times with a non-exponential distribution but the same mean. The three concepts reversibility, insensitivity and

product form are intimately related. There is by now a huge literature, but some issues remain complicated.

Theorem 9.2. (*insensitivity of $M/GI/N/0$ model*) *The steady-state distribution of the number of customers in the $M/GI/N/0$ model depends on the general service-time distribution only through its mean.*

We apply a generalization of Theorem 6.1 (Kelly's lemma) to the case of a non-discrete state space. We guess the steady-state distribution of the Markov process (to be specified below) and the reverse-time Markov process, and then show that the detailed balance equation holds. By the same elementary reasoning that completes the proof provided that the generalization of the equation $\alpha Q = 0$ is still valid. The validity of such characterizations of the steady-state distribution gets much more complicated with non-discrete state spaces, but the ideas are similar. See §4.9 of Ethier and Kurtz (1986) for supporting details for that step. (In that greater generality, the rate matrix Q is the **generator**, often denoted by A . For the general state space, it is assumed that A generates a strongly continuous contraction semi-group on an appropriate space and the martingale problem for A is well posed. Then the existence of a stationary distribution α is equivalent to $\alpha Q = 0$, expressed as $\int Af d\alpha = 0$ for all f in the core of A .) See Whitt (1980) for an alternative continuity argument to establish the insensitivity and Burman (1981) for a direct argument using generators.

The main line of reasoning is the same as we have been using. The process describing the system state can be made a continuous-time Markov process (but with uncountable state space) if we include the ages of the customers in service. A state of this Markov process can be written as

$$(n; x) \equiv (n; (x_1, x_2, \dots, x_n)), \quad (9.16)$$

where, without loss of generality, we take the ages to be expressed in increasing order, so that $x_1 < x_2 < \dots < x_n$.

The steady-state distribution turns out to be just what we would hope for: The steady-state number in system is the same as for the $M/M/N/0$ model. Given that there are n customers in service, the ages of these service times in process (and the residual lifetimes) in steady state are independently distributed, each with cdf G_e , the equilibrium excess distribution associated with the service-time cdf G , which has density

$$g_e(x) \equiv \frac{G^c(x)}{ES}, \quad (9.17)$$

where $G^c(x) \equiv 1 - G(x)$. Taking account of the order requirement on the ages, which introduces the factor $n!$, we obtain the steady state distribution of the Markov process.

Theorem 9.3. (*steady-state of $M/GI/N/0$ Markov process*) *For the $M/GI/N/0$ model, if the service-time cdf has a pdf g , then the steady-state distribution is*

$$\alpha(n; x) \equiv \alpha(n; (x_1, x_2, \dots, x_n)) = \alpha_{M/M/N/0}(n) \times n! \prod_{i=1}^n g_{ie}(x_i). \quad (9.18)$$

Proof. We guess that the claim is correct and identify the reverse-time Markov process and apply the generalization of Kelly's lemma. The reverse-time Markov process is also an $M/GI/N/0$ Markov process with Poisson arrival process, but in reverse time the forward-time ages become residual lifetimes or excess variables. There are two possible events: (i) an arrival or (ii) a service completion. An arrival takes the state $(n; x)$ into $(n + 1; (0, x))$ for x

of dimension n . A service completion takes the state $(n; x)$ into $(n - 1; e_i(x))$, where x is of dimension n and $e_i(x)$ is the $n - 1$ -dimension vector with the i^{th} component of x removed. Let the hazard rate of the service-time cdf be

$$h(x) \equiv \frac{g(x)}{G^c(x)}, \quad x > 0. \quad (9.19)$$

The two detailed flow equations are

$$\alpha(n; x)\lambda = \alpha(n + 1; (0, x)) \quad (9.20)$$

$$\alpha(n; x)h(x_i) = \alpha(n - 1; e_i(x))\lambda g(x_i) \quad (9.21)$$

These can be seen to hold, as we now show. First, for (9.20), we observe that the left side is the flow rate from state $(n; x)$ to state $(n + 1; (0, x))$ due to an arrival, while the right side is the flow rate of the reverse time process in steady state in the opposite direction. Recall that, for the reverse time process, the times are interpreted as remaining times until the next service completion. Hence, a service completion is imminent only if the minimal component of the vector is 0. Hence, $\alpha(n + 1; (0, x))$ is indeed the rate of flow of the reverse time process from the state $(n + 1; (0, x))$ due to a service completion, which takes the system to state $(n; x)$.

It remains to show that the equation is satisfied for our candidate steady-state distribution. To see that it is, write $\alpha(n + 1; (0, x))$ as $\alpha(n; x)$ multiplied by other terms. From above, we see that

$$\alpha(n + 1; (0, x)) = \left(\frac{\lambda}{(n + 1)\mu} \right) (n + 1)g_e(0) \times \alpha(n; x), \quad (9.22)$$

but

$$g_e(0) = \mu G^c(0) = \mu \quad \text{and} \quad \left(\frac{\lambda}{(n + 1)\mu} \right) (n + 1)\mu = \lambda, \quad (9.23)$$

so that (9.20) holds.

We now turn to (9.21). First, observe that the left side is the flow rate from state $(n; x)$ to state $(n - 1; e_i(x))$ due to a service completion of the customer with age x_i . In state $(n; x)$, the intensity of such a service completion is precisely $h(x_i)$. On the other hand, the right side is the flow rate from state $(n - 1; e_i(x))$ to state $(n; x)$ in the reverse-time process due to an arrival of a new customer with required service time x_i . We then see that equality holds in (9.21) because

$$\frac{\alpha(n; x)}{\alpha(n - 1; e_i(x))} = \left(\frac{\lambda}{n\mu} \right) n g_e(x_i) = \frac{\lambda g_e(x_i)}{\mu} = \lambda G^c(x_i). \quad (9.24)$$

Then note that

$$\frac{\lambda g(x_i)}{h(x_i)} = \lambda G^c(x_i). \quad (9.25)$$

Combining (9.24) and (9.25), we see that indeed (9.21) holds. ■

10. Regularity and Irregularity in Infinite-State CTMC's

Most of the theory for finite-state CTMC's extends to infinite-state CTMC's, but regularity conditions are required. Otherwise, the infinite state space introduces complications. See Chung (1967) for an early account and Chapter II and §III.2 of Asmussen (2003) for a recent account. We give an overview of the theory, without giving any proofs.

10.1. Motivating Examples

Consider CTMC's on the nonnegative integers defined by transition rates $Q_{i,j}$, where we specify $Q_{i,j}$ for all j with $j \neq i$, for all i . Describe the evolution of the system in each case.

Example 10.1. (*moving instantaneously*) Let Q be defined on the nonnegative integers by

$$Q_{i,j} = i + j \quad \text{for all } i \geq 0 \quad \text{and } j \geq -i \quad \text{with } j \neq 0.$$

Example 10.2. (*moving up VERY fast*) Consider the pure-birth process on the nonnegative integers defined by

$$\lambda_i \equiv Q_{i,i+1} = (i+1)^2 \quad \text{for all } i \geq 0.$$

Example 10.3. (*moving up relentlessly*) Consider the pure-birth process on the nonnegative integers defined by

$$\lambda_i \equiv Q_{i,i+1} = (i+1) \quad \text{for all } i \geq 0.$$

Example 10.4. (*moving up and down rapidly for large values*) Consider the BD process on the nonnegative integers defined by

$$\lambda_i \equiv Q_{i,i+1} = Q_{i,i-1} = (i+1) \equiv \mu_i \quad \text{for all } i \geq 0.$$

Example 10.5. (*moving up and down even more rapidly for large values*) Consider the BD process on the nonnegative integers defined by

$$\lambda_i \equiv Q_{i,i+1} = Q_{i,i-1} \equiv \mu_i = (i+1)^2 \quad \text{for all } i \geq 0.$$

Example 10.6. (*subtle variation*) Consider the BD process on the nonnegative integers defined by

$$\lambda_i \equiv Q_{i,i+1} = (i+1)^4 \quad \text{and} \quad \mu_i \equiv Q_{i,i-1} = i^2(i+1)^2 \quad \text{for all } i \geq 0.$$

Example 10.1 is really bad; the total rate of transitions from every state is infinite. We have instantaneous transitions. We avoid that by assuming that

$$\sum_{j,j \neq i} Q_{i,j} < \infty \quad \text{for all } i.$$

If that condition is satisfied, then we call the CTMC a jump process. None of the other examples have that problem. However, Example 10.2 has an explosion. That CTMC has infinitely many transitions in finite time. It explodes in finite time. Example 10.3 is better; the CTMC diverges to infinity, but it does so over time. We say the CTMC is regular if it does not explode, i.e., if it has only finitely many transitions in finite time.

The remaining examples are regular pure jump BD processes. Example 10.4 is recurrent but not positive recurrent, because $\sum_{k=0}^{\infty} r_k = \infty$. (Here $r_k = 1/(k+1)$.) In contrast, Example 10.5 is positive recurrent, because $\sum_{k=0}^{\infty} r_k < \infty$. (Here $r_k = 1/(k+1)^2$.) Corollary 10.1 below shows that both these BD processes are regular.

We conclude with the observation that the story can be genuinely complicated. We can have $\sum_{k=0}^{\infty} r_k < \infty$ without the process being positive recurrent or even recurrent. That is illustrated by Example 10.6. Again Corollary 10.1 below shows that this BD process is regular. However, Theorem 10.7 below shows that this BD process is not recurrent and then of course not positive recurrent either. (Here we have $r_k = 1/(k+1)^2$ but $\lambda_k = (k+1)^4$ and $\lambda_k r_k = (k+1)^2$.) Incidentally, Theorem 10.7 below shows that the BD processes in Examples 10.4 and 10.5 are recurrent because $\lambda_k r_k = 1$ for all $k \geq 0$.

10.2. Instantaneous Transitions, Explosions and the Minimal Construction

Instantaneous transitions occur in 0 time. This pathology is avoided by direct assumption. We directly assume that for each initial state i that the process remains in state i for an exponential length of time with positive finite mean. That is achieved by assuming that

$$-Q_{i,i} \equiv \sum_{j,j \neq i} Q_{i,j} < \infty \quad \text{for all } i. \quad (10.1)$$

Condition (10.1) is trivially satisfied if the state space is finite, but not otherwise. If we assume condition (10.1), which we henceforth do, then say we are considering a Markov **jump process**. Then the time spent in state i on each visit is exponentially distributed with mean $1/|Q_{i,i}|$.

It is important that a pure-jump Markov process always has the **strong Markov property**. (As with all these forms of weirdness, it is most important to realize that exceptions actually exist, without the special structure.) The Markov property states that the probability of each future event given the present state and additional information about the past, even the entire past history, depends only on the present state. The strong Markov property states that the Markov property still holds for random times, provided that those random times are **stopping times**. A random time T is a stopping time relative to a stochastic process $\{X(t) : t \geq 0\}$ if the event $\{T \leq t\}$ depends only on $\{X(s) : s \leq t\}$ for each $t > 0$. See §I.1 and §II.1 of Asmussen (2003). The strong Markov property states that the probability of each future event given the state at a random stopping time and additional information about the past (prior to that random time), even the entire past history, depends only on the present state at that random time. The detailed argument depends on the precise definition of conditional expectation and conditional probability; e.g. see §9.1 of Chung (2001).

Explosions occur when the process has infinitely many transitions in finite time. For Markov jump processes, this can only occur if the process diverges to infinity in finite time; we call that an **explosion**. An explosion is possible, even for Markov jump processes. A CTMC is called **regular** if the number of transitions from each state in finite time is finite with probability 1; i.e., a regular CTMC has no explosions.

For a pure birth process, divergence to infinity turns out to occur if and only if the expected time to go from any state to infinity is finite.

Theorem 10.1. (*explosions in a pure birth process*) *A pure birth process is explosive (diverges to infinity in finite time) if and only if the mean time to diverge is finite, i.e.,*

$$\sum_{n=1}^{\infty} (1/\lambda_n) < \infty. \quad (10.2)$$

Example 10.7. *Here are simple examples: A pure birth process is explosive is $\lambda_n = cn^2$, $n \geq 1$, but not if $\lambda_n = cn$, $n \geq 1$.*

The CTMC can be well defined even with explosions. The **minimal construction** for a CTMC is based on exponential holding times in each state and a DTMC transition matrix at transition times, but with the added feature that the process is absorbed in an extra “death state” Δ after infinitely many transitions have taken place. Let S_n be the time of the n^{th} transition. Let

$$S_{\infty} \equiv \sup_{n \geq 1} S_n, \quad (10.3)$$

which is less than or equal to infinity. let the process be defined by

$$X(t) \equiv \Delta \quad \text{if } t \geq S_{\infty}. \quad (10.4)$$

Theorem 10.2. (*Kolmogorov ODE's for minimal construction*) *The minimal construction yields a solution to the Kolmogorov forward and backward ODE's.*

The extra state Δ plays no role in a regular CTMC.

10.3. Conditions for Regularity and Recurrence

We are considering irreducible pure-jump CTMC's. First we want conditions for regularity; i.e., we want to ensure that there are no explosions.

Theorem 10.3. (*Reuter's condition*) *A pure-jump CTMC is regular if and only if the only nonnegative bounded solution y to the matrix equation*

$$Qy = y \tag{10.5}$$

is the 0 vector, with $y_j = 0$ for all j .

See Proposition II.3.3 on p. 47 of Asmussen (2003).

Theorem 10.4. (*Kolmogorov ODE's*) *The Kolmogorov forward and backward ODE's have unique identical solutions for any regular pure-jump CTMC.*

Let Y_n be the state of the CTMC at the n^{th} transition. The following is Proposition II.2.3 of Asmussen (2003).

Theorem 10.5. (*regularity*) *An irreducible pure-jump CTMC is not regular if and only if*

$$\sum_{n=1}^{\infty} |Q_{Y_n, Y_n}| < \infty.$$

Corollary 10.1. (*sufficient conditions for regularity*) *Sufficient conditions for regularity are (i) $\sup_i |Q_{i,i}| < \infty$ or (ii) the DTMC $\{Y_n : n \geq 0\}$ is recurrent.*

Proof. Observe that $|Q_{Y_n, Y_n}| \rightarrow \infty$ if the process is not regular. If the process is not regular, then $|Q_{Y_n, Y_n}| \rightarrow \infty$ as $n \rightarrow \infty$. If $\{Y_n\}$ is recurrent and $Y_0 = i$, then $Q_{i,i}$ is a limit point of the sequence $\{|Q_{Y_n, Y_n}|\}$, which prevents $|Q_{Y_n, Y_n}| \rightarrow \infty$ as $n \rightarrow \infty$. ■

But we also want a proper steady-state distribution. A regular CTMC is **recurrent** if the process returns to each state w.p.1 after it leaves. A recurrent CTMC is **positive recurrent** if the expected time to return to each state after it leaves is finite. We necessarily have the following implications.

$$\text{positive recurrence} \quad \Rightarrow \quad \text{recurrence} \quad \Rightarrow \quad \text{regularity.}$$

Theorem 10.6. (*positive recurrence and the stationary distribution*) *A regular irreducible CTMC is positive recurrent if and only if the equation $\alpha Q = 0$ has a unique solution α that is a probability vector. In that case, α is the unique stationary vector.*

The following coincides with Proposition III.2.1 of Asmussen (2003), but is expressed slightly differently.

Theorem 10.7. (*regularity and recurrence of BD processes*) An irreducible BD process is recurrent, and thus regular, if and only if

$$\sum_{n=0}^{\infty} (\lambda_n r_n)^{-1} = \infty, \quad (10.6)$$

where $r_0 \equiv 1$ and

$$r_n \equiv \frac{\lambda_0 \times \lambda_1 \times \cdots \times \lambda_{n-1}}{\mu_1 \times \mu_2 \times \cdots \times \mu_n}, \quad n \geq 1. \quad (10.7)$$

Corollary 10.2. (*sufficient condition for recurrence in BD processes*) If $\limsup_{n \rightarrow \infty} \lambda_n / \mu_n = c < 1$, then an irreducible BD process is recurrent, and thus regular.

Theorem 10.8. (*positive recurrence of BD processes*) An irreducible BD process is positive recurrent if and only if it is recurrent and

$$\sum_{n=0}^{\infty} r_n < \infty \quad (10.8)$$

for r_n defined above. In which case, there is a unique stationary distribution and limiting distribution, which coincide and satisfy

$$\alpha_j \equiv \lim_{t \rightarrow \infty} P(X(t) = j) = \frac{r_j}{\sum_{n=0}^{\infty} r_n}, \quad j \geq 0. \quad (10.9)$$

Corollary 10.3. (*finite-state BD processes*) An irreducible finite-state BD process is always positive recurrent.

Theorem 10.9. (*reversibility*) The stationary version (obtained by taking $P(X(0) = j) = \alpha_j$ for all j) of a positive recurrent irreducible BD process is time reversible, because

$$\alpha_i \lambda_i = \alpha_{i+1} \mu_{i+1} \quad \text{for all } i \geq 0. \quad (10.10)$$

11. More on Reversible CTMC's and Birth-and-Death Processes

In this section we discuss additional ways to exploit the structure of reversibility. Much more is in the cited books by Kelly (1979) and Keilson (1979) and others.

11.1. Spectral Representation in Reversible Markov Chains

Among Markov processes, reversible Markov processes are distinguished by symmetry in the transition probabilities that lead to explicit spectral representations associated with **spectral theory**. This structure is easy to see in finite-state Markov chains, involving finite matrices. In particular, time-reversible CTMC's such as BD processes have spectral representations where the eigenvalues and eigenvectors are all real. See Chapter 3 of Keilson (1979); see Aldous and Fill (2002) for advanced material. These spectral representations provide explicit representations for the transient transition probabilities and explicit representations for the rate of convergence to the steady-state probability vector.

There is a corresponding algebraic approach to DTMC's too. The Perron-Frobenius theory of positive matrices can be applied. The rough idea is to diagonalize the transition matrix P in DTMC's or $P(t)$ in CTMC's; i.e., for the DTMC with transition matrix P , we write

$$P = U \Lambda U^{-1}, \quad (11.1)$$

where the diagonal elements of Λ are the eigenvalues of P , while U and U^{-1} are made up of the associated eigenvectors. Then

$$P^n = U\Lambda^n U^{-1}, \quad n \geq 1. \quad (11.2)$$

Reversibility plays a critical role here, because this algebraic structure occurs in a relatively simple form for reversible Markov chains. For simplicity, suppose that the state space is finite of dimension m . We now give additional details, following Keilson (1979). First, we consider an irreducible DTMC with transition matrix P and unique stationary probability vector π , satisfying $\pi = \pi P$. Let π_D be the $m \times m$ diagonal matrix with i^{th} diagonal element π_i and all off-diagonal elements 0. Then observe the following.

Proposition 5. (*reversibility in P characterized by symmetry*) *An irreducible finite-state DTMC P is reversible if and only if the $m \times m$ matrix $S \equiv \pi_D P$ is a symmetric matrix; i.e., if and only if $S_{i,j} = S_{j,i}$ for all i and j .*

We then can apply the finite **spectral theorem** for real symmetric matrices. The finite-dimensional spectral theorem says that any symmetric matrix whose entries are real can be diagonalized by an orthogonal matrix, i.e., we can write (11.1) where U and U^{-1} are real and nonsingular with $UU^{-1} = I$ and Λ is a diagonal matrix with real eigenvalues. The columns of U are right eigenvector of P while the rows of U^{-1} are left eigenvectors of P . We can then apply the Perron-Frobenius theory to conclude that there is one eigenvalue taking the value 1 and all the other eigenvalues satisfy $|\lambda_i| < 1$. We thus have the explicit representations

$$P_{i,j} = \sum_{k=1}^m U_{i,k} U_{k,j}^{-1} \lambda_k \quad \text{and} \quad P_{i,j}^n = \sum_{k=1}^m U_{i,k} U_{k,j}^{-1} \lambda_k^n = \pi_j + \sum_{k=1}^{m-1} a_{i,j} \lambda_k^n, \quad (11.3)$$

where λ_k are real numbers and, in the last expression, $a_{i,j}$ are real numbers and $|\lambda_k| < 1$, so that the last expression gives an explicit rate of convergence to steady state. The rate of convergence is primarily determined by the **spectral gap**, how much smaller is the absolute value of the eigenvalue with the largest absolute value strictly less than 1; e.g., see Aldous and Fill (2002). For large n , that term will dominate.

A corresponding story holds for CTMC's. In particular, we have the following.

Theorem 11.1. (*rate of convergence to steady state*) *For a reversible positive recurrent irreducible CTMC, such as a BD process,*

$$P_{i,j}(t) = \alpha_j + \sum_{l=1}^{m-1} a_{i,j} e^{-r_l t}, \quad t \geq 0, \quad (11.4)$$

where α_j is the steady-state probability, $r_j > 0$ for all j and $a_{i,j}$ is a real number for all i and j .

Proof. We can apply the result for DTMC's in (11.3) above by using uniformization, as discussed in §3.4 of these notes. We can represent a CTMC as a Poisson randomization of a DTMC. We can write the CTMC as

$$X(t) = Y_{N(t)}, \quad t \geq 0, \quad (11.5)$$

where $\{Y_n : n \geq 0\}$ is a DTMC and $\{N(t) : t \geq 0\}$ is a Poisson process. We let the rate of the Poisson process be r , where $r > |Q_{i,i}|$ for all i . We let the transition matrix of the DTMC be $P_{i,j} \equiv Q_{i,j}/r$ for all $i \neq j$. That is we let $P = I + r^{-1}Q$. We then have

$$P_{i,j}(t) = \sum_{k=0}^{\infty} P_{i,j}^k \frac{e^{-rt}(rt)^k}{k!}, \quad (11.6)$$

as in (3.28) of the lecture notes. Hence the spectral representation for the DTMC carries over directly to the CTMC

$$\begin{aligned} P_{i,j}(t) &= \sum_{k=0}^{\infty} P_{i,j}^k \frac{e^{-rt}(rt)^k}{k!} \\ &= \sum_{k=0}^{\infty} \left(\sum_{l=1}^m U_{i,l} U_{l,j}^{-1} \lambda_l^k \right) \frac{e^{-rt}(rt)^k}{k!} \\ &= \sum_{l=1}^m U_{i,l} U_{l,j}^{-1} \sum_{k=0}^{\infty} \lambda_l^k \frac{e^{-rt}(rt)^k}{k!} \\ &= \sum_{l=1}^m U_{i,l} U_{l,j}^{-1} e^{-rt} \sum_{k=0}^{\infty} \frac{(\lambda_l rt)^k}{k!} \\ &= \sum_{l=1}^m U_{i,l} U_{l,j}^{-1} e^{-r(1-\lambda_l)t}. \end{aligned} \quad (11.7)$$

Since $|\lambda_i| \leq 1$ for all i and equals 1 for only one i , we have

$$P_{i,j}(t) = \alpha_j + \sum_{l=1}^{m-1} a_{i,j} e^{-r_l t}, \quad t \geq 0, \quad (11.8)$$

where $r_j > 0$ for all j and $a_{i,j}$ is a real number for all i and j . ■

11.2. Fitting BD Processes to Data

A BD process can be fit to data from any stochastic process that makes all transitions up one or down one. Let $T_i(t)$ be the total time spent in state i in the interval $[0, t]$; Let $A_i(t)$ be the number of transitions up one from state i in the interval $[0, t]$; and Let $D_i(t)$ be the number of transitions down one from state i in the interval $[0, t]$. Then define estimated birth and death rates by

$$\bar{\lambda}_i(t) \equiv \frac{A_i(t)}{T_i(t)} \quad \text{and} \quad \bar{\mu}_i(t) \equiv \frac{D_i(t)}{T_i(t)}. \quad (11.9)$$

Let the estimated steady-state distribution be

$$\bar{\alpha}_i(t) \equiv \frac{T_i(t)}{t}, \quad i \geq 0. \quad (11.10)$$

We might say that the BD model fits the data well if the estimated steady-state probability vector $\bar{\alpha}$ agrees closely with the theoretical steady-state distribution based on the estimated birth and death rates, using formula (10.9). However, **the estimated steady-state probability vector $\bar{\alpha}$ automatically agrees very closely with the theoretical steady-state probability vector based on the estimated birth and death rates.** See Whitt (2012) and references therein.

In particular, Theorem 1 of Whitt (2012) shows that $\bar{\alpha}$ coincides exactly with the theoretical steady-state probability vector based on the estimated birth and death rates if the system ends in the same state it starts. Otherwise, there is likelihood-ratio stochastic order (which implies ordinary stochastic order), depending on the ordering of the initial and final states. As illustrated by Corollary 4.1 in the cited reference, it is possible to show, under regularity conditions, that the difference between the two probability vectors goes to 0 as the amount of data increases. Note that this holds **without any model assumptions** beyond having all transitions be up one or down one. In particular, the model need not be Markovian and the behavior could be highly time-dependent. The arrival rate might be highly time-dependent, such as sinusoidal. Nevertheless, a BD model fit to the data will necessarily produce a steady-state distribution that matches the long-run average performance. Of course, the long-run average performance may not match what happens at any particular time.

11.3. Comparing BD processes

Two BD process can be compared using a sample-path stochastic ordering if the smaller one has lower birth rates and higher death rates. The sample path stochastic ordering is a special construction putting both processes on the same underlying probability space, making them highly dependent, but leaving their individual distributions as stochastic processes unchanged. The sample-path construction allows us to make strong stochastic comparisons.

The desired conclusion is some form of stochastic order; see Ch. 9 of Ross (1996), Müller and Stoyan (2002) and Whitt (1981) for background. One (real-valued) random variable X_1 is said to be **stochastically less than or equal to** another X_2 , denoted by $X_1 \leq_{st} X_2$, if

$$P(X_1 > c) \leq P(X_2 > c) \quad \text{for all } c.$$

The idea behind the sample path stochastic ordering is contained in the following basic **coupling** result.

Theorem 11.2. (*coupling for stochastic order*) *If $X_1 \leq_{st} X_2$, then there exist random variables $Y_i, i = 1, 2$, such that $P(Y_1 \leq Y_2) = 1$ and $Y_i \stackrel{d}{=} X_i$, i.e., Y_i is distributed the same as X_i , $i = 1, 2$.*

Proof. Let F_i be the cdf of X_i . Let U be a random variable uniformly distributed on $[0, 1]$ and let

$$Y_i \equiv F_i^{-1}(U) \equiv \inf \{t : F_i(t) > U\}. \quad (11.11)$$

Then, by Lemma 11.1 below, $Y_i \stackrel{d}{=} X_i, i = 1, 2$, and $P(Y_1 \leq Y_2) = 1$. ■

Lemma 11.1. (*generating a random variable*) *Let U be a random variable uniformly distributed on $[0, 1]$ and let F be a cdf. Then $F^{-1}(U)$ has cdf F .*

Proof. For simplicity, we shall only do the proof for the case in which the cdf F is continuous and strictly increasing, so that it has an inverse, F^{-1} , with the properties

$$F^{-1}(F(x)) = x \quad \text{and} \quad F(F^{-1}(t)) = t$$

for all x and t with $0 < t < 1$. Hence, we can use the random variable $F^{-1}(U)$ because, for all x ,

$$P(F^{-1}(U) \leq x) = P(U \leq F(x)) = F(x).$$

However, the result holds in general. ■

Corollary 11.1. (*alternative characterization of stochastic order*) Stochastic order $X_1 \leq_{st} X_2$, holds if and only if $E[g(X_1)] \leq E[g(X_2)]$ for all nondecreasing bounded real-valued functions g .

Proof. Apply Theorem 11.2 to get the random variables Y_i with $Y_i \stackrel{d}{=} X_i$, $i = 1, 2$, and $P(Y_1 \leq Y_2) = 1$. It is then immediate that $P(g(Y_1) \leq g(Y_2)) = 1$ for the nondecreasing function g , which in turn implies that $E[g(Y_1)] \leq E[g(Y_2)]$. However, since $Y_i \stackrel{d}{=} X_i$, $g(Y_i) \stackrel{d}{=} g(X_i)$ and $E[g(Y_i)] = E[g(X_i)]$. Thus, also $E[g(X_1)] \leq E[g(X_2)]$ for all nondecreasing bounded real-valued functions g . ■

For birth-and-death processes, we have the following analog of Theorem 11.2. This is a variant of the comparisons in Whitt (1981).

Theorem 11.3. (*sample path stochastic order for BD processes*) For $i = 1, 2$, let $\lambda_k^{(i)}$, $k \geq 0$ and $\mu_k^{(i)}$, $k \geq 1$, be birth and death rates for two BD processes $\{Y_i(t) : t \geq 0\}$. If $Y_1(0) \leq_{st} Y_2(0)$ and

$$\lambda_k^{(1)} \leq \lambda_k^{(2)} \quad \text{and} \quad \mu_k^{(1)} \geq \mu_k^{(2)} \quad \text{for all } k, \quad (11.12)$$

then there exists special versions of these BD processes $\{X^{(i)}(t) : t \geq 0\}$ constructed on the same sample space such that each separately has the correct probability law as a BD process, while

$$P(X^{(1)}(t) \leq X^{(2)}(t) \quad \text{for all } t \geq 0) = 1. \quad (11.13)$$

Proof. The idea now is to directly perform a sample path construction to get an ordering of the entire sample paths. We first apply Theorem 11.2 to get $X_1(0) \leq X_2(0)$ w.p.1, where these random variables have the proper initial distribution. Then, to do the rest of the construction, use thinning of a Poisson process whenever the two processes are in the same state i . Let potential transitions be generated with a Poisson process having rate $(\lambda_i^{(1)} \vee \lambda_i^{(2)}) + (\mu_i^{(1)} \vee \mu_i^{(2)})$. Make the upper process have a birth whenever the lower process has a birth; make the lower process have a death whenever the upper process has a death. In that way the two processes each have the given probability law, but the sample paths are ordered w.p.1. ■

Here is a typical corollary. Let $T_{j,k}^{(i)}$ be the first passage time from state j to state k in BD process i .

Corollary 11.2. (*stochastic order of first passage for BD processes*) Consider two BD processes with birth and death rates ordered as in (11.12). If $j \leq k$, then

$$T_{j,k}^{(2)} \leq_{st} T_{j,k}^{(1)}. \quad (11.14)$$

If $j \geq k$, then

$$T_{j,k}^{(1)} \leq_{st} T_{j,k}^{(2)}. \quad (11.15)$$

Proof. We establish only (8.15). Apply Theorem 11.3 to get (8.14). If $j \leq k$, then for that special construction get

$$P(T_{j,k}^{(2)} \leq T_{j,k}^{(1)}) = 1. \quad (11.16)$$

As an immediate consequence, get

$$P(T_{j,k}^{(2)} > a) \leq P(T_{j,k}^{(1)} > a) \quad \text{for all } a, \quad (11.17)$$

which is equivalent to the stated conclusion. Notice that the last statement applies to the distribution of the two processes viewed separately; i.e., it no longer depends on the special construction. ■

Example 11.1. (*one fast server versus many slow servers*) We want to understand the efficiency tradeoff between a single-server queue having one fast servers as opposed to a many-server queue with many slow servers. Specifically, consider a $M/M/1$ and $M/M/s$ queueing models, both having a Poisson arrival process with rate λ . Let the individual service rate in the $M/M/s$ model be μ and let the individual service rate in the $M/M/1$ model be $s\mu$. In both models the traffic intensity is $\rho \equiv \lambda/s\mu$; assume that $\rho < 1$, so that the models are stable and steady-state distributions exist. Let $Q_k(t)$ be the number of customers in the system with k servers, where $k = 1$ or s . Suppose that the initial number of customers in the system is the same for both. We want to show that

$$Q_1(t) \leq_{st} Q_s(t) \quad \text{for all } t \geq 0. \quad (11.18)$$

Let $T_k(t)$ be the time that an arrival spends in the system before completing service. We want to show that

$$E[T_1(t)] \leq E[T_s(t)] \quad \text{for all } t \geq 0. \quad (11.19)$$

Answer: Observe that we can apply Theorem 11.3 to get $P(Q_1(t) \leq Q_s(t)) = 1$ for all $t = 1$, because the birth rates are identical and the death rates are ordered

$$\mu_k^{(1)} = s\mu > (k \wedge s)\mu = \mu_k^{(s)}, \quad k \geq 0.$$

To treat the expected time in the system, note that they are ordered too, starting from the same number k an arrival finds in the system:

$$E[T_1(t)|Q_1(t) = k] = \frac{k+1}{s\mu},$$

while

$$E[T_s(t)|Q_s(t) = k] = \frac{1}{\mu}, \quad k \leq s-1; \quad E[T_s(t)|Q_s(t) = k] = \frac{(k-s+1)}{s\mu} + \frac{1}{\mu}, \quad k \geq s.$$

For $k < s-1$,

$$E[T_1(t)|Q_1(t) = k] = \frac{k+1}{s\mu} < \frac{1}{\mu} = E[T_s(t)|Q_s(t) = k];$$

for $k \geq s-1$, we have equality. For $k \geq s$,

$$E[T_s(t)|Q_s(t) = k] = E[T_s(t)|Q_s(t) = k] = \frac{k+1}{s\mu}.$$

11.4. First-Passage Times in BD Processes

The first passage time from state i to state j , $T_{i,j}$, can be expressed as the sum of the first passage times to neighboring states. If $i < j$ then

$$T_{i,j} = T_{i,i+1} + T_{i+1,i+2} + \cdots + T_{j-1,j} \quad (11.20)$$

If $i > j$ then

$$T_{i,j} = T_{i,i-1} + T_{i-1,i-2} + \cdots + T_{j+1,j} \quad (11.21)$$

where the sum in each case is over independent random variables.

The first passage time up to the nearest neighbor has a relatively simple form, because it suffices to consider a finite-state absorbing CTMC; e.g., see Keilson (1979). The first passage time down is more complicated with an infinite state space. **Its Laplace transform can**

be computed using continued fractions, which can be used to calculate the distribution, using numerical inversion of Laplace transforms; see Abate and Whitt (1999) (The numerical inversion algorithm requires computing the Laplace transform for a modest number of complex arguments, e.g., about 50. Those required transform values can in turn be computed by algorithms for calculating (infinite) continued fractions.)

We now give the construction of the Laplace transform recursion for first passage times down. For $i \geq 1$, let T_i be the first passage time down from state i to state $i - 1$. Let $\hat{f}_i(s)$ be its Laplace transform, i.e.,

$$\hat{f}_i(s) \equiv E[e^{-sT_i}] = \int_0^\infty e^{-sx} f_{T_i}(x) dx. \quad (11.22)$$

To develop a recursion, we consider the first transition from state i . With probability $\lambda_i/(\lambda_i + \mu_i)$, the process moves up to state $i + 1$; with probability $\mu_i/(\lambda_i + \mu_i)$, the process moves down to state $i - 1$. If the process moves down then the first passage is complete. If the process moves up, then it must move down to i from $i + 1$ and then move down from i to $i - 1$. Recall that the time of the first transition is independent of the location (basic property of minimum of two independent exponential random variables). Let L be the location of the first transition from i and let T be the time of that transition. Thus, as in (4.7) of Abate and Whitt (1999),

$$\begin{aligned} \hat{f}_i(s) &\equiv E[e^{-sT_i}] = P(L = i - 1)E[e^{-sT}] + P(L = i + 1) \left(E[e^{-sT}] \hat{f}_{i+1}(s) \hat{f}_i(s) \right) \\ &= \left(\frac{\mu_i}{\lambda_i + \mu_i} \right) \left(\frac{\lambda_i + \mu_i}{\lambda_i + \mu_i + s} \right) \\ &\quad \left(\frac{\lambda_i}{\lambda_i + \mu_i} \right) \left(\frac{\lambda_i + \mu_i}{\lambda_i + \mu_i + s} \right) \hat{f}_{i+1}(s) \hat{f}_i(s), \end{aligned} \quad (11.23)$$

from which we obtain, by simple algebra, the recursive relation

$$\hat{f}_i(s) = \frac{\mu_i}{\lambda_i + \mu_i + s - \lambda_i \hat{f}_{i+1}(s)}. \quad (11.24)$$

This would be a finite recursion if there were only finitely many states, but the recursion never ends if there are infinitely many states. Nevertheless, for conventional BD processes, we expect a finite limit.

To put this in one of the standard continued fraction (CF) representations, we want no constant factor before $\hat{f}_{i+1}(s)$ in (11.24). To understand the basic recursion, write out the next step:

$$\hat{f}_i(s) = \frac{\mu_i}{\lambda_i + \mu_i + s - \frac{\lambda_i \mu_{i+1}}{\lambda_{i+1} + \mu_{i+1} + s - \lambda_{i+1} \hat{f}_{i+2}(s)}}. \quad (11.25)$$

and the step after that

$$\hat{f}_i(s) = \frac{\mu_i}{\lambda_i + \mu_i + s - \frac{\lambda_i \mu_{i+1}}{\lambda_{i+1} + \mu_{i+1} + s - \frac{\lambda_{i+1} \mu_{i+2}}{\lambda_{i+2} + \mu_{i+2} + s - \lambda_{i+2} \hat{f}_{i+3}(s)}}}. \quad (11.26)$$

To obtain a clean orderly representation, we rewrite the last version as

$$\hat{f}_i(s) = -\frac{1}{\lambda_{i-1}} \left(\frac{-\lambda_{i-1} \mu_i}{\lambda_i + \mu_i + s + \frac{-\lambda_i \mu_{i+1}}{\lambda_{i+1} + \mu_{i+1} + s + \frac{-\lambda_{i+1} \mu_{i+2}}{\lambda_{i+2} + \mu_{i+2} + s - \lambda_{i+2} \hat{f}_{i+3}(s)}}} \right). \quad (11.27)$$

We then express the result as

$$w_i = c_i \Phi_{n=i}^{\infty} \frac{a_n}{b_n} \quad \text{or} \quad w_i = c_i \left(\frac{a_i}{b_i +} \frac{a_{i+1}}{b_{i+1} +} \frac{a_{i+2}}{b_{i+2} +} \frac{a_{i+3}}{b_{i+3} +} \dots \right) \quad (11.28)$$

for

$$c_i \equiv -\frac{1}{\lambda_{i-1}}, \quad a_n \equiv -\lambda_{n-1}\mu_n \quad \text{and} \quad b_n \equiv \lambda_n + \mu_n + s, \quad n \geq i. \quad (11.29)$$

For a CF (sometimes called a generalized CF, because a CF can be expressed in more than one way), we write

$$w = \Phi_{n=1}^{\infty} \frac{a_n}{b_n} \quad \text{or} \quad w = \frac{a_1}{b_1 +} \frac{a_2}{b_2 +} \frac{a_3}{b_3 +} \frac{a_4}{b_4 +} \dots \quad (11.30)$$

There is a relatively simple recursion for calculating the successive approximants due to Euler in 1737, namely,

$$w_n = \frac{P_n}{Q_n}, \quad (11.31)$$

where $P_0 = 0$, $P_1 = a_1$, $Q_0 = 1$, $Q_1 = b_1$ and

$$P_n = b_n P_{n-1} + a_n P_{n-2} \quad \text{and} \quad Q_n = b_n Q_{n-1} + a_n Q_{n-2}, \quad n \geq 2. \quad (11.32)$$

Example 11.2. A continued fraction for π is (11.28) with

$$a_1 = 4, \quad b_1 = 1, \quad a_n = (n-1)^2 \quad \text{and} \quad b_n = 2 \quad \text{for all } n \geq 2. \quad (11.33)$$

A continued fraction for e is $b_0 + w$ for w in (11.29) with $b_0 = 2$,

$$a_n = 1 \quad \text{for } n \geq 1 \quad \text{and} \quad \{b_k : k \geq 1\} = \{1, 2, 1, 1, 4, 1, 1, 6, 1, 1, 8, 1, 1, 10, 1, 1, 12, \dots\} \quad (11.34)$$

This is sequence A003417 in the online encyclopedia of sequences (OEIS).

12. Some Next Steps

There is much more material in the references. Among the important next steps are **diffusion approximations** and other **asymptotic methods** in which the scale of the model grows, via large customer populations, multiple classes of customers or high dimension; e.g., many queues in a queueing network; e.g., see Pang et al. (2007), Whitt (1985) and Kelly (1991).

References

- [1] Abate, J., W. Whitt. 1999. Computing Laplace Transforms for Numerical Inversion Via Continued Fractions. *INFORMS Journal on Computing* 11, 394–405.
- [2] Aldous, D. J., J. Fill. 2002 *Reversible Markov Chains and Random Walks on Graphs*, University of California Press.
- [3] Asmussen, S. 2003. *Applied Probability and Queues*, second edition, Springer.
- [4] Baskett, F., K. M. Chandy, R. R. Muntz and F. G. Palacios. 1975. Open, closed and mixed networks of queues with different classes of customers. *J. ACM* 22, 248–260.
- [5] Burman, D. Y. 1981. Insensitivity in Queueing Systems. *Advances in Appl. Prob.* 13, 846–859.
- [6] Chen, H. and D. D. Yao. 2001. *Fundamentals of Queueing Networks*, Springer.
- [7] Choudhury, G. L., K. K. Leung and W. Whitt. 1995a. Calculating Normalization Constants of Closed Queueing Networks by Numerically Inverting Their Generating Functions. *J. ACM* 42, 935–970.
- [8] Choudhury, G. L., K. K. Leung and W. Whitt. 1995b. An Algorithm to Compute Blocking Probabilities in Multi-Rate Multi-Class Multi-Resource Loss Models. *Advances in Appl. Prob.* 27, 1104–1143.
- [9] Chung, K. L. 1967. *Markov Chains with Stationary Transition Probabilities*, second edition, Springer.
- [10] Chung, K. L. 2001. *A Course in Probability Theory*, third edition, Springer.
- [11] El Taha, M. and S. Stidham, Jr. 1999. *Sample-Path Analysis of Queueing Systems*, Kluwer.
- [12] Ethier, S. N. and T. G. Kurtz. 1986. *Markov Processes: Characterization and Convergence*, Wiley.
- [13] Glynn, P. W. 1989. A GSMP formalism for discrete-event systems. *Proceedings IEEE* 77, 14–23.
- [14] Jagerman, D. L. 1974. Some properties of the Erlang loss function. *Bell System Technical Journal* 53, 525–551.
- [15] Keilson, J. 1979. *Markov Chain Models – Rarity and Exponentiality*, Springer.
- [16] Kelly, F. P. 1979. *Reversibility and Stochastic Networks*, Wiley.
- [17] Kelly, F. P. 1991. Loss networks. *Ann. Appl. Prob.* 1, 319–378.
- [18] Melamed, B. and W. Whitt. 1990. On arrivals that see time averages. *Operations Research* 38, 156–172.
- [19] Moler, C. and C. Van Loan. 2003. Nineteen dubious ways to compute the exponential of a matrix, twenty-five years later. *SIAM Review* 45, 3–49.
- [20] Müller, A., D. Stoyan. 2002. *Comparison Methods for Stochastic Models and Risks*, Wiley.

- [21] Pang, G., R. Talreja and W. Whitt. 2007. Martingale Proofs of Many-Server Heavy-Traffic Limits for Markovian Queues. *Probability Surveys* 4, 193–267.
- [22] Ross, S. M. 2010. *Introduction to Probability Models*, tenth ed., Academic Press.
- [23] Ross, S. M. 1996. *Stochastic Processes*, second ed., Wiley.
- [24] Sauer, C. H. and K. M. Chandy. 1981. *Computer Systems Performance Modeling*, Prentice Hall.
- [25] Serfozo, R. F. 1999. *Introduction to Stochastic Networks*, Springer.
- [26] Simmons, G. F. 1991. *Differential Equations with Applications and Historical Notes*, second edition, McGraw-Hill.
- [27] van Dijk, N. 1993. *Queueing Networks and Product Forms*, Wiley.
- [28] Walrand, J. 1988. *Introduction to Queueing Networks*, Prentice-Hall.
- [29] Whitt, W. 1980. Continuity of Generalized Semi-Markov Processes. *Mathematics of Operations Research* 5, 494–501.
- [30] Whitt, W. 1981. Comparing counting processes and queues. *Advances Appl. Prob.* 13, 207–220.
- [31] Whitt, W. 1984. Open and Closed Models of Networks of Queues. *AT&T Bell laboratories Tech. J.* 63, 1911–1979.
- [32] Whitt, W. 1985. Blocking When Service is Required from Several Facilities Simultaneously. *AT&T Technical Journal* 64, 1807–1856.
- [33] Whitt, W. 2002. The Erlang B and C formulas: problems and solutions. Lecture notes. Available at: <http://www.columbia.edu/~ww2040/allpapers.html>
- [34] Whitt, W. 2012. Fitting Birth-and-Death Queueing Models to Data. *Statistics and Probability Letters*, 82, 998–1004.
- [35] Whittle, P. 1986. *Systems in Stochastic Equilibrium*, Wiley.
- [36] Wolff, R. W. 1989. *Stochastic Modeling and the Theory of Queues*, Prentice-Hall.