

---

## Chapter 8

### Queueing Theory

---

#### 8.1. Introduction

In this chapter we will study a class of models in which customers arrive in some random manner at a service facility. Upon arrival they are made to wait in queue until it is their turn to be served. Once served they are generally assumed to leave the system. For such models we will be interested in determining, among other things, such quantities as the average number of customers in the system (or in the queue) and the average time a customer spends in the system (or spends waiting in the queue).

In Section 8.2 we derive a series of basic queueing identities which are of great use in analyzing queueing models. We also introduce three different sets of limiting probabilities which correspond to what an arrival sees, what a departure sees, and what an outside observer would see.

In Section 8.3 we deal with queueing systems in which all of the defining probability distributions are assumed to be exponential. For instance, the simplest such model is to assume that customers arrive in accordance with a Poisson process (and thus the interarrival times are exponentially distributed) and are served one at a time by a single server who takes an exponentially distributed length of time for each service. These exponential queueing models are special examples of continuous-time Markov chains and so can be analyzed as in Chapter 6. However, at the cost of a (very) slight amount of repetition we shall not assume the reader to be familiar with the material of Chapter 6, but rather we shall redevelop any needed material. Specifically we shall derive anew (by a heuristic argument) the formula for the limiting probabilities.

In Section 8.4 we consider models in which customers move randomly among a network of servers. The model of Section 8.4.1 is an open system in which customers are allowed to enter and depart the system, whereas the one studied in Section 8.4.2 is closed in the sense that the set of customers in the system is constant over time.

In Section 8.5 we study the model  $M/G/1$ , which while assuming Poisson arrivals, allows the service distribution to be arbitrary. To analyze this model we first introduce in Section 8.5.1 the concept of work, and then use this concept in Section 8.5.2 to help analyze this system. In Section 8.5.3 we derive the average amount of time that a server remains busy between idle periods.

In Section 8.6 we consider some variations of the model  $M/G/1$ . In particular in Section 8.6.1 we suppose that bus loads of customers arrive according to a Poisson process and that each bus contains a random number of customers. In Section 8.6.2 we suppose that there are two different classes of customers—with type 1 customers receiving service priority over type 2.

In Section 8.7 we consider a model with exponential service times but where the interarrival times between customers is allowed to have an arbitrary distribution. We analyze this model by use of an appropriately defined Markov chain. We also derive the mean length of a busy period and of an idle period for this model.

In the final section of the chapter we talk about multiserver systems. We start with loss systems, in which arrivals, finding all servers busy, are assumed to depart and as such are lost to the system. This leads to the famous result known as Erlang's loss formula, which presents a simple formula for the number of busy servers in such a model when the arrival process is Poisson and the service distribution is general. We then discuss multiserver systems in which queues are allowed. However, except in the case where exponential service times are assumed, there are very few explicit formulas for these models. We end by presenting an approximation for the average time a customer waits in queue in a  $k$ -server model which assumes Poisson arrivals but allows for a general service distribution.

## 8.2. Preliminaries

In this section we will derive certain identities which are valid in the great majority of queueing models.

### 8.2.1. Cost Equations

Some fundamental quantities of interest for queueing models are

$L$ , the average number of customers in the system;

$L_Q$ , the average number of customers waiting in queue;

$W$ , the average amount of time a customer spends in the system;

$W_Q$ , the average amount of time a customer spends waiting in queue.

A large number of interesting and useful relationships between the preceding and other quantities of interest can be obtained by making use of the following idea: Imagine that entering customers are forced to pay money (according to some rule) to the system. We would then have the following basic cost identity

$$\begin{aligned} &\text{average rate at which the system earns} \\ &= \lambda_a \times \text{average amount an entering customer pays} \end{aligned} \quad (8.1)$$

where  $\lambda_a$  is defined to be average arrival rate of entering customers. That is, if  $N(t)$  denotes the number of customer arrivals by time  $t$ , then

$$\lambda_a = \lim_{t \rightarrow \infty} \frac{N(t)}{t}$$

We now present an heuristic proof of Equation (8.1).

**Heuristic Proof of Equation (8.1)** Let  $T$  be a fixed large number. In two different ways, we will compute the average amount of money the system has earned by time  $T$ . On one hand, this quantity approximately can be obtained by multiplying the average rate at which the system earns by the length of time  $T$ . On the other hand, we can approximately compute it by multiplying the average amount paid by an entering customer by the average number of customers entering by time  $T$  (and this latter factor is approximately  $\lambda_a T$ ). Hence, both sides of Equation (8.1) when multiplied by  $T$  are approximately equal to the average amount earned by  $T$ . The result then follows by letting  $T \rightarrow \infty$ .\*

By choosing appropriate cost rules, many useful formulas can be obtained as special cases of Equation (8.1). For instance, by supposing that each customer pays \$1 per unit time while in the system, Equation (8.1) yields the so-called Little's formula,

$$L = \lambda_a W \quad (8.2)$$

\* This can be made into a rigorous proof provided we assume that the queueing process is regenerative in the sense of Section 7.5. Most models, including all the ones in this chapter, satisfy this condition.

This follows since, under this cost rule, the rate at which the system earns is just the number in the system, and the amount a customer pays is just equal to its time in the system.

Similarly if we suppose that each customer pays \$1 per unit time while in queue, then Equation (8.1) yields

$$L_Q = \lambda_a W_Q \quad (8.3)$$

By supposing the cost rule that each customer pays \$1 per unit time while in service we obtain from Equation (8.1) that the

$$\text{average number of customers in service} = \lambda_a E[S] \quad (8.4)$$

where  $E[S]$  is defined as the average amount of time a customer spends in service.

It should be emphasized that Equations (8.1) through (8.4) are valid for almost all queuing models regardless of the arrival process, the number of servers, or queue discipline.

### 8.2.2. Steady-State Probabilities

Let  $X(t)$  denote the number of customers in the system at time  $t$  and define  $P_n$ ,  $n \geq 0$ , by

$$P_n = \lim_{t \rightarrow \infty} P\{X(t) = n\}$$

where we assume the above limit exists. In other words,  $P_n$  is the limiting or long-run probability that there will be exactly  $n$  customers in the system. It is sometimes referred to as the *steady-state probability* of exactly  $n$  customers in the system. It also usually turns out that  $P_n$  equals the (long-run) proportion of time that the system contains exactly  $n$  customers. For example, if  $P_0 = 0.3$ , then in the long-run, the system will be empty of customers for 30 percent of the time. Similarly,  $P_1 = 0.2$  would imply that for 20 percent of the time the system would contain exactly one customer.\*

Two other sets of limiting probabilities are  $\{a_n, n \geq 0\}$  and  $\{d_n, n \geq 0\}$ , where

$a_n$  = proportion of customers that find  $n$   
in the system when they arrive, and

$d_n$  = proportion of customers leaving behind  $n$   
in the system when the depart

\* A sufficient condition for the validity of the dual interpretation of  $P_n$  is that the queuing process be regenerative.

That is,  $P_n$  is the proportion of time during which there are  $n$  in the system;  $a_n$  is the proportion of arrivals that find  $n$ ; and  $d_n$  is the proportion of departures that leave behind  $n$ . That these quantities need not always be equal is illustrated by the following example.

**Example 8.1** Consider a queuing model in which all customers have service times equal to 1, and where the times between successive customers are always greater than 1 [for instance, the interarrival times could be uniformly distributed over  $(1, 2)$ ]. Hence, as every arrival finds the system empty and every departure leaves it empty, we have

$$a_0 = d_0 = 1$$

However,

$$P_0 \neq 1$$

as the system is not always empty of customers. ♦

It was, however, no accident that  $a_n$  equaled  $d_n$  in the previous example. That arrivals and departures always see the same number of customers is always true as is shown in the next proposition.

**Proposition 8.1** In any system in which customers arrive one at a time and are served one at a time

$$a_n = d_n, \quad n \geq 0$$

**Proof** An arrival will see  $n$  in the system whenever the number in the system goes from  $n$  to  $n + 1$ ; similarly, a departure will leave behind  $n$  whenever the number in the system goes from  $n + 1$  to  $n$ . Now in any interval of time  $T$  the number of transitions from  $n$  to  $n + 1$  must equal to within 1 the number from  $n + 1$  to  $n$ . [For instance, if transitions from 2 to 3 occur 10 times, then 10 times there must have been a transition back to 2 from a higher state (namely, 3).] Hence, the rate of transitions from  $n$  to  $n + 1$  equals the rate from  $n + 1$  to  $n$ ; or, equivalently, the rate at which arrivals find  $n$  equals the rate at which departures leave  $n$ . The result now follows since the overall arrival rate must equal the overall departure rate (what goes in eventually goes out.) ♦

Hence, on the average, arrivals and departures always see the same number of customers. However, as Example 8.1 illustrates, they do not, in general, see the time averages. One important exception where they do is in the case of Poisson arrivals.

**Proposition 8.2** Poisson arrivals always see time averages. In particular, for Poisson arrivals,

$$P_n = a_n$$

To understand why Poisson arrivals always see time averages, consider an arbitrary Poisson arrival. If we knew that it arrived at time  $t$ , then the conditional distribution of what it sees upon arrival is the same as the unconditional distribution of the system state at time  $t$ . For knowing that an arrival occurs at time  $t$  gives us no information about what occurred prior to  $t$ . (Since the Poisson process has independent increments, knowing that an event occurred at time  $t$  does not affect the distribution of what occurred prior to  $t$ .) Hence, an arrival would just see the system according to the limiting probabilities.

Contrast the foregoing with the situation of Example 8.1 where knowing that an arrival occurred at time  $t$  tells us a great deal about the past; in particular it tells us that there have been no arrivals in  $(t - 1, t)$ . Thus, in this case, we cannot conclude that the distribution of what an arrival at time  $t$  observes is the same as the distribution of the system state at time  $t$ .

For a second argument as to why Poisson arrivals see time averages, note that the total time the system is in state  $n$  by time  $T$  is (roughly)  $P_n T$ . Hence, as Poisson arrivals always arrive at rate  $\lambda$  no matter what the system state, it follows that the number of arrivals in  $[0, T]$  that find the system in state  $n$  is (roughly)  $\lambda P_n T$ . In the long run, therefore, the rate at which arrivals find the system in state  $n$  is  $\lambda P_n$  and, as  $\lambda$  is the overall arrival rate, it follows that  $\lambda P_n / \lambda = P_n$  is the proportion of arrivals that find the system in state  $n$ .

## 8.3. Exponential Models

### 8.3.1. A Single-Server Exponential Queuing System

Suppose that customers arrive at a single-server service station in accordance with a Poisson process having rate  $\lambda$ . That is, the times between successive arrivals are independent exponential random variables having mean  $1/\lambda$ . Each customer, upon arrival, goes directly into service if the server is free and, if not, the customer joins the queue. When the server finishes serving a customer, the customer leaves the system, and the next customer in line, if there is any, enters service. The successive service times are assumed to be independent exponential random variables having mean  $1/\mu$ .

The above is called the  $M/M/1$  queue. The two  $M$ 's refer to the fact that both the interarrival and service distributions are exponential (and thus

memoryless, or Markovian), and the 1 to the fact that there is a single server. To analyze it, we shall begin by determining the limiting probabilities  $P_n$ , for  $n = 0, 1, \dots$ . To do so, think along the following lines. Suppose that we have an infinite number of rooms numbered  $0, 1, 2, \dots$ , and suppose that we instruct an individual to enter room  $n$  whenever there are  $n$  customers in the system. That is, he would be in room 2 whenever there are two customers in the system; and if another were to arrive, then he would leave room 2 and enter room 3. Similarly, if a service would take place he would leave room 2 and enter room 1 (as there would now be only one customer in the system).

Now suppose that in the long-run our individual is seen to have entered room 1 at the rate of ten times an hour. Then at what rate must he have left room 1? Clearly, at this same rate of ten times an hour. For the total number of times that he enters room 1 must be equal to (or one greater than) the total number of times he leaves room 1. This sort of argument thus yields the general principle which will enable us to determine the state probabilities. Namely, for each  $n \geq 0$ , the rate at which the process enters state  $n$  equals the rate at which it leaves state  $n$ . Let us now determine these rates. Consider first state 0. When in state 0 the process can leave only by an arrival as clearly there cannot be a departure when the system is empty. Since the arrival rate is  $\lambda$  and the proportion of time that the process is in state 0 is  $P_0$ , it follows that the rate at which the process leaves state 0 is  $\lambda P_0$ . On the other hand, state 0 can only be reached from state 1 via a departure. That is, if there is a single customer in the system and he completes service, then the system becomes empty. Since the service rate is  $\mu$  and the proportion of time that the system has exactly one customer is  $P_1$ , it follows that the rate at which the process enters state 0 is  $\mu P_1$ .

Hence, from our rate-equality principle we get our first equation,

$$\lambda P_0 = \mu P_1$$

Now consider state 1. The process can leave this state either by an arrival (which occurs at rate  $\lambda$ ) or a departure (which occurs at rate  $\mu$ ). Hence, when in state 1, the process will leave this state at a rate of  $\lambda + \mu$ .<sup>\*</sup> Since the proportion of time the process is in state 1 is  $P_1$ , the rate at which the process leaves state 1 is  $(\lambda + \mu)P_1$ . On the other hand, state 1 can be entered either from state 0 via an arrival or from state 2 via a departure. Hence, the rate at which the process enters state 1 is  $\lambda P_0 + \mu P_2$ . Because the reasoning

<sup>\*</sup> If one event occurs at rate  $\lambda$  and another occurs at rate  $\mu$ , then the total rate at which either event occurs is  $\lambda + \mu$ . Suppose one man earns \$2 per hour and another earns \$3 per hour; then together they clearly earn \$5 per hour.

for other states is similar, we obtain the following set of equations:

$$\begin{array}{ll} \text{State} & \text{Rate at which the process leaves} = \text{rate at which it enters} \\ 0 & \lambda P_0 = \mu P_1 \\ n, n \geq 1 & (\lambda + \mu)P_n = \lambda P_{n-1} + \mu P_{n+1} \end{array} \quad (8.5)$$

The set of Equations (8.5) which balances the rate at which the process enters each state with the rate at which it leaves that state is known as *balance equations*.

In order to solve Equations (8.5), we rewrite them to obtain

$$P_1 = \frac{\lambda}{\mu} P_0,$$

$$P_{n+1} = \frac{\lambda}{\mu} P_n + \left( P_n - \frac{\lambda}{\mu} P_{n-1} \right), \quad n \geq 1$$

Solving in terms of  $P_0$  yields

$$P_0 = P_0,$$

$$P_1 = \frac{\lambda}{\mu} P_0,$$

$$P_2 = \frac{\lambda}{\mu} P_1 + \left( P_1 - \frac{\lambda}{\mu} P_0 \right) = \frac{\lambda}{\mu} P_1 = \left( \frac{\lambda}{\mu} \right)^2 P_0,$$

$$P_3 = \frac{\lambda}{\mu} P_2 + \left( P_2 - \frac{\lambda}{\mu} P_1 \right) = \frac{\lambda}{\mu} P_2 = \left( \frac{\lambda}{\mu} \right)^3 P_0,$$

$$P_4 = \frac{\lambda}{\mu} P_3 + \left( P_3 - \frac{\lambda}{\mu} P_2 \right) = \frac{\lambda}{\mu} P_3 = \left( \frac{\lambda}{\mu} \right)^4 P_0,$$

$$P_{n+1} = \frac{\lambda}{\mu} P_n + \left( P_n - \frac{\lambda}{\mu} P_{n-1} \right) = \frac{\lambda}{\mu} P_n = \left( \frac{\lambda}{\mu} \right)^{n+1} P_0$$

To determine  $P_0$  we use the fact that the  $P_n$  must sum to 1, and thus

$$1 = \sum_{n=0}^{\infty} P_n = \sum_{n=0}^{\infty} \left( \frac{\lambda}{\mu} \right)^n P_0 = \frac{P_0}{1 - \lambda/\mu}$$

or

$$P_0 = 1 - \frac{\lambda}{\mu},$$

$$P_n = \left( \frac{\lambda}{\mu} \right)^n \left( 1 - \frac{\lambda}{\mu} \right), \quad n \geq 1 \quad (8.6)$$

Notice that for the preceding equations to make sense, it is necessary for  $\lambda/\mu$  to be less than 1. For otherwise  $\sum_{n=0}^{\infty} (\lambda/\mu)^n$  would be infinite and all the  $P_n$  would be 0. Hence, we shall assume that  $\lambda/\mu < 1$ . Note that it is quite intuitive that there would be no limiting probabilities if  $\lambda > \mu$ . For suppose that  $\lambda > \mu$ . Since customers arrive at a Poisson rate  $\lambda$ , it follows that the expected total number of arrivals by time  $t$  is  $\lambda t$ . On the other hand, what is the expected number of customers served by time  $t$ ? If there were always customers present, then the number of customers served would be a Poisson process having rate  $\mu$  since the time between successive services would be independent exponentials having mean  $1/\mu$ . Hence, the expected number of customers served by time  $t$  is no greater than  $\mu t$ ; and, therefore, the expected number in the system at time  $t$  is at least

$$\lambda t - \mu t = (\lambda - \mu)t$$

Now if  $\lambda > \mu$ , then the above number goes to infinity as  $t$  becomes large. That is,  $\lambda/\mu > 1$ , the queue size increases without limit and there will be no limiting probabilities. Note also that the condition  $\lambda/\mu < 1$  is equivalent to the condition that the mean service time be less than the mean time between successive arrivals. This is the general condition that must be satisfied for limited probabilities to exist in most single-server queueing systems.

Now let us attempt to express the quantities  $L$ ,  $L_Q$ ,  $W$ , and  $W_Q$  in terms of the limiting probabilities  $P_n$ . Since  $P_n$  is the long-run probability that the system contains exactly  $n$  customers, the average number of customers in the system clearly is given by

$$\begin{aligned} L &= \sum_{n=0}^{\infty} n P_n \\ &= \sum_{n=0}^{\infty} n \left( \frac{\lambda}{\mu} \right)^n \left( 1 - \frac{\lambda}{\mu} \right) \\ &= \frac{\lambda}{\mu - \lambda} \end{aligned} \quad (8.7)$$

where the last equation followed upon application of the algebraic identity

$$\sum_{n=0}^{\infty} n x^n = \frac{x}{(1-x)^2}$$

The quantities  $W$ ,  $W_Q$ , and  $L_Q$  now can be obtained with the help of Equations (8.2) and (8.3). That is, since  $\lambda_a = \lambda$ , we have from Equation (8.7)

that

$$\begin{aligned}
 W &= \frac{L}{\lambda} \\
 &= \frac{1}{\mu - \lambda}, \\
 W_Q &= W - E[S] \\
 &= W - \frac{1}{\mu} \\
 &= \frac{\lambda}{\mu(\mu - \lambda)}, \\
 L_Q &= \lambda W_Q \\
 &= \frac{\lambda^2}{\mu(\mu - \lambda)} \tag{8.8}
 \end{aligned}$$

**Example 8.2** Suppose that customers arrive at a Poisson rate of one per every 12 minutes, and that the service time is exponential at a rate of one service per 8 minutes. What are  $L$  and  $W$ ?

**Solution:** Since  $\lambda = \frac{1}{12}$ ,  $\mu = \frac{1}{8}$ , we have

$$L = 2, \quad W = 24$$

Hence, the average number of customers in the system is two, and the average time a customer spends in the system is 24 minutes.

Now suppose that the arrival rate increases 20 percent to  $\lambda = \frac{1}{10}$ . What is the corresponding change in  $L$  and  $W$ ? Again using Equations (8.7), we get

$$L = 4, \quad W = 40$$

Hence, an increase of 20 percent in the arrival rate *doubled* the average number of customers in the system.

To understand this better, write Equations (8.7) as

$$\begin{aligned}
 L &= \frac{\lambda/\mu}{1 - \lambda/\mu}, \\
 W &= \frac{1/\mu}{1 - \lambda/\mu}
 \end{aligned}$$

From these equations we can see that when  $\lambda/\mu$  is near 1, a slight increase in  $\lambda/\mu$  will lead to a large increase in  $L$  and  $W$ . ♦

**A Technical Remark** We have used the fact that if one event occurs at an exponential rate  $\lambda$ , and another independent event at an exponential rate  $\mu$ , then together they occur at an exponential rate  $\lambda + \mu$ . To check this formally, let  $T_1$  be the time at which the first event occurs, and  $T_2$  the time at which the second event occurs. Then

$$P\{T_1 \leq t\} = 1 - e^{-\lambda t},$$

$$P\{T_2 \leq t\} = 1 - e^{-\mu t}$$

Now if we are interested in the time until either  $T_1$  or  $T_2$  occurs, then we are interested in  $T = \min(T_1, T_2)$ . Now

$$\begin{aligned}
 P\{T \leq t\} &= 1 - P\{T > t\} \\
 &= 1 - P\{\min(T_1, T_2) > t\}
 \end{aligned}$$

However,  $\min(T_1, T_2) > t$  if and only if both  $T_1$  and  $T_2$  are greater than  $t$ ; hence,

$$\begin{aligned}
 P\{T \geq t\} &= 1 - P\{T_1 > t, T_2 > t\} \\
 &= 1 - P\{T_1 > t\}P\{T_2 > t\} \\
 &= 1 - e^{-\lambda t}e^{-\mu t} \\
 &= 1 - e^{-(\lambda + \mu)t}
 \end{aligned}$$

Thus,  $T$  has an exponential distribution with rate  $\lambda + \mu$ , and we are justified in adding the rates. ♦

Let  $W^*$  denote the amount of time an arbitrary customer spends in the system. To obtain the distribution of  $W^*$ , we condition on the number in the system when the customer arrives. This yields

$$\begin{aligned}
 P\{W^* \leq a\} &= \sum_{n=0}^{\infty} P\{W^* \leq a \mid n \text{ in the system when he arrives}\} \\
 &\quad \times P\{n \text{ in the system when he arrives}\} \tag{8.9}
 \end{aligned}$$

Now consider the amount of time that our customer must spend in the system if there are already  $n$  customers present when he arrives. If  $n = 0$ , then his time in the system will just be his service time. When  $n \geq 1$ , there will be one customer in service and  $n - 1$  waiting in line ahead of our arrival. The customer in service might have been in service for some time, but due to the lack of memory of the exponential distribution (see Section 5.2), it follows that our arrival would have to wait an exponential amount of time with rate  $\mu$  for this customer to complete service. As he also would have to wait an exponential amount of time for each of the other  $n - 1$

customers in line, it follows, upon adding his own service time, that the amount of time that a customer must spend in the system if there are already  $n$  customers present when he arrives is the sum of  $n + 1$  independent and identically distributed exponential random variables with rate  $\mu$ . But it is known (see Section 5.2.3) that such a random variable has a gamma distribution with parameters  $(n + 1, \mu)$ . That is,

$$P\{W^* \leq a \mid n \text{ in the system when he arrives}\} \\ = \int_0^a \mu e^{-\mu t} \frac{(\mu t)^n}{n!} dt$$

Because

$$P\{n \text{ in the system when he arrives}\} = P_n \quad (\text{since Poisson arrivals}) \\ = \left(\frac{\lambda}{\mu}\right)^n \left(1 - \frac{\lambda}{\mu}\right)$$

we have from Equation (8.9) and the preceding that

$$P\{W^* \leq a\} = \sum_{n=0}^{\infty} \int_0^a \mu e^{-\mu t} \frac{(\mu t)^n}{n!} dt \left(\frac{\lambda}{\mu}\right)^n \left(1 - \frac{\lambda}{\mu}\right) \\ = \int_0^a (\mu - \lambda) e^{-\mu t} \sum_{n=0}^{\infty} \frac{(\lambda t)^n}{n!} dt \quad (\text{by interchanging}) \\ = \int_0^a (\mu - \lambda) e^{-\mu t} e^{\lambda t} dt \\ = \int_0^a (\mu - \lambda) e^{-(\mu - \lambda)t} dt \\ = 1 - e^{-(\mu - \lambda)a}$$

In other words,  $W^*$ , the amount of time a customer spends in the system, is an exponential random variable with rate  $\mu - \lambda$ . (As a check, we note that  $E[W^*] = 1/(\mu - \lambda)$  which checks with Equation (8.8) since  $W = E[W^*]$ .)

**Remark** Another argument as to why  $W^*$  is exponential with rate  $\mu - \lambda$  is as follows. If we let  $N$  denote the number of customers in the system as seen by an arrival, then this arrival will spend  $N + 1$  service times in the system before departing. Now,

$$P\{N + 1 = j\} = P\{N = j - 1\} = (\lambda/\mu)^{j-1} (1 - \lambda/\mu), \quad j \geq 1$$

In words, the number of services that have to be completed before our arrival departs is a geometric random variable with parameter  $1 - \lambda/\mu$ . Therefore, after each service completion our customer will be the one departing with probability  $1 - \lambda/\mu$ . Thus, no matter how long the customer has already spent in the system, the probability he will depart in the next  $h$  time units is  $\mu h + o(h)$ , the probability that a service ends in that time, multiplied by  $1 - \lambda/\mu$ . That is, the customer will depart in the next  $h$  time units with probability  $(\mu - \lambda)h + o(h)$ ; which says that the hazard rate function of  $W^*$  is the constant  $\mu - \lambda$ . But only the exponential has a constant hazard rate, and so we can conclude that  $W^*$  is exponential with rate  $\mu - \lambda$ .

### 8.3.2. A Single-Server Exponential Queuing System Having Finite Capacity

In the previous model, we assumed that there was no limit on the number of customers that could be in the system at the same time. However, in reality there is always a finite system capacity  $N$ , in the sense that there can be no more than  $N$  customers in the system at any time. By this, we mean that if an arriving customer finds that there are already  $N$  customers present, then he does not enter the system.

As before, we let  $P_n$ ,  $0 \leq n \leq N$ , denote the limiting probability that there are  $n$  customers in the system. The rate-equality principle yields the following set of balance equations:

State	Rate at which the process leaves = rate at which it enters
0	$\lambda P_0 = \mu P_1$
$1 \leq n \leq N - 1$	$(\lambda + \mu) P_n = \lambda P_{n-1} + \mu P_{n+1}$
N	$\mu P_N = \lambda P_{N-1}$

The argument for state 0 is exactly as before. Namely, when in state 0, the process will leave only via an arrival (which occurs at rate  $\lambda$ ) and hence the rate at which the process leaves state 0 is  $\lambda P_0$ . On the other hand, the process can enter state 0 only from state 1 via a departure; hence, the rate at which the process enters state 0 is  $\mu P_1$ . The equation for states  $n$ , where  $1 \leq n < N$ , is the same as before. The equation for state  $N$  is different because now state  $N$  can only be left via a departure since an arriving customer will not enter the system when it is in state  $N$ ; also, state  $N$  can now only be entered from state  $N - 1$  (as there is no longer a state  $N + 1$ ) via an arrival.

To solve, we again rewrite the preceding system of equations:

$$\begin{aligned} P_1 &= \frac{\lambda}{\mu} P_0, \\ P_{n+1} &= \frac{\lambda}{\mu} P_n + \left( P_n - \frac{\lambda}{\mu} P_{n-1} \right), \quad 1 \leq n \leq N-1 \\ P_N &= \frac{\lambda}{\mu} P_{N-1} \end{aligned}$$

which, solving in terms of  $P_0$ , yields

$$\begin{aligned} P_1 &= \frac{\lambda}{\mu} P_0, \\ P_2 &= \frac{\lambda}{\mu} P_1 + \left( P_1 - \frac{\lambda}{\mu} P_0 \right) = \frac{\lambda}{\mu} P_1 = \left( \frac{\lambda}{\mu} \right)^2 P_0, \\ P_3 &= \frac{\lambda}{\mu} P_2 + \left( P_2 - \frac{\lambda}{\mu} P_1 \right) = \frac{\lambda}{\mu} P_2 = \left( \frac{\lambda}{\mu} \right)^3 P_0, \\ &\vdots \\ P_{N-1} &= \frac{\lambda}{\mu} P_{N-2} + \left( P_{N-2} - \frac{\lambda}{\mu} P_{N-3} \right) = \left( \frac{\lambda}{\mu} \right)^{N-1} P_0, \\ P_N &= \frac{\lambda}{\mu} P_{N-1} = \left( \frac{\lambda}{\mu} \right)^N P_0 \end{aligned} \quad (8.10)$$

By using the fact that  $\sum_{n=0}^N P_n = 1$ , we obtain

$$\begin{aligned} 1 &= P_0 \sum_{n=0}^N \left( \frac{\lambda}{\mu} \right)^n \\ &= P_0 \left[ \frac{1 - (\lambda/\mu)^{N+1}}{1 - \lambda/\mu} \right] \end{aligned}$$

or

$$P_0 = \frac{(1 - \lambda/\mu)}{1 - (\lambda/\mu)^{N+1}}$$

and hence from Equation (8.10) we obtain

$$P_n = \frac{(\lambda/\mu)^n (1 - \lambda/\mu)}{1 - (\lambda/\mu)^{N+1}}, \quad n = 0, 1, \dots, N \quad (8.11)$$

Note that in this case, there is no need to impose the condition that  $\lambda/\mu < 1$ . The queue size is, by definition, bounded so there is no possibility of its increasing indefinitely.

As before,  $L$  may be expressed in terms of  $P_n$  to yield

$$\begin{aligned} L &= \sum_{n=0}^N n P_n \\ &= \frac{(1 - \lambda/\mu)}{1 - (\lambda/\mu)^{N+1}} \sum_{n=0}^N n \left( \frac{\lambda}{\mu} \right)^n \end{aligned}$$

which after some algebra yields

$$L = \frac{\lambda[1 + N(\lambda/\mu)^{N+1} - (N+1)(\lambda/\mu)^N]}{(\mu - \lambda)(1 - (\lambda/\mu)^{N+1})} \quad (8.12)$$

In deriving  $W$ , the expected amount of time a customer spends in the system, we must be a little careful about what we mean by a customer. Specifically, are we including those "customers" who arrive to find the system full and thus do not spend any time in the system? Or, do we just want the expected time spent in the system by a customer that actually entered the system? The two questions lead, of course, to different answers. In the first case, we have  $\lambda_a = \lambda$ ; whereas in the second case, since the fraction of arrivals that actually enter the system is  $1 - P_N$ , it follows that  $\lambda_a = \lambda(1 - P_N)$ . Once it is clear what we mean by a customer,  $W$  can be obtained from

$$W = \frac{L}{\lambda_a}$$

**Example 8.3** Suppose that it costs  $c\mu$  dollars per hour to provide service at a rate  $\mu$ . Suppose also that we incur a gross profit of  $A$  dollars for each customer served. If the system has a capacity  $N$ , what service rate  $\mu$  maximizes our total profit?

**Solution:** To solve this, suppose that we use rate  $\mu$ . Let us determine the amount of money coming in per hour and subtract from this the amount going out each hour. This will give us our profit per hour, and we can choose  $\mu$  so as to maximize this.

Now, potential customers arrive at a rate  $\lambda$ . However, a certain proportion of them do not join the system; namely, those who arrive when there are  $N$  customers already in the system. Hence, since  $P_N$  is the proportion of time that the system is full, it follows that entering customers arrive at a rate of  $\lambda(1 - P_N)$ . Since each customer pays  $\$A$ , it follows that money comes in at an hourly rate of  $\lambda(1 - P_N)A$  and since it goes out at an hourly rate of  $c\mu$ , it follows that our total profit per



hour is given by

$$\begin{aligned}\text{Profit per hour} &= \lambda(1 - P_N)A - c\mu \\ &= \lambda A \left[ 1 - \frac{(\lambda/\mu)^N(1 - \lambda/\mu)}{1 - (\lambda/\mu)^{N+1}} \right] - c\mu \\ &= \frac{\lambda A [1 - (\lambda/\mu)^N]}{1 - (\lambda/\mu)^{N+1}} - c\mu\end{aligned}$$

For instance if  $N = 2$ ,  $\lambda = 1$ ,  $A = 10$ ,  $c = 1$ , then

$$\begin{aligned}\text{Profit per hour} &= \frac{10[1 - (1/\mu)^2]}{1 - (1/\mu)^3} - \mu \\ &= \frac{10(\mu^3 - \mu)}{\mu^3 - 1} - \mu\end{aligned}$$

in order to maximize profit we differentiate to obtain

$$\frac{d}{d\mu} [\text{Profit per hour}] = 10 \frac{(2\mu^3 - 3\mu^2 + 1)}{(\mu^3 - 1)^2} - 1$$

The value of  $\mu$  that maximizes our profit now can be obtained by equating to zero and solving numerically. ♦

In the previous two models, it has been quite easy to define the state of the system. Namely, it was defined as the number of people in the system. Now we shall consider some examples where a more detailed state space is necessary.

### 8.3.3. A Shoeshine Shop

Consider a shoeshine shop consisting of two chairs. Suppose that an entering customer first will go to chair 1. When his work is completed in chair 1, he will go either to chair 2 if that chair is empty or else wait in chair 1 until chair 2 becomes empty. Suppose that a potential customer will enter this shop as long as chair 1 is empty. (Thus, for instance, a potential customer might enter even if there is a customer in chair 2).

If we suppose that potential customers arrive in accordance with a Poisson process at rate  $\lambda$ , and that the service times for the two chairs are independent and have respective exponential rates of  $\mu_1$  and  $\mu_2$ , then

- what proportion of potential customers enters the system?
- what is the mean number of customers in the system?

- what is the average amount of time that an entering customer spends in the system?

To begin we must first decide upon an appropriate state space. It is clear that the state of the system must include more information than merely the number of customers in the system. For instance, it would not be enough to specify that there is one customer in the system as we would also have to know which chair he was in. Further, if we only know that there are two customers in the system, then we would not know if the man in chair 1 is still being served or if he is just waiting for the person in chair 2 to finish. To account for these points, the following state space, consisting of the five states,  $(0, 0)$ ,  $(1, 0)$ ,  $(0, 1)$ ,  $(1, 1)$ , and  $(b, 1)$ , will be used. The states have the following interpretation:

State	Interpretation
$(0, 0)$	There are no customers in the system.
$(1, 0)$	There is one customer in the system, and he is in chair 1.
$(0, 1)$	There is one customer in the system, and he is in chair 2.
$(1, 1)$	There are two customers in the system, and both are presently being served.
$(b, 1)$	There are two customers in the system, but the customer in the first chair has completed his work in that chair and is waiting for the second chair to become free.

It should be noted that when the system is in state  $(b, 1)$ , the person in chair 1, though not being served, is nevertheless "blocking" potential arrivals from entering the system.

As a prelude to writing down the balance equations, it is usually worthwhile to make a transition diagram. This is done by first drawing a circle for each state and then drawing an arrow labeled by the rate at which the process goes from one state to another. The transition diagram for this model is shown in Figure 8.1. The explanation for the diagram is as follows:

The arrow from state  $(0, 0)$  to state  $(1, 0)$  which is labeled  $\lambda$  means that when the process is in state  $(0, 0)$ , that is, when the system is empty, then it goes to state  $(1, 0)$  at a rate  $\lambda$ , that is via an arrival. The arrow from  $(0, 1)$  to  $(1, 1)$  is similarly explained.

When the process is in state  $(1, 0)$ , it will go to state  $(0, 1)$  when the customer in chair 1 is finished and this occurs at a rate  $\mu_1$ ; hence the arrow from  $(1, 0)$  to  $(0, 1)$  labeled  $\mu_1$ . The arrow from  $(1, 1)$  to  $(b, 1)$  is similarly explained.

When in state  $(b, 1)$  the process will go to state  $(0, 1)$  when the customer in chair 2 completes his service (which occurs at rate  $\mu_2$ ); hence the arrow from  $(b, 1)$  to  $(0, 1)$  labeled  $\mu_2$ . Also when in state  $(1, 1)$  the process will

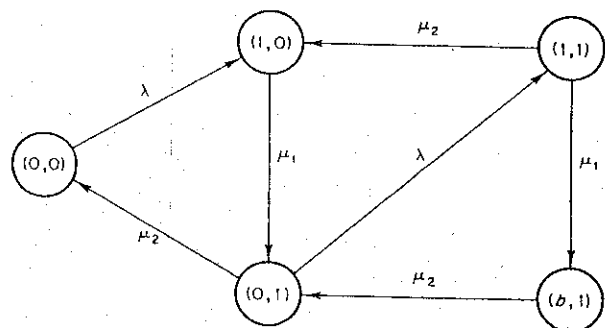


Figure 8.1. A transition diagram.

go to state (1, 0) when the man in chair 2 finishes and hence the arrow from (1, 1) to (1, 0) labeled  $\mu_2$ . Finally, if the process is in state (0, 1), then it will go to state (0, 0) when the man in chair 2 completes his service, hence the arrow from (0, 1) to (0, 0) labeled  $\mu_2$ .

Because there are no other possible transitions, this completes the transition diagram.

To write the balance equations we equate the sum of the arrows (multiplied by the probability of the states where they originate) coming into a state with the sum of the arrows (multiplied by the probability of the state) going out of that state. This gives

State	Rate that the process leaves = rate that it enters
(0, 0)	$\lambda P_{00} = \mu_2 P_{01}$
(1, 0)	$\mu_1 P_{10} = \lambda P_{00} + \mu_2 P_{11}$
(0, 1)	$(\lambda + \mu_2) P_{01} = \mu_1 P_{10} + \mu_2 P_{b1}$
(1, 1)	$(\mu_1 + \mu_2) P_{11} = \lambda P_{01}$
(b, 1)	$\mu_2 P_{b1} = \mu_1 P_{11}$

These along with the equation

$$P_{00} + P_{10} + P_{01} + P_{11} + P_{b1} = 1$$

may be solved to determine the limiting probabilities. Though it is easy to solve the preceding equations, the resulting solutions are quite involved and hence will not be explicitly presented. However, it is easy to answer our questions in terms of these limiting probabilities. First, since a potential customer will enter the system when the state is either (0, 0) or (0, 1), it follows that the proportion of customers entering the system is  $P_{00} + P_{01}$ .

Secondly, since there is one customer in the system whenever the state is (0, 1) or (1, 0) and two customers in the system whenever the state is (1, 1) or (b, 1), it follows that  $L$ , the average number in the system, is given by

$$L = P_{01} + P_{10} + 2(P_{11} + P_{b1})$$

To derive the average amount of time that an entering customer spends in the system, we use the relationship  $W = L/\lambda_a$ . Since a potential customer will enter the system when in state (0, 0) or (0, 1), it follows that  $\lambda_a = \lambda(P_{00} + P_{01})$  and hence

$$W = \frac{P_{01} + P_{10} + 2(P_{11} + P_{b1})}{\lambda(P_{00} + P_{01})}$$

**Example 8.4** (a) If  $\lambda = 1, \mu_1 = 1, \mu_2 = 2$ , then calculate the preceding quantities of interest.

(b) If  $\lambda = 1, \mu_1 = 2, \mu_2 = 1$ , then calculate the preceding.

**Solution:** (a) Solving the balance equations yields

$$P_{00} = \frac{12}{37}, \quad P_{10} = \frac{16}{37}, \quad P_{11} = \frac{2}{37}, \quad P_{01} = \frac{6}{37}, \quad P_{b1} = \frac{1}{37}$$

Hence,

$$L = \frac{28}{37}, \quad W = \frac{28}{18}$$

(b) Solving the balance equations yields

$$P_{00} = \frac{3}{11}, \quad P_{10} = \frac{2}{11}, \quad P_{11} = \frac{1}{11}, \quad P_{b1} = \frac{2}{11}, \quad P_{01} = \frac{3}{11}$$

Hence,

$$L = 1, \quad W = \frac{11}{6} \quad \blacklozenge$$

### 8.3.4. A Queueing System with Bulk Service

In this model, we consider a single-server exponential queueing system in which the server is able to serve two customers at the same time. Whenever the server completes a service, he then serves the next two customers at the same time. However, if there is only one customer in line, then he serves that customer by himself. We shall assume that his service time is exponential at rate  $\mu$  whether he is serving one or two customers. As usual, we suppose that customers arrive at an exponential rate  $\lambda$ . One example of such a system might be an elevator or a cable car which can take at most two passengers at any time.

It would seem that the state of the system would have to tell us not only how many customers there are in the system, but also whether one or two

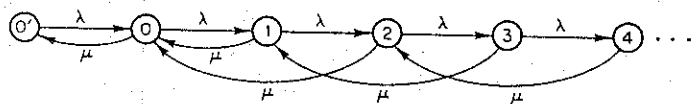


Figure 8.2.

are presently being served. However, it turns out that we can solve the problem easier not by concentrating on the number of customers in the system, but rather on the number in *queue*. So let us define the state as the number of customers waiting in queue, with two states when there is no one in queue. That is, let us have as a state space  $0', 0, 1, 2, \dots$ , with the interpretation

State	Interpretation
$0'$	No one in service
$0$	Server busy; no one waiting
$n, n > 0$	$n$ customers waiting

The transition diagram is shown in Figure 8.2 and the balance equations are

State	Rate at which the process leaves = rate at which it enters
$0'$	$\lambda P_{0'} = \mu P_0$
$0$	$(\lambda + \mu)P_0 = \lambda P_{0'} + \mu P_1 + \mu P_2$
$n, n \geq 1$	$(\lambda + \mu)P_n = \lambda P_{n-1} + \mu P_{n+2}$

Now the set of equations

$$(\lambda + \mu)P_n = \lambda P_{n-1} + \mu P_{n+2} \quad n = 1, 2, \dots \quad (8.13)$$

has a solution of the form

$$P_n = \alpha^n P_0$$

To see this, substitute the preceding in Equation (8.13) to obtain

$$(\lambda + \mu)\alpha^n P_0 = \lambda \alpha^{n-1} P_0 + \mu \alpha^{n+2} P_0$$

or

$$(\lambda + \mu)\alpha = \lambda + \mu\alpha^3$$

Solving this for  $\alpha$  yields the three roots:

$$\alpha = 1, \quad \alpha = \frac{-1 - \sqrt{1 + 4\lambda/\mu}}{2}, \quad \text{and} \quad \alpha = \frac{-1 + \sqrt{1 + 4\lambda/\mu}}{2}$$

As the first two are clearly not possible, it follows that

$$\alpha = \frac{\sqrt{1 + 4\lambda/\mu} - 1}{2}$$

Hence,

$$P_n = \alpha^n P_0,$$

$$P_{0'} = \frac{\mu}{\lambda} P_0$$

where the bottom equation follows from the first balance equation. (We can ignore the second balance equation as one of these equations is always redundant.) To obtain  $P_0$ , we use

$$P_0 + P_{0'} + \sum_{n=1}^{\infty} P_n = 1$$

or

$$P_0 \left[ 1 + \frac{\mu}{\lambda} + \sum_{n=1}^{\infty} \alpha^n \right] = 1$$

or

$$P_0 \left[ \frac{1}{1 - \alpha} + \frac{\mu}{\lambda} \right] = 1$$

or

$$P_0 = \frac{\lambda(1 - \alpha)}{\lambda + \mu(1 - \alpha)}$$

and thus

$$P_n = \frac{\alpha^n \lambda(1 - \alpha)}{\lambda + \mu(1 - \alpha)}, \quad n \geq 0$$

$$P_{0'} = \frac{\mu(1 - \alpha)}{\lambda + \mu(1 - \alpha)} \quad (8.14)$$

where

$$\alpha = \frac{\sqrt{1 + 4\lambda/\mu} - 1}{2}$$

Note that for the preceding to be valid we need  $\alpha < 1$ , or equivalently  $\lambda/\mu < 2$ , which is intuitive since the maximum service rate is  $2\mu$ , which must be larger than the arrival rate  $\lambda$  to avoid overloading the system.

All the relevant quantities of interest now can be determined. For instance, to determine the proportion of customers that are served alone,

we first note that the rate at which customers are served alone is  $\lambda P_0 + \mu P_1$ , since when the system is empty a customer will be served alone upon the next arrival and when there is one customer in queue he will be served alone upon a departure. As the rate at which customers are served is  $\lambda$ , it follows that

$$\begin{aligned} \text{proportion of customers that are served alone} &= \frac{\lambda P_0 + \mu P_1}{\lambda} \\ &= P_0 + \frac{\mu}{\lambda} P_1 \end{aligned}$$

Also,

$$\begin{aligned} L_Q &= \sum_{n=1}^{\infty} n P_n \\ &= \frac{\lambda(1-\alpha)}{\lambda + \mu(1-\alpha)} \sum_{n=1}^{\infty} n \alpha^n \quad \text{from Equation (8.14)} \\ &= \frac{\lambda \alpha}{(1-\alpha)[\lambda + \mu(1-\alpha)]} \quad \text{by algebraic identity } \sum_{n=1}^{\infty} n \alpha^n = \frac{\alpha}{(1-\alpha)^2} \end{aligned}$$

and

$$\begin{aligned} W_Q &= \frac{L_Q}{\lambda}, \\ W &= W_Q + \frac{1}{\mu}, \\ L &= \lambda W \end{aligned}$$

### 8.4. Network of Queues

#### 8.4.1. Open Systems

Consider a two-server system in which customers arrive at a Poisson rate  $\lambda$  at server 1. After being served by server 1 they then join the queue in front of server 2. We suppose there is infinite waiting space at both servers. Each server serves one customer at a time with server  $i$  taking an exponential time with rate  $\mu_i$  for a service,  $i = 1, 2$ . Such a system is called a *tandem* or *sequential* system (see Figure 8.3).

To analyze this system we need to keep track of the number of customers at server 1 and the number at server 2. So let us define the state by the pair  $(n, m)$ —meaning that there are  $n$  customers at server 1 and  $m$  at server 2.

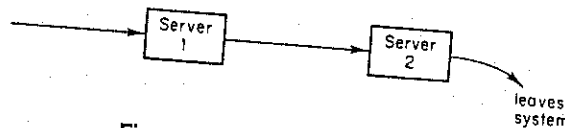


Figure 8.3. A tandem queue.

The balance equations are

State	Rate that the process leaves = rate that it enters
0, 0	$\lambda P_{0,0} = \mu_2 P_{0,1}$
$n, 0; n > 0$	$(\lambda + \mu_1) P_{n,0} = \mu_2 P_{n,1} + \lambda P_{n-1,0}$
$0, m; m > 0$	$(\lambda + \mu_2) P_{0,m} = \mu_2 P_{0,m+1} + \mu_1 P_{1,m-1}$
$n, m; nm > 0$	$(\lambda + \mu_1 + \mu_2) P_{n,m} = \mu_2 P_{n,m+1} + \mu_1 P_{n+1,m-1} + \lambda P_{n-1,m}$

(8.15)

Rather than directly attempting to solve these (along with the equation  $\sum_{n,m} P_{n,m} = 1$ ) we shall guess a solution and then verify that it indeed satisfies the preceding. We first note that the situation at server 1 is just as in an  $M/M/1$  model. Similarly, as it was shown in Section 6.6 that the departure process of an  $M/M/1$  queue is a Poisson process with rate  $\lambda$ , it follows that what server 2 faces is also an  $M/M/1$  queue. Hence, the probability that there are  $n$  customers at server 1 is

$$P\{n \text{ at server 1}\} = \left(\frac{\lambda}{\mu_1}\right)^n \left(1 - \frac{\lambda}{\mu_1}\right)$$

and, similarly,

$$P\{m \text{ at server 2}\} = \left(\frac{\lambda}{\mu_2}\right)^m \left(1 - \frac{\lambda}{\mu_2}\right)$$

Now if the numbers of customers at servers 1 and 2 were independent random variables, then it would follow that

$$P_{n,m} = \left(\frac{\lambda}{\mu_1}\right)^n \left(1 - \frac{\lambda}{\mu_1}\right) \left(\frac{\lambda}{\mu_2}\right)^m \left(1 - \frac{\lambda}{\mu_2}\right) \quad (8.16)$$

To verify that  $P_{n,m}$  is indeed equal to the preceding (and thus that the number of customers at server 1 is independent of the number at server 2), all we need do is verify that the preceding satisfies the set of Equations (8.15)—this suffices since we know that the  $P_{n,m}$  are the unique solution of Equations (8.15). Now, for instance, if we consider the first equation of (8.15), we need to show that

$$\lambda \left(1 - \frac{\lambda}{\mu_1}\right) \left(1 - \frac{\lambda}{\mu_2}\right) = \mu_2 \left(1 - \frac{\lambda}{\mu_1}\right) \left(\frac{\lambda}{\mu_2}\right) \left(1 - \frac{\lambda}{\mu_2}\right)$$

which is easily verified. We leave it as an exercise to show that the  $P_{n,m}$ , as given by Equation (8.16), satisfy all of the Equations (8.15), and are thus the limiting probabilities.

From the preceding we see that  $L$ , the average number of customers in the system, is given by

$$\begin{aligned} L &= \sum_{n,m} (n+m)P_{n,m} \\ &= \sum_n n \left(\frac{\lambda}{\mu_1}\right)^n \left(1 - \frac{\lambda}{\mu_1}\right) + \sum_m m \left(\frac{\lambda}{\mu_2}\right)^m \left(1 - \frac{\lambda}{\mu_2}\right) \\ &= \frac{\lambda}{\mu_1 - \lambda} + \frac{\lambda}{\mu_2 - \lambda} \end{aligned}$$

and from this we see that the average time a customer spends in the system is

$$W = \frac{L}{\lambda} = \frac{1}{\mu_1 - \lambda} + \frac{1}{\mu_2 - \lambda}$$

**Remarks** (i) The result (Equation 8.15) could have been obtained as a direct consequence of the time reversibility of an  $M/M/1$  (see Section 6.6). For not only does time reversibility imply that the output from server 1 is a Poisson process, but it also implies (Exercise 26 of Chapter 6) that the number of customers at server 1 is independent of the past departure times from server 1. As these past departure times constitute the arrival process to server 2, the independence of the numbers of customers in the two systems follows.

(ii) Since a Poisson arrival sees time averages, it follows that in a tandem queue the numbers of customers an arrival (to server 1) sees at the two servers are independent random variables. However, it should be noted that this does not imply that the waiting times of a given customer at the two servers are independent. For a counter example suppose that  $\lambda$  is very small with respect to  $\mu_1 = \mu_2$ ; and thus almost all customers have zero wait in queue at both servers. However, given that the wait in queue of a customer at server 1 is positive, his wait in queue at server 2 also will be positive with probability at least as large as  $\frac{1}{2}$  (why?). Hence, the waiting times in queue are not independent. Remarkably enough, however, it turns out that the total times (that is, service time plus wait in queue) that an arrival spends at the two servers are indeed independent random variables.

The preceding result can be substantially generalized. To do so, consider a system of  $k$  servers. Customers arrive from outside the system to server  $i$ ,  $i = 1, \dots, k$ , in accordance with independent Poisson process at rate  $r_i$ ; they

then join the queue at  $i$  until their turn at service comes. Once a customer is served by server  $i$ , he then joins the queue in front of server  $j$ ,  $j = 1, \dots, k$ , with probability  $P_{ij}$ . Hence,  $\sum_{j=1}^k P_{ij} \leq 1$ , and  $1 - \sum_{j=1}^k P_{ij}$  represents the probability that a customer departs the system after being served by server  $i$ .

If we let  $\lambda_j$  denote the total arrival rate of customers to server  $j$ , then the  $\lambda_j$  can be obtained as the solution of

$$\lambda_j = r_j + \sum_{i=1}^k \lambda_i P_{ij}, \quad i = 1, \dots, k \quad (8.17)$$

Equation (8.17) follows since  $r_j$  is the arrival rate of customers to  $j$  coming from outside the system and, as  $\lambda_i$  is the rate at which customers depart server  $i$  (rate in must equal rate out),  $\lambda_i P_{ij}$  is the arrival rate to  $j$  of those coming from server  $i$ .

It turns out that the number of customers at each of the servers is independent and of the form

$$P\{n \text{ customers at server } j\} = \left(\frac{\lambda_j}{\mu_j}\right)^n \left(1 - \frac{\lambda_j}{\mu_j}\right), \quad n \geq 1$$

where  $\mu_j$  is the exponential service rate at server  $j$  and the  $\lambda_j$  are the solution to Equation (8.17). Of course, it is necessary that  $\lambda_j/\mu_j < 1$  for all  $j$ . To prove this, we first note that it is equivalent to asserting that the limiting probabilities  $P(n_1, n_2, \dots, n_k) = P\{n_j \text{ at server } j, j = 1, \dots, k\}$  are given by

$$P(n_1, n_2, \dots, n_k) = \prod_{j=1}^k \left(\frac{\lambda_j}{\mu_j}\right)^{n_j} \left(1 - \frac{\lambda_j}{\mu_j}\right) \quad (8.18)$$

which can be verified by showing that it satisfies the balance equations for this model.

The average number of customers in the system is

$$\begin{aligned} L &= \sum_{j=1}^k \text{average number at server } j \\ &= \sum_{j=1}^k \frac{\lambda_j}{\mu_j - \lambda_j} \end{aligned}$$

The average time a customer spends in the system can be obtained from  $L = \lambda W$  with  $\lambda = \sum_{j=1}^k r_j$ . (Why not  $\lambda = \sum_{j=1}^k \lambda_j$ ?) This yields

$$W = \frac{\sum_{j=1}^k \lambda_j / (\mu_j - \lambda_j)}{\sum_{j=1}^k r_j}$$

**Remarks** The result embodied in Equation (8.18) is rather remarkable in that it says that the distribution of the number of customers at server  $i$  is the same as in an  $M/M/1$  system with rates  $\lambda_i$  and  $\mu_i$ . What is remarkable is that in the network model the arrival process at node  $i$  need not be a Poisson process. For if there is a possibility that a customer may visit a server more than once (a situation called *feedback*), then the arrival process will not be Poisson. An easy example illustrating this is to suppose that there is a single server whose service rate is very large with respect to the arrival rate from outside. Suppose also that with probability  $p = 0.9$  a customer upon completion of service is fed back into the system. Hence, at an arrival time epoch there is a large probability of another arrival in a short time (namely, the feedback arrival); whereas at an arbitrary time point there will be only a very slight chance of an arrival occurring shortly (since  $\lambda$  is so very small). Hence, the arrival process does not possess independent increments and so cannot be Poisson. In fact even though it is straightforward to verify Equation (8.18) there does not appear to be, at present, any simple explanation as to why it is, in fact, true.

Thus, we see that when feedback is allowed the steady-state probabilities of the number of customers at any given station have the same distribution as in an  $M/M/1$  model even though the model is not  $M/M/1$ . (Presumably such quantities as the joint distribution of the number at the station at two different time points will not be the same as for an  $M/M/1$ .)

**Example 8.5** Consider a system of two servers where customers from outside the system arrive at server 1 at a Poisson rate 4 and at server 2 at a Poisson rate 5. The service rates of 1 and 2 are respectively 8 and 10. A customer upon completion of service at server 1 is equally likely to go to server 2 or to leave the system (i.e.,  $P_{11} = 0$ ,  $P_{12} = \frac{1}{2}$ ); whereas a departure from server 2 will go 25 percent of the time to server 1 and will depart the system otherwise (i.e.,  $P_{21} = \frac{1}{4}$ ,  $P_{22} = 0$ ). Determine the limiting probabilities,  $L$ , and  $W$ .

**Solution:** The total arrival rates to servers 1 and 2—call them  $\lambda_1$  and  $\lambda_2$ —can be obtained from Equation (8.17). That is, we have

$$\lambda_1 = 4 + \frac{1}{4}\lambda_2,$$

$$\lambda_2 = 5 + \frac{1}{2}\lambda_1$$

implying that

$$\lambda_1 = 6, \quad \lambda_2 = 8$$

Hence,

$$P\{n \text{ at server 1, } m \text{ at server 2}\} = \left(\frac{3}{4}\right)^n \frac{1}{4} \left(\frac{4}{5}\right)^m \frac{1}{5} \\ = \frac{1}{20} \left(\frac{3}{4}\right)^n \left(\frac{4}{5}\right)^m$$

and

$$L = \frac{6}{8-6} + \frac{8}{10-8} = 7,$$

$$W = \frac{L}{9} = \frac{7}{9} \quad \blacklozenge$$

### 8.4.2. Closed Systems

The queueing systems described in Section 8.4.1 are called *open systems* since customers are able to enter and depart the system. A system in which new customers never enter and existing ones never depart is called a *closed system*.

Let us suppose that we have  $m$  customers moving among a system of  $k$  servers. When a customer completes service at server  $i$ , she then joins the queue in front of server  $j$ ,  $j = 1, \dots, k$ , with probability  $P_{ij}$ , where we now suppose that  $\sum_{j=1}^k P_{ij} = 1$  for all  $i = 1, \dots, k$ . That is,  $\mathbf{P} = [P_{ij}]$  is Markov transition probability matrix, which we shall assume is irreducible. Let  $\pi = (\pi_1, \dots, \pi_k)$  denote the stationary probabilities for this Markov chain; that is,  $\pi$  is the unique positive solution of

$$\pi_j = \sum_{i=1}^k \pi_i P_{ij}, \\ \sum_{j=1}^k \pi_j = 1 \tag{8.19}$$

If we denote the average arrival rate (or equivalently the average service completion rate) at server  $j$  by  $\lambda_m(j)$ ,  $j = 1, \dots, k$  then, analogous to Equation (8.17), the  $\lambda_m(j)$  satisfy

$$\lambda_m(j) = \sum_{i=1}^k \lambda_m(i) P_{ij}$$

Hence, from (8.19) we can conclude that

$$\lambda_m(j) = \lambda_m \pi_j, \quad j = 1, 2, \dots, k \tag{8.20}$$

where

$$\lambda_m = \sum_{j=1}^k \lambda_m(j) \tag{8.21}$$

From Equation (8.21), we see that  $\lambda_m$  is the average service completion rate of the entire system, that is, it is the system *throughput* rate.\*

If we let  $P_m(n_1, n_2, \dots, n_k)$  denote the limiting probabilities

$$P_m(n_1, n_2, \dots, n_k) = P\{n_j \text{ customers at server } j, j = 1, \dots, k\}$$

then, by verifying that they satisfy the balance equation, it can be shown that

$$P_m(n_1, n_2, \dots, n_k) = \begin{cases} K_m \prod_{j=1}^k (\lambda_m(j)/\mu_j)^{n_j}, & \text{if } \sum_{j=1}^k n_j = m \\ 0, & \text{otherwise} \end{cases}$$

But from Equation (8.20) we thus obtain that

$$P_m(n_1, n_2, \dots, n_k) = \begin{cases} C_m \prod_{j=1}^k (\pi_j/\mu_j)^{n_j}, & \text{if } \sum_{j=1}^k n_j = m \\ 0, & \text{otherwise} \end{cases} \quad (8.22)$$

where

$$C_m = \left[ \sum_{\substack{n_1, \dots, n_k \\ \sum n_j = m}} \prod_{j=1}^k (\pi_j/\mu_j)^{n_j} \right]^{-1} \quad (8.23)$$

Equation (8.22) is not as useful as one might suppose, for in order to utilize it we must determine the normalizing constant  $C_m$  given by Equation (8.23) which requires summing the products  $\prod_{j=1}^k (\pi_j/\mu_j)^{n_j}$  over all the feasible

vectors  $(n_1, \dots, n_k)$ :  $\sum_{j=1}^k n_j = m$ . Hence, since there are  $\binom{m+k-1}{m}$  vectors this is only computationally feasible for relatively small values of  $m$  and  $k$ .

We will now present an approach that will enable us to determine recursively many of the quantities of interest in this model without first computing the normalizing constants. To begin, consider a customer who has just left server  $i$  and is headed to server  $j$ , and let us determine the probability of the system as seen by this customer. In particular, let us determine the probability that this customer observes, at that moment,  $n_l$  customers at server  $l$ ,  $l = 1, \dots, k$ ,  $\sum_{l=1}^k n_l = m - 1$ . This is done

\* We are using the notation of  $\lambda_m(j)$  and  $\lambda_m$  to indicate the dependence on the number of customers in the closed system. This will be used in recursive relations we will develop.

as follows:

$$\begin{aligned} &P\{\text{customer observes } n_l \text{ at server } l, \\ & \quad l = 1, \dots, k \mid \text{customer goes from } i \text{ to } j\} \\ &= \frac{P\{\text{state is } (n_1, \dots, n_i + 1, \dots, n_j, \dots, n_k), \text{ customer goes from } i \text{ to } j\}}{P\{\text{customer goes from } i \text{ to } j\}} \\ &= \frac{P_m(n_1, \dots, n_i + 1, \dots, n_j, \dots, n_k) \mu_i P_{ij}}{\sum_{n: \sum n_j = m-1} P_m(n_1, \dots, n_i + 1, \dots, n_k) \mu_i P_{ij}} \\ &= \frac{(\pi_i/\mu_i) \prod_{j=1}^k (\pi_j/\mu_j)^{n_j}}{K} \quad \text{from (8.22)} \\ &= C \prod_{j=1}^k (\pi_j/\mu_j)^{n_j} \end{aligned}$$

where  $C$  does not depend on  $n_1, \dots, n_k$ . But because the above is a probability density on the set of vectors  $(n_1, \dots, n_k)$ ,  $\sum_{j=1}^k n_j = m - 1$ , it follows from (8.22) that it must equal  $P_{m-1}(n_1, \dots, n_k)$ . Hence,

$$\begin{aligned} &P\{\text{customer observes } n_l \text{ at server } l, \\ & \quad l = 1, \dots, k \mid \text{customer goes from } i \text{ to } j\} \\ &= P_{m-1}(n_1, \dots, n_k), \quad \sum_{l=1}^k n_l = m - 1 \end{aligned} \quad (8.24)$$

As (8.24) is true for all  $i$ , we thus have proven the following proposition, known as the arrival theorem.

**Proposition 8.3** (The Arrival Theorem). In the closed network system with  $m$  customers, the system as seen by arrivals to server  $j$  is distributed as the stationary distribution in the same network system when there are only  $m - 1$  customers.

Denote by  $L_m(j)$  and  $W_m(j)$  the average number of customers and the average time a customer spends at server  $j$  when there are  $m$  customers in the network. Upon conditioning on the number of customers found at server  $j$  by an arrival to that server, it follows that

$$\begin{aligned} W_m(j) &= \frac{1 + E_m[\text{number at server } j \text{ as seen by an arrival}]}{\mu_j} \\ &= \frac{1 + L_{m-1}(j)}{\mu_j} \end{aligned} \quad (8.25)$$

where the last equality follows from the arrival theorem. Now when there are  $m - 1$  customers in the system, then, from Equation (8.20),  $\lambda_{m-1}(j)$ , the average arrival rate to server  $j$ , satisfies

$$\lambda_{m-1}(j) = \lambda_{m-1} \pi_j$$

Now, applying the basic cost identity Equation (8.1) with the cost rule being that each customer in the network system of  $m - 1$  customers pays one unit time while at server  $j$ , we obtain

$$L_{m-1}(j) = \lambda_{m-1} \pi_j W_{m-1}(j) \quad (8.26)$$

Using Equation (8.25), this yields

$$W_m(j) = \frac{1 + \lambda_{m-1} \pi_j W_{m-1}(j)}{\mu_j} \quad (8.27)$$

Also using the fact that  $\sum_{j=1}^k L_{m-1}(j) = m - 1$  (why?) we obtain, from Equation (8.26):

$$m - 1 = \lambda_{m-1} \sum_{j=1}^k \pi_j W_{m-1}(j)$$

or

$$\lambda_{m-1} = \frac{m - 1}{\sum_{i=1}^k \pi_i W_{m-1}(i)} \quad (8.28)$$

Hence, from Equation (8.27), we obtain the recursion

$$W_m(j) = \frac{1}{\mu_j} + \frac{(m - 1) \pi_j W_{m-1}(j)}{\mu_j \sum_{i=1}^k \pi_i W_{m-1}(i)} \quad (8.29)$$

Starting with the stationary probabilities  $\pi_j, j = 1, \dots, k$ , and  $W_1(j) = 1/\mu_j$  we can now use Equation (8.29) to determine recursively  $W_2(j), W_3(j), \dots, W_m(j)$ . We can then determine the throughput rate  $\lambda_m$  by using Equation (8.28), and this will determine  $L_m(j)$  by Equation (8.26). This recursive approach is called *mean value analysis*.

**Example 8.6** Consider a  $k$ -server network in which the customers move in a cyclic permutation. That is,

$$P_{i,i+1} = 1, \quad i = 1, 2, \dots, k - 1, \quad P_{k,1} = 1$$

Let us determine the average number of customers at server  $j$  when there are two customers in the system. Now, for this network

$$\pi_i = 1/k, \quad i = 1, \dots, k$$

and as

$$W_1(j) = \frac{1}{\mu_j}$$

we obtain from Equation (8.29) that

$$\begin{aligned} W_2(j) &= \frac{1}{\mu_j} + \frac{(1/k)(1/\mu_j)}{\mu_j \sum_{i=1}^k (1/k)(1/\mu_i)} \\ &= \frac{1}{\mu_j} + \frac{1}{\mu_j^2 \sum_{i=1}^k 1/\mu_i} \end{aligned}$$

Hence, from Equation (8.28),

$$\lambda_2 = \frac{2}{\sum_{i=1}^k \frac{1}{k} W_2(i)} = \frac{2k}{\sum_{i=1}^k \left( \frac{1}{\mu_i} + \frac{1}{\mu_i^2 \sum_{i=1}^k 1/\mu_i} \right)}$$

and finally, using Equation (8.26),

$$\begin{aligned} L_2(j) &= \lambda_2 \frac{1}{k} W_2(j) \\ &= \frac{2 \left( \frac{1}{\mu_j} + \frac{1}{\mu_j^2 \sum_{i=1}^k 1/\mu_i} \right)}{\sum_{i=1}^k \left( \frac{1}{\mu_i} + \frac{1}{\mu_i^2 \sum_{i=1}^k 1/\mu_i} \right)} \quad \blacklozenge \end{aligned}$$

Another approach to learning about the stationary probabilities specified by Equation (8.22), which finesses the computational difficulties of computing the constant  $C_m$ , is to use the Gibbs sampler of Section 4.9 to generate a Markov chain having these stationary probabilities. To begin, note that since there are always a total of  $m$  customers in the system, Equation (8.22) may equivalently be written as a joint mass function of the numbers of customers at each of the servers  $1, \dots, k - 1$ , as follows:

$$\begin{aligned} P_m(n_1, \dots, n_{k-1}) &= C_m (\pi_k / \mu_k)^{m - \sum_{j=1}^{k-1} n_j} \prod_{j=1}^{k-1} (\pi_j / \mu_j)^{n_j} \\ &= K \prod_{j=1}^{k-1} (a_j)^{n_j}, \quad \sum_{j=1}^{k-1} n_j \leq m \end{aligned}$$

where  $a_j = (\pi_j \mu_k) / (\pi_k \mu_j)$ ,  $j = 1, \dots, k - 1$ . Now, if  $\mathbf{N} = (N_1, \dots, N_{k-1})$  has the preceding joint mass function then,

$$\begin{aligned} P\{N_i = n \mid N_1 = n_1, \dots, N_{i-1} = n_{i-1}, N_{i+1} = n_{i+1}, \dots, N_{k-1} = n_{k-1}\} \\ &= \frac{P_m(n_1, \dots, n_{i-1}, n, n_{i+1}, \dots, n_{k-1})}{\sum_r P_m(n_1, \dots, n_{i-1}, r, n_{i+1}, \dots, n_{k-1})} \\ &= C a_i^n, \quad n \leq m - \sum_{j \neq i} n_j \end{aligned}$$



It follows from the preceding that we may use the Gibbs sampler to generate the values of a Markov chain whose limiting probability mass function is  $P_m(n_1, \dots, n_{k-1})$  as follows:

1. Let  $(n_1, \dots, n_{k-1})$  be arbitrary nonnegative integers satisfying  $\sum_{j=1}^{k-1} n_j \leq m$ .
2. Generate a random variable  $I$  that is equally likely to be any of  $1, \dots, k-1$ .
3. If  $I = i$ , set  $s = m - \sum_{j \neq i} n_j$ , and generate the value of a random variable  $X$  having probability mass function

$$P\{X = n\} = C a_i^n, \quad n = 0, \dots, s$$

4. Let  $n_i = X$  and go to step 2.

The successive values of the state vector  $(n_1, \dots, n_{k-1}, m - \sum_{j=1}^{k-1} n_j)$  constitute the sequence of states of a Markov chain with the limiting distribution  $P_m$ . All quantities of interest can be estimated from this sequence. For instance, the average of the values of the  $j$ th coordinate of these vectors will converge to the mean number of individuals at station  $j$ , the proportion of vectors whose  $j$ th coordinate is less than  $r$  will converge to the limiting probability that the number of individuals at station  $j$  is less than  $r$ , and so on.

### 8.5. The System M/G/1

#### 8.5.1. Preliminaries: Work and Another Cost Identity

For an arbitrary queueing system, let us define the work in the system at any time  $t$  to be the sum of the remaining service times of all customers in the system at time  $t$ . For instance, suppose there are three customers in the system—the one in service having been there for three of his required five units of service time, and both people in queue having service times of six units. Then the work at that time is  $2 + 6 + 6 = 14$ . Let  $V$  denote the (time) average work in the system.

Now recall the fundamental cost Equation (8.1), which states that the

$$\begin{aligned} &\text{average rate at which the system earns} \\ &= \lambda_a \times \text{average amount a customer pays} \end{aligned}$$

and consider the following cost rule: *Each customer pays at a rate of  $y$ /unit time when his remaining service time is  $y$ , whether he is in queue or in service.* Thus, the rate at which the system earns is just the work in the system; so the basic identity yields that

$$V = \lambda_a E[\text{amount paid by a customer}]$$

Now, let  $S$  and  $W_Q^*$  denote respectively the service time and the time a given customer spends waiting in queue. Then, since the customer pays at a constant rate of  $S$  per unit time while he waits in queue and at a rate of  $S - x$  after spending an amount of time  $x$  in service, we have

$$E[\text{amount paid by a customer}] = E\left[SW_Q^* + \int_0^S (S - x) dx\right]$$

and thus

$$V = \lambda_a E[SW_Q^*] + \frac{\lambda_a E[S^2]}{2} \tag{8.30}$$

It should be noted that the preceding is a basic queueing identity [like Equations (8.2)–(8.4)] and as such valid in almost all models. In addition, if a customer's service time is independent of his wait in queue (as is usually, but not always the case),<sup>†</sup> then we have from Equation (8.30) that

$$V = \lambda_a E\{S\}W_Q + \frac{\lambda_a E[S^2]}{2} \tag{8.31}$$

#### 8.5.2. Application of Work to M/G/1

The M/G/1 model assumes (i) Poisson arrivals at rate  $\lambda$ ; (ii) a general service distribution; and (iii) a single server. In addition, we will suppose that customers are served in the order of their arrival.

Now, for an arbitrary customer in an M/G/1 system,

$$\text{Customer's wait in queue} = \text{work in the system when he arrives} \tag{8.32}$$

this follows since there is only a single server (think about it!). Taking expectations of both sides of Equation (8.32) yields

$$W_Q = \text{average work as seen by an arrival}$$

But, due to Poisson arrivals, the average work as seen by an arrival will equal  $V$ , the time average work in the system. Hence, for the model M/G/1,

$$W_Q = V$$

<sup>†</sup> For an example where it is not true, see Section 8.6.2.

The preceding in conjunction with the identity

$$V = \lambda E[S]W_Q + \frac{\lambda E[S^2]}{2}$$

yields the so-called Pollaczek-Khinchine formula,

$$W_Q = \frac{\lambda E[S^2]}{2(1 - \lambda E[S])} \quad (8.33)$$

where  $E[S]$  and  $E[S^2]$  are the first two moments of the service distribution.

The quantities  $L$ ,  $L_Q$ , and  $W$  can be obtained from Equation (8.33) as

$$\begin{aligned} L_Q &= \lambda W_Q = \frac{\lambda^2 E[S^2]}{2(1 - \lambda E[S])}, \\ W &= W_Q + E[S] = \frac{\lambda E[S^2]}{2(1 - \lambda E[S])} + E[S], \\ L &= \lambda W = \frac{\lambda^2 E[S^2]}{2(1 - \lambda E[S])} + \lambda E[S] \end{aligned} \quad (8.34)$$

**Remarks** (i) For the preceding quantities to be finite, we need  $\lambda E[S] < 1$ . This condition is intuitive since we know from renewal theory that if the server was always busy, then the departure rate would be  $1/E[S]$  (see Section 7.3), which must be larger than the arrival rate  $\lambda$  to keep things finite.

(ii) Since  $E[S^2] = \text{Var}(S) + (E[S])^2$ , we see from Equations (8.33) and (8.34) that, for fixed mean service time,  $L$ ,  $L_Q$ ,  $W$ , and  $W_Q$  all increase as the variance of the service distribution increases.

(iii) Another approach to obtain  $W_Q$  is presented in Exercise 34.

### 8.5.3. Busy Periods

The system alternates between idle periods (when there are no customers in the system, and so the server is idle) and busy periods (when there is at least one customer in the system, and so the server is busy).

Let us denote by  $I_n$  and  $B_n$ , respectively, the lengths of the  $n$ th idle and the  $n$ th busy period,  $n \geq 1$ . Hence, in the first  $\sum_{j=1}^n (I_j + B_j)$  time units the server will be idle for a time  $\sum_{j=1}^n I_j$ , and so the proportion of time that the server will be idle, which of course is just  $P_0$ , can be expressed as

$$\begin{aligned} P_0 &= \text{proportion of idle time} \\ &= \lim_{n \rightarrow \infty} \frac{I_1 + \cdots + I_n}{I_1 + \cdots + I_n + B_1 + \cdots + B_n} \end{aligned}$$

Now it is easy to see that the  $I_1, I_2, \dots$  are independent and identically distributed as are the  $B_1, B_2, \dots$ . Hence, by dividing the numerator and the denominator of the right side of the above by  $n$ , and then applying the strong law of large numbers, we obtain

$$\begin{aligned} P_0 &= \lim_{n \rightarrow \infty} \frac{(I_1 + \cdots + I_n)/n}{(I_1 + \cdots + I_n)/n + (B_1 + \cdots + B_n)/n} \\ &= \frac{E[I]}{E[I] + E[B]} \end{aligned} \quad (8.35)$$

where  $I$  and  $B$  represent idle and busy time random variables.

Now  $I$  represents the time from when a customer departs and leaves the system empty until the next arrival. Hence, from Poisson arrivals, it follows that  $I$  is exponential with rate  $\lambda$ , and so

$$E[I] = \frac{1}{\lambda} \quad (8.36)$$

To compute  $P_0$ , we note from Equation (8.4) (obtained from the fundamental cost equation by supposing that a customer pays at a rate of one per unit time while in service) that

$$\text{average number of busy servers} = \lambda E[S]$$

However, as the left-hand side of the above equals  $1 - P_0$  (why?), we have

$$P_0 = 1 - \lambda E[S] \quad (8.37)$$

and, from Equations (8.35)–(8.37),

$$1 - \lambda E[S] = \frac{1/\lambda}{1/\lambda + E[B]}$$

or

$$E[B] = \frac{E[S]}{1 - \lambda E[S]}$$

Another quantity of interest is  $C$ , the number of customers served in a busy period. The mean of  $C$  can be computed by noting that, on the average, for every  $E[C]$  arrivals exactly one will find the system empty (namely, the first customer in the busy period). Hence,

$$a_0 = \frac{1}{E[C]}$$

and, as  $a_0 = P_0 = 1 - \lambda E[S]$  because of Poisson arrivals, we see that

$$E[C] = \frac{1}{1 - \lambda E[S]}$$

## 8.6. Variations on the M/G/1

### 8.6.1. The M/G/1 with Random-Sized Batch Arrivals

Suppose that, as in the M/G/1, arrivals occur in accordance with a Poisson process having rate  $\lambda$ . But now suppose that each arrival consists not of a single customer but of a random number of customers. As before there is a single server whose service times have distribution  $G$ .

Let us denote by  $\alpha_j, j \geq 1$ , the probability that an arbitrary batch consists of  $j$  customers; and let  $N$  denote a random variable representing the size of a batch and so  $P\{N = j\} = \alpha_j$ . Since  $\lambda_a = \lambda E(N)$ , the basic formula for work [Equation (8.31)] becomes

$$V = \lambda E[N] \left[ E(S)W_Q + \frac{E(S^2)}{2} \right] \quad (8.38)$$

To obtain a second equation relating  $V$  to  $W_Q$ , consider an average customer. We have that

$$\begin{aligned} \text{his wait in queue} &= \text{work in system when he arrives} \\ &+ \text{his waiting time due to those in his batch} \end{aligned}$$

Taking expectations and using the fact that Poisson arrivals see time averages yields

$$\begin{aligned} W_Q &= V + E[\text{waiting time due to those in his batch}] \\ &= V + E[W_B] \end{aligned} \quad (8.39)$$

Now,  $E(W_B)$  can be computed by conditioning on the number in the batch, but we must be careful. For the probability that our average customer comes from a batch of size  $j$  is *not*  $\alpha_j$ . For  $\alpha_j$  is the proportion of batches which are of size  $j$ , and if we pick a customer at random, it is more likely that he comes from a larger rather than a smaller batch. (For instance, suppose  $\alpha_1 = \alpha_{100} = \frac{1}{2}$ , then half the batches are of size 1 but 100/101 of the customers will come from a batch of size 100!)

To determine the probability that our average customer came from a batch of size  $j$  we reason as follows: Let  $M$  be a large number. Then of the first  $M$  batches approximately  $M\alpha_j$  will be of size  $j, j \geq 1$ , and thus there would have been approximately  $jM\alpha_j$  customers that arrived in a batch of size  $j$ . Hence, the proportion of arrivals in the first  $M$  batches that were from batches of size  $j$  is approximately  $jM\alpha_j / \sum_j jM\alpha_j$ . This proportion

becomes exact as  $M \rightarrow \infty$ , and so we see that

$$\begin{aligned} \text{proportion of customers from batches of size } j &= \frac{j\alpha_j}{\sum_j j\alpha_j} \\ &= \frac{j\alpha_j}{E[N]} \end{aligned}$$

We are now ready to compute  $E(W_B)$ , the expected wait in queue due to others in the batch:

$$E[W_B] = \sum_j E[W_B | \text{batch of size } j] \frac{j\alpha_j}{E[N]} \quad (8.40)$$

Now if there are  $j$  customers in his batch, then our customer would have to wait for  $i - 1$  of them to be served if he was  $i$ th in line among his batch members. As he is equally likely to be either 1st, 2nd, ..., or  $j$ th in line we see that

$$\begin{aligned} E[W_B | \text{batch is of size } j] &= \sum_{i=1}^j (i-1)E(S) \frac{1}{j} \\ &= \frac{j-1}{2} E[S] \end{aligned}$$

Substituting this in Equation (8.40) yields

$$\begin{aligned} E[W_B] &= \frac{E[S]}{2E[N]} \sum_j (j-1)j\alpha_j \\ &= \frac{E[S](E[N^2] - E[N])}{2E[N]} \end{aligned}$$

and from Equations (8.38) and (8.39) we obtain

$$W_Q = \frac{E[S](E[N^2] - E[N])/2E[N] + \lambda E[N]E[S]^2/2}{1 - \lambda E[N]E[S]}$$

**Remarks** (i) Note that the condition for  $W_Q$  to be finite is that

$$\lambda E(N) < \frac{1}{E[S]}$$

which again says that the arrival rate must be less than the service rate (when the server is busy).

(ii) For fixed value of  $E[N]$ ,  $W_Q$  is increasing in  $\text{Var}[N]$ , again indicating that "single-server queues do not like variation."

(iii) The other quantities  $L$ ,  $L_Q$ , and  $W$  can be obtained by using

$$W = W_Q + E[S],$$

$$L = \lambda_a W = \lambda E[N]W,$$

$$L_Q = \lambda E[N]W_Q$$

### 8.6.2. Priority Queues

Priority queuing systems are ones in which customers are classified into types and then given service priority according to their type. Consider the situation where there are two types of customers, which arrive according to independent Poisson processes with respective rates  $\lambda_1$  and  $\lambda_2$ , and have service distributions  $G_1$  and  $G_2$ . We suppose that type 1 customers are given service priority, in that service will never begin on a type 2 customer if a type 1 is waiting. However, if a type 2 is being served and a type 1 arrives, we assume that the service of the type 2 is continued until completion. That is, there is no preemption once service has begun.

Let  $W_Q^i$  denote the average wait in queue of a type  $i$  customer,  $i = 1, 2$ . Our objective is to compute the  $W_Q^i$ .

First, note that the total work in the system at any time would be exactly the same no matter what priority rule was employed (as long as the server is always busy whenever there are customers in the system). This is so since the work will always decrease at a rate of one per unit time when the server is busy (no matter who is in service) and will always jump by the service time of an arrival. Hence, the work in the system is exactly as it would be if there was no priority rule but rather a first-come, first-served (called FIFO) ordering. However, under FIFO the above model is just M/G/1 with

$$\lambda = \lambda_1 + \lambda_2$$

$$G(x) = \frac{\lambda_1}{\lambda} G_1(x) + \frac{\lambda_2}{\lambda} G_2(x) \quad (8.41)$$

which follows since the combination of two independent Poisson processes is itself a Poisson process whose rate is the sum of the rates of the component processes. The service distribution  $G$  can be obtained by conditioning on which priority class the arrival is from—as is done in Equation (8.41).

Hence, from the results of Section 8.5, it follows that  $V$ , the average work in the priority queueing system, is given by

$$\begin{aligned} V &= \frac{\lambda E[S^2]}{2(1 - \lambda E[S])} \\ &= \frac{\lambda((\lambda_1/\lambda)E[S_1^2] + (\lambda_2/\lambda)E[S_2^2])}{2[1 - \lambda((\lambda_1/\lambda)E[S_1] + (\lambda_2/\lambda)E[S_2])]} \\ &= \frac{\lambda_1 E[S_1^2] + \lambda_2 E[S_2^2]}{2(1 - \lambda_1 E[S_1] - \lambda_2 E[S_2])} \end{aligned} \quad (8.42)$$

where  $S_i$  has distribution  $G_i$ ,  $i = 1, 2$ .

Continuing in our quest for  $W_Q^i$ , let us note that  $S$  and  $W_Q^*$ , the service and wait in queue of an arbitrary customer, are not independent in the priority model since knowledge about  $S$  gives us information as to the type of customer which in turn gives us information about  $W_Q^*$ . To get around this we will compute separately the average amount of type 1 and type 2 work in the system. Denoting  $V^i$  as the average amount of type  $i$  work we have, exactly as in Section 8.5.1,

$$V^i = \lambda_i E[S_i] W_Q^i + \frac{\lambda_i E[S_i^2]}{2}, \quad i = 1, 2 \quad (8.43)$$

If we define

$$V_Q^i = \lambda_i E[S_i] W_Q^i,$$

$$V_S^i = \frac{\lambda_i E[S_i^2]}{2}$$

then we may interpret  $V_Q^i$  as the average amount of type  $i$  work in queue, and  $V_S^i$  as the average amount of type  $i$  work in service (why?).

Now we are ready to compute  $W_Q^1$ . To do so, consider an arbitrary type 1 arrival. Then

his delay = amount of type 1 work in the system when he arrives  
+ amounts of type 2 work in service when he arrives

Taking expectations and using the fact that Poisson arrivals see time averages yields

$$\begin{aligned} W_Q^1 &= V^1 + V_S^2 \\ &= \lambda_1 E[S_1] W_Q^1 + \frac{\lambda_1 E[S_1^2]}{2} + \frac{\lambda_2 E[S_2^2]}{2} \end{aligned} \quad (8.44)$$

or

$$W_Q^1 = \frac{\lambda_1 E[S_1^2] + \lambda_2 E[S_2^2]}{2(1 - \lambda_1 E[S_1])} \quad (8.45)$$

To obtain  $W_Q^2$  we first note that since  $V = V^1 + V^2$ , we have from Equations (8.42) and (8.43) that

$$\begin{aligned} \frac{\lambda_1 E[S_1^2] + \lambda_2 E[S_2^2]}{2(1 - \lambda_1 E[S_1] - \lambda_2 E[S_2])} &= \lambda_1 E[S_1] W_Q^1 + \lambda_2 E[S_2] W_Q^2 \\ &+ \frac{\lambda_1 E[S_1^2]}{2} + \frac{\lambda_2 E[S_2^2]}{2} \\ &= W_Q^1 + \lambda_2 E[S_2] W_Q^2 \quad \text{[from Equation (8.44)]} \end{aligned}$$

Now, using Equation (8.45), we obtain

$$\lambda_2 E[S_2] W_Q^2 = \frac{\lambda_1 E[S_1^2] + \lambda_2 E[S_2^2]}{2} \left[ \frac{1}{1 - \lambda_1 E[S_1] - \lambda_2 E[S_2]} - \frac{1}{1 - \lambda_1 E[S_1]} \right]$$

or

$$W_Q^2 = \frac{\lambda_1 E[S_1^2] + \lambda_2 E[S_2^2]}{2(1 - \lambda_1 E[S_1] - \lambda_2 E[S_2])(1 - \lambda_1 E[S_1])} \quad (8.46)$$

**Remarks** (i) Note that from Equation (8.45), the condition for  $W_Q^1$  to be finite is that  $\lambda_1 E[S_1] < 1$ , which is independent of the type 2 parameters. (Is this intuitive?) For  $W_Q^2$  to be finite, we need, from Equation (8.46), that

$$\lambda_1 E[S_1] + \lambda_2 E[S_2] < 1$$

Since the arrival rate of all customers is  $\lambda = \lambda_1 + \lambda_2$ , and the average service time of a customer is  $(\lambda_1/\lambda)E[S_1] + (\lambda_2/\lambda)E[S_2]$ , the preceding condition is just that the average arrival rate be less than the average service rate.

(ii) If there are  $n$  types of customers, we can solve for  $V^j$ ,  $j = 1, \dots, n$ ; in a similar fashion. First, note that the total amount of work in the system of customers of types  $1, \dots, j$  is independent of the internal priority rule concerning types  $1, \dots, j$  and only depends on the fact that each of them is given priority over any customers of types  $j + 1, \dots, n$ . (Why is this? Reason it out!) Hence,  $V^1 + \dots + V^j$  is the same as it would be if types  $1, \dots, j$  were considered as a single type I priority class and types  $j + 1, \dots, n$  as a single type II priority class. Now, from Equations (8.43) and (8.45),

$$V^1 = \frac{\lambda_1 E[S_1^2] + \lambda_1 \lambda_{II} E[S_1] E[S_{II}^2]}{2(1 - \lambda_1 E[S_1])}$$

where

$$\begin{aligned} \lambda_I &= \lambda_1 + \dots + \lambda_j, \\ \lambda_{II} &= \lambda_{j+1} + \dots + \lambda_n, \\ E[S_I] &= \sum_{i=1}^j \frac{\lambda_i}{\lambda_I} E[S_i], \\ E[S_I^2] &= \sum_{i=1}^j \frac{\lambda_i}{\lambda_I} E[S_i^2], \\ E[S_{II}^2] &= \sum_{i=j+1}^n \frac{\lambda_i}{\lambda_{II}} E[S_i^2] \end{aligned}$$

Hence, as  $V^1 = V^1 + \dots + V^j$ , we have an expression for  $V^1 + \dots + V^j$ , for each  $j = 1, \dots, n$ , which then can be solved for the individual  $V^1, V^2, \dots, V^n$ . We now can obtain  $W_Q^i$  from Equation (8.43). The result of all this (which we leave for an exercise) is that

$$W_Q^i = \frac{\lambda_1 E[S_1^2] + \dots + \lambda_n E[S_n^2]}{2 \prod_{j=i-1}^i (1 - \lambda_1 E[S_1] - \dots - \lambda_j E[S_j])}, \quad i = 1, \dots, n \quad (8.47)$$

### 8.7. The Model G/M/1

The model G/M/1 assumes that the times between successive arrivals have an arbitrary distribution  $G$ . The service times are exponentially distributed with rate  $\mu$  and there is a single server.

The immediate difficulty in analyzing this model stems from the fact that the number of customers in the system is not informative enough to serve as a state space. For in summarizing what has occurred up to the present we would need to know not only the number in the system, but also the amount of time that has elapsed since the last arrival (since  $G$  is not memoryless). (Why need we not be concerned with the amount of time the person being served has already spent in service?) To get around this problem we shall only look at the system when a customer arrives; and so let us define  $X_n$ ,  $n \geq 1$ , by

$X_n \equiv$  the number in the system as seen by the  $n$ th arrival

It is easy to see that the process  $\{X_n, n \geq 1\}$  is a Markov chain. To compute the transition probabilities  $P_{ij}$  for this Markov chain let us first note that, as long as there are customers to be served, the number of services in any length of time  $t$  is a Poisson random variable with mean  $\mu t$ . This is true since the time between successive services is exponential and, as we

know, this implies that the number of services thus constitutes a Poisson process. Hence,

$$P_{i,i+1-j} = \int_0^\infty e^{-\mu t} \frac{(\mu t)^j}{j!} dG(t), \quad j = 0, 1, \dots, i$$

which follows since if an arrival finds  $i$  in the system, then the next arrival will find  $i + 1$  minus the number served, and the probability that  $j$  will be served is easily seen to equal the right side of the above (by conditioning on the time between the successive arrivals).

The formula for  $P_{i0}$  is a little different (it is the probability that at least  $i + 1$  Poisson events occur in a random length of time having distribution  $G$ ) and can be obtained from

$$P_{i0} = 1 - \sum_{j=0}^i P_{i,i+1-j}$$

The limiting probabilities  $\pi_k$ ,  $k = 0, 1, \dots$ , can be obtained as the unique solution of

$$\pi_k = \sum_i \pi_i P_{ik}, \quad k \geq 0$$

$$\sum_k \pi_k = 1$$

which, in this case, reduce to

$$\pi_k = \sum_{i=k-1}^\infty \pi_i \int_0^\infty e^{-\mu t} \frac{(\mu t)^{i+1-k}}{(i+1-k)!} dG(t), \quad k \geq 1$$

$$\sum_0^\infty \pi_k = 1 \tag{8.48}$$

(We have not included the equation  $\pi_0 = \sum \pi_i P_{i0}$  since one of the equations is always redundant.)

To solve the above, let us try a solution of the form  $\pi_k = c\beta^k$ . Substitution into Equation (8.48) leads to

$$c\beta^k = c \sum_{i=k-1}^\infty \beta^i \int_0^\infty e^{-\mu t} \frac{(\mu t)^{i+1-k}}{(i+1-k)!} dG(t)$$

$$= c \int_0^\infty e^{-\mu t} \beta^{k-1} \sum_{i=k-1}^\infty \frac{(\beta \mu t)^{i+1-k}}{(i+1-k)!} dG(t) \tag{8.49}$$

However,

$$\sum_{i=k-1}^\infty \frac{(\beta \mu t)^{i+1-k}}{(i+1-k)!} = \sum_{j=0}^\infty \frac{(\beta \mu t)^j}{j!}$$

$$= e^{\beta \mu t}$$

and thus Equation (8.49) reduces to

$$\beta^k = \beta^{k-1} \int_0^\infty e^{-\mu t(1-\beta)} dG(t)$$

or

$$\beta = \int_0^\infty e^{-\mu t(1-\beta)} dG(t) \tag{8.50}$$

The constant  $c$  can be obtained from  $\sum_k \pi_k = 1$ , which implies that

$$c \sum_0^\infty \beta^k = 1$$

or

$$c = 1 - \beta$$

As the  $\pi_k$  is the unique solution to Equation (8.48), and  $\pi_k = (1 - \beta)\beta^k$  satisfies, it follows that

$$\pi_k = (1 - \beta)\beta^k, \quad k = 0, 1, \dots$$

where  $\beta$  is the solution of Equation (8.50). [It can be shown that if the mean of  $G$  is greater than the mean service time  $1/\mu$ , then there is a unique value of  $\beta$  satisfying Equation (8.50) which is between 0 and 1.] The exact value of  $\beta$  usually can only be obtained by numerical methods.

As  $\pi_k$  is the limiting probability that an arrival sees  $k$  customers, it is just the  $a_k$  as defined in Section 8.2. Hence,

$$a_k = (1 - \beta)\beta^k, \quad k \geq 0 \tag{8.51}$$

We can obtain  $W$  by conditioning on the number in the system when a customer arrives. This yields

$$W = \sum_k E[\text{time in system} \mid \text{arrival sees } k](1 - \beta)\beta^k$$

$$= \sum_k \frac{k+1}{\mu} (1 - \beta)\beta^k \quad (\text{Since if an arrival sees } k, \text{ then he spends } k+1 \text{ service periods in the system.})$$

$$= \frac{1}{\mu(1 - \beta)} \quad \left( \text{by using } \sum_0^\infty kx^k = \frac{x}{(1-x)^2} \right)$$

and

$$W_Q = W - \frac{1}{\mu} = \frac{\beta}{\mu(1 - \beta)},$$

$$L = \lambda W = \frac{\lambda}{\mu(1 - \beta)}, \tag{8.52}$$

$$L_Q = \lambda W_Q = \frac{\lambda\beta}{\mu(1 - \beta)}$$

where  $\lambda$  is the reciprocal of the mean interarrival time. That is,

$$\frac{1}{\lambda} = \int_0^{\infty} x dG(x)$$

In fact, in exactly the same manner as shown for the  $M/M/1$  in Section 8.3.1 and Exercise 4 we can show that

$W^*$  is exponential with rate  $\mu(1 - \beta)$ ,

$$W_Q^* = \begin{cases} 0 & \text{with probability } 1 - \beta \\ \text{exponential with rate } \mu(1 - \beta) & \text{with probability } \beta \end{cases}$$

where  $W^*$  and  $W_Q^*$  are the amounts of time that a customer spends in system and queue, respectively (their means are  $W$  and  $W_Q$ ).

Whereas  $a_k = (1 - \beta)\beta^k$  is the probability that an arrival sees  $k$  in the system, it is not equal to the proportion of time during which there are  $k$  in the system (since the arrival process is not Poisson). To obtain the  $P_k$  we first note that the rate at which the number in the system changes from  $k - 1$  to  $k$  must equal the rate at which it changes from  $k$  to  $k - 1$  (why?). Now the rate at which it changes from  $k - 1$  to  $k$  is equal to the arrival rate  $\lambda$  multiplied by the proportion of arrivals finding  $k - 1$  in the system. That is,

$$\text{rate number in system goes from } k - 1 \text{ to } k = \lambda a_{k-1}$$

Similarly, the rate at which the number in the system changes from  $k$  to  $k - 1$  is equal to the proportion of time during which there are  $k$  in the system multiplied by the (constant) service rate. That is,

$$\text{rate number in system goes from } k \text{ to } k - 1 = P_k \mu$$

Equating these rates yields

$$P_k = \frac{\lambda}{\mu} a_{k-1}, \quad k \geq 1$$

and so, from Equation (8.51),

$$P_k = \frac{\lambda}{\mu} (1 - \beta)\beta^{k-1}, \quad k \geq 1$$

and, as  $P_0 = 1 - \sum_{k=1}^{\infty} P_k$ , we obtain

$$P_0 = 1 - \frac{\lambda}{\mu}$$

**Remark** In the foregoing analysis we guessed at a solution of the stationary probabilities of the Markov chain of the form  $\pi_k = c\beta^k$ , then verified such a solution by substituting in the stationary Equation (8.48). However, it could have been argued directly that the stationary probabilities of the Markov chain are of this form. To do so, define  $\beta_i$  to be the expected number of times that state  $i + 1$  is visited in the Markov chain between two successive visits to state  $i$ ,  $i \geq 0$ . Now it is not difficult to see (and we will let the reader argue it out for him or herself) that

$$\beta_0 = \beta_1 = \beta_2 = \dots = \beta$$

Now it can be shown by using renewal reward processes that

$$\begin{aligned} \pi_{i+1} &= \frac{E[\text{number of visits to state } i + 1 \text{ in an } i - i \text{ cycle}]}{E[\text{number of transitions in an } i - i \text{ cycle}]} \\ &= \frac{\beta_i}{1/\pi_i} \end{aligned}$$

and so,

$$\pi_{i+1} = \beta_i \pi_i = \beta \pi_i, \quad i \geq 0$$

implying, since  $\sum_0^{\infty} \pi_i = 1$ , that

$$\pi_i = \beta^i (1 - \beta), \quad i \geq 0$$

### 8.7.1. The $G/M/1$ Busy and Idle Periods

Suppose that an arrival has just found the system empty—and so initiates a busy period—and let  $N$  denote the number of customers served in that busy period. Since the  $N$ th arrival (after the initiator of the busy period) will also find the system empty, it follows that  $N$  is the number of transitions for the Markov chain (of Section 8.7) to go from state 0 to state 0. Hence,  $1/E[N]$  is the proportion of transitions that take the Markov chain into state 0; or equivalently, it is the proportion of arrivals that find the system empty. Therefore,

$$E[N] = \frac{1}{a_0} = \frac{1}{1 - \beta}$$

Also, as the next busy period begins after the  $N$ th interarrival, it follows that the cycle time (that is, the sum of a busy and idle period) is equal to the time until the  $N$ th interarrival. In other words, the sum of a busy and idle period can be expressed as the sum of  $N$  interarrival times. Thus, if  $T_i$  is the

$i$ th interarrival time after the busy period begins, then

$$\begin{aligned} E[\text{Busy}] + E[\text{Idle}] &= E\left[\sum_{i=1}^N T_i\right] \\ &= E[N]E[T] \quad (\text{by Wald's equation}) \\ &= \frac{1}{\lambda(1-\beta)} \end{aligned} \tag{8.53}$$

For a second relation between  $E[\text{Busy}]$  and  $E[\text{Idle}]$ , we can use the same argument as in Section 8.5.3 to conclude that

$$1 - P_0 = \frac{E[\text{Busy}]}{E[\text{Idle}] + E[\text{Busy}]}$$

and since  $P_0 = 1 - \lambda/\mu$ , we obtain, upon combining this with (8.53), that

$$\begin{aligned} E[\text{Busy}] &= \frac{1}{\mu(1-\beta)}, \\ E[\text{Idle}] &= \frac{\mu - \lambda}{\lambda\mu(1-\beta)} \end{aligned}$$

### 8.8. Multiserver Queues

By and large, systems that have more than one server are much more difficult to analyze than those with a single server. In Section 8.8.1 we start first with a Poisson arrival system in which no queue is allowed, and then consider in Section 8.8.2 the infinite capacity  $M/M/k$  system. For both of these models we are able to present the limiting probabilities. In Section 8.8.3 we consider the model  $G/M/k$ . The analysis here is similar to that of the  $G/M/1$  (Section 7) except that in place of a single quantity  $\beta$  given as the solution of an integral equation, we have  $k$  such quantities. We end in Section 8.8.4 with the model  $M/G/k$  for which unfortunately our previous technique (used in  $M/G/1$ ) no longer enables us to derive  $W_Q$ , and we content ourselves with an approximation.

#### 8.8.1. Erlang's Loss System

A loss system is a queueing system in which arrivals that find all servers busy do not enter but rather are lost to the system. The simplest such system is the  $M/M/k$  loss system in which customers arrive according to a Poisson

process having rate  $\lambda$ , enter the system if at least one of the  $k$  servers is free, and then spend an exponential amount of time with rate  $\mu$  being served. The balance equations for this system are

State	Rate leave = rate enter
0	$\lambda P_0 = \mu P_1$
1	$(\lambda + \mu)P_1 = 2\mu P_2 + \lambda P_0$
2	$(\lambda + 2\mu)P_2 = 3\mu P_3 + \lambda P_1$
$i, 0 < i < k$	$(\lambda + i\mu)P_i = (i+1)\mu P_{i+1} + \lambda P_{i-1}$
$k$	$k\mu P_k = \lambda P_{k-1}$

Rewriting gives

$$\begin{aligned} \lambda P_0 &= \mu P_1, \\ \lambda P_1 &= 2\mu P_2, \\ \lambda P_2 &= 3\mu P_3, \\ &\vdots \\ \lambda P_{k-1} &= k\mu P_k \end{aligned}$$

or

$$\begin{aligned} P_1 &= \frac{\lambda}{\mu} P_0 \\ P_2 &= \frac{\lambda}{2\mu} P_1 = \frac{(\lambda/\mu)^2}{2} P_0, \\ P_3 &= \frac{\lambda}{3\mu} P_2 = \frac{(\lambda/\mu)^3}{3!} P_0, \\ &\vdots \\ P_k &= \frac{\lambda}{k\mu} P_{k-1} = \frac{(\lambda/\mu)^k}{k!} P_0 \end{aligned}$$

and using  $\sum_0^k P_i = 1$ , we obtain

$$P_i = \frac{(\lambda/\mu)^i / i!}{\sum_{j=0}^k (\lambda/\mu)^j / j!}, \quad i = 0, 1, \dots, k$$

Since  $E[S] = 1/\mu$ , where  $E[S]$  is the mean service time, the preceding can be written as

$$P_i = \frac{(\lambda E[S])^i / i!}{\sum_{j=0}^k (\lambda E[S])^j / j!}, \quad i = 0, 1, \dots, k \tag{8.54}$$



Consider now the same system except that the service distribution is general—that is, consider the  $M/G/k$  with no queue allowed. This model is sometimes called the *Erlang loss system*. It can be shown (though the proof is advanced) that Equation (8.54) (which is called Erlang's loss formula) remains valid for this more general system.

### 8.8.2. The $M/M/k$ Queue

The  $M/M/k$  infinite capacity queue can be analyzed by the balance equation technique. We leave it for the reader to verify that

$$P_i = \begin{cases} \frac{(\lambda/\mu)^i}{i!} \frac{k\mu}{k\mu - \lambda} \sum_{i=0}^{k-1} \frac{(\lambda/\mu)^i}{i!} + \frac{(\lambda/\mu)^k}{k!} \frac{k\mu}{k\mu - \lambda}, & i \leq k \\ \frac{(\lambda/k\mu)^i k^k}{k!} P_0, & i > k \end{cases}$$

We see from the preceding that we need to impose the condition  $\lambda < k\mu$ .

### 8.8.3. The $G/M/k$ Queue

In this model we again suppose that there are  $k$  servers, each of which serves at an exponential rate  $\mu$ . However, we now allow the time between successive arrivals to have an arbitrary distribution  $G$ . To ensure that a steady-state (or limiting) distribution exists, we assume the condition  $1/\mu_G < k\mu$  where  $\mu_G$  is the mean of  $G$ .\*

The analysis for this model is similar to that presented in Section 8.7 for the case  $k = 1$ . Namely, to avoid having to keep track of the time since the last arrival, we look at the system only at arrival epochs. Once again, if we define  $X_n$  as the number in the system at the moment of the  $n$ th arrival, then  $\{X_n, n \geq 0\}$  is a Markov chain.

To derive the transition probabilities of the Markov chain, it helps to first note the relationship

$$X_{n+1} = X_n + 1 - Y_n, \quad n \geq 0$$

\* It follows from renewal theory (Proposition 7.1) that customers arrive at rate  $1/\mu_G$ , and as the maximum service rate is  $k\mu$ , we clearly need that  $1/\mu_G < k\mu$  for limiting probabilities to exist.

where  $Y_n$  denotes the number of departures during the interarrival time between the  $n$ th and  $(n+1)$ st arrival. The transition probabilities  $P_{ij}$  can now be calculated as follows:

#### Case (i) $j > i + 1$ .

In this case it easily follows that  $P_{ij} = 0$ .

#### Case (ii) $j \leq i + 1 \leq k$ .

In this case if an arrival finds  $i$  in the system, then as  $i < k$  the new arrival will also immediately enter service. Hence, the next arrival will find  $j$  if of the  $i+1$  services exactly  $i+1-j$  are completed during the interarrival time. Conditioning on the length of this interarrival time yields

$$\begin{aligned} P_{ij} &= P\{i+1-j \text{ of } i+1 \text{ services are completed in an interarrival time}\} \\ &= \int_0^\infty P\{i+1-j \text{ of } i+1 \text{ are completed} \mid \text{interarrival time is } t\} dG(t) \\ &= \int_0^\infty \binom{i+1}{j} (1 - e^{-\mu t})^{i+1-j} (e^{-\mu t})^j dG(t) \end{aligned}$$

where the last equality follows since the number of service completions in a time  $t$  will have a binomial distribution.

#### Case (iii) $i+1 \geq j \geq k$

To evaluate  $P_{ij}$  in this case we first note that when all servers are busy, the departure process is a Poisson process with rate  $k\mu$  (why?). Hence, again conditioning on the interarrival time we have

$$\begin{aligned} P_{ij} &= P\{i+1-j \text{ departures}\} \\ &= \int_0^\infty P\{i+1-j \text{ departures in time } t\} dG(t) \\ &= \int_0^\infty e^{-k\mu t} \frac{(k\mu t)^{i+1-j}}{(i+1-j)!} dG(t) \end{aligned}$$

#### Case (iv) $i+1 \geq k > j$

In this case since when all servers are busy the departure process is a Poisson process, it follows that the length of time until there will only be  $k$  in the system will have a gamma distribution with parameters  $i+1-k$ ,  $k\mu$  (the time until  $i+1-k$  event of a Poisson process with rate  $k\mu$  occur is gamma distributed with parameters  $i+1-k$ ,  $k\mu$ ). Conditioning first on the interarrival time and then on the time until there are only  $k$  in the system

(call this latter random variable  $T_k$ ) yields

$$\begin{aligned}
 P_{ij} &= \int_0^\infty P\{i + 1 - j \text{ departures in time } t\} dG(t) \\
 &= \int_0^\infty \int_0^t P\{i + 1 - j \text{ departures in } t \mid T_k = s\} k\mu e^{-k\mu s} \frac{(k\mu s)^{i-k}}{(i-k)!} ds dG(t) \\
 &= \int_0^\infty \int_0^t \binom{k}{j} (1 - e^{-\mu(t-s)})^{k-j} (e^{-\mu(t-s)})^j k\mu e^{-k\mu s} \frac{(k\mu s)^{i-k}}{(i-k)!} ds dG(t)
 \end{aligned}$$

where the last equality follows since of the  $k$  people in service at time  $s$  the number whose service will end by time  $t$  is binomial with parameters  $k$  and  $1 - e^{-\mu(t-s)}$ .

We now can verify either by a direct substitution into the equations  $\pi_j = \sum_i \pi_i P_{ij}$ , or by the same argument as presented in the remark at the end of Section 8.7, that the limiting probabilities of this Markov chain are of the form

$$\pi_{k-1+j} = c\beta^j, \quad j = 0, 1, \dots$$

Substitution into any of the equations  $\pi_j = \sum_i \pi_i P_{ij}$  when  $j > k$  yields that  $\beta$  is given as the solution of

$$\beta = \int_0^\infty e^{-k\mu t(1-\beta)} dG(t)$$

The values  $\pi_0, \pi_1, \dots, \pi_{k-2}$ , can be obtained by recursively solving the first  $k - 1$  of the steady-state equations, and  $c$  can then be computed by using  $\sum_0^\infty \pi_i = 1$ .

If we let  $W_Q^*$  denote the amount of time that a customer spends in queue, then in exactly the same manner as in  $G/M/1$  we can show that

$$W_Q^* = \begin{cases} 0, & \text{with probability } \sum_0^{k-1} \pi_i = 1 - \frac{c\beta}{1-\beta} \\ \text{Exp}(k\mu(1-\beta)), & \text{with probability } \sum_k^\infty \pi_i = \frac{c\beta}{1-\beta} \end{cases}$$

where  $\text{Exp}(k\mu(1-\beta))$  is an exponential random variable with rate  $k\mu(1-\beta)$ .

### 8.8.4. The $M/G/k$ Queue

In this section we consider the  $M/G/k$  system in which customers arrive at a Poisson rate  $\lambda$  and are served by any of  $k$  servers, each of whom has the service distribution  $G$ . If we attempt to mimic the analysis presented in

Section 8.5 for the  $M/G/1$  system, then we would start with the basic identity

$$V = \lambda E[S] W_Q + \lambda E[S^2]/2 \tag{8.55}$$

and then attempt to derive a second equation relating  $V$  and  $W_Q$ .

Now if we consider an arbitrary arrival, then we have the following identity:

$$\begin{aligned}
 &\text{work in system when customer arrives} \\
 &= k \times \text{time customer spends in queue} + R \tag{8.56}
 \end{aligned}$$

where  $R$  is the sum of the remaining service times of all other customers in service at the moment when our arrival enters service.

The foregoing follows since while the arrival is waiting in queue, work is being processed at a rate  $k$  per unit time (since all servers are busy). Thus, an amount of work  $k \times$  time in queue is processed while he waits in queue. Now, all of this work was present when he arrived and in addition the remaining work on those still being served when he enters service was also present when he arrived—so we obtain Equation (8.56). For an illustration, suppose that there are three servers all of whom are busy when the customer arrives. Suppose, in addition, that there are no other customers in the system and also that the remaining service times of the three people in service are 3, 6, and 7. Hence, the work seen by the arrival is  $3 + 6 + 7 = 16$ . Now the arrival will spend 3 time units in queue, and at the moment he enters service, the remaining times of the other two customers are  $6 - 3 = 3$  and  $7 - 3 = 4$ . Hence,  $R = 3 + 4 = 7$  and as a check of Equation (8.56) we see that  $16 = 3 \times 3 + 7$ .

Taking expectations of Equation (8.55) and using the fact that Poisson arrivals see time averages, we obtain

$$V = kW_Q + E[R]$$

which, along with Equation (8.55), would enable us to solve for  $W_Q$  if we could compute  $E[R]$ . However there is no known method for computing  $E[R]$  and in fact, there is no known exact formula for  $W_Q$ . The following approximation for  $W_Q$  was obtained in Reference 6 by using the foregoing approach and then approximating  $E[R]$ :

$$W_Q \approx \frac{\lambda^k E[S^2] (E[S])^{k-1}}{2(k-1)!(k-\lambda E[S])^2 \left[ \sum_{n=0}^{k-1} \frac{(\lambda E[S])^n}{n!} + \frac{(\lambda E[S])^k}{(k-1)!(k-\lambda E[S])} \right]} \tag{8.57}$$

The preceding approximation has been shown to be quite close to the  $W_Q$  when the service distribution is gamma. It is also exact when  $G$  is exponential.

## Exercises

1. For the  $M/M/1$  queue, compute

- the expected number of arrivals during a service period and
- the probability that no customers arrive during a service period.

Hint: "Condition."

\*2. Machines in a factory break down at an exponential rate of six per hour. There is a single repairman who fixes machines at an exponential rate of eight per hour. The cost incurred in lost production when machines are out of service is \$10 per hour per machine. What is the average cost rate incurred due to failed machines?

3. The manager of a market can hire either Mary or Alice. Mary, who gives service at an exponential rate of 20 customers per hour, can be hired at a rate of \$3 per hour. Alice, who gives service at an exponential rate of 30 customers per hour, can be hired at a rate of \$ $C$  per hour. The manager estimates that, on the average, each customer's time is worth \$1 per hour and should be accounted for in the model. If customers arrive at a Poisson rate of 10 per hour, then

- what is the average cost per hour if Mary is hired? if Alice is hired?
- find  $C$  if the average cost per hour is the same for Mary and Alice.

4. For the  $M/M/1$  queue, show that the probability that a customer spends an amount of time  $x$  or less in queue is given by

$$1 - \frac{\lambda}{\mu}, \quad \text{if } x = 0$$

$$1 - \frac{\lambda}{\mu} + \frac{\lambda}{\mu}(1 - e^{-(\mu-\lambda)x}), \quad \text{if } x > 0$$

5. Two customers move about among three servers. Upon completion of service at server  $i$ , the customer leaves that server and enters service at whichever of the other two servers is free. (Therefore, there are always two busy servers.) If the service times at server  $i$  are exponential with rate  $\mu_i$ ,  $i = 1, 2, 3$ , what proportion of time is server  $i$  idle?

\*6. Show that  $W$  is smaller in an  $M/M/1$  model having arrivals at rate  $\lambda$  and service at rate  $2\mu$  than it is in a two-server  $M/M/2$  model with arrivals at rate  $\lambda$  and with each server at rate  $\mu$ . Can you give an intuitive explanation for this result? Would it also be true for  $W_Q$ ?

7. A group of  $n$  customers moves around among two servers. Upon completion of service, the served customer then joins the queue (or enters service if the server is free) at the other server. All service times are exponential with rate  $\mu$ . Find the proportion of time that there are  $j$  customers at server 1,  $j = 0, \dots, n$ .

8. A facility produces items according to a Poisson process with rate  $\lambda$ . However, it has shelf space for only  $k$  items and so it shuts down production whenever  $k$  items are present. Customers arrive at the facility according to a Poisson process with rate  $\mu$ . Each customer wants one item and will immediately depart either with the item or empty handed if there is no item available.

- Find the proportion of customers that go away empty handed.
- Find the average time that an item is on the shelf.
- Find the average number of items on the shelf.

Suppose now that when a customer does not find any available items it joins the "customers' queue" as long as there are no more than  $n - 1$  other customers waiting at that time. If there are  $n$  waiting customers then the new arrival departs without an item.

- Set up the balance equations.
- In terms of the solution of the balance equations, what is the average number of customers in the system.

9. A group of  $m$  customers frequents a single-server station in the following manner. When a customer arrives, he or she either enters service if the server is free or joins the queue otherwise. Upon completing service the customer departs the system, but then returns after an exponential time with rate  $\theta$ . All service times are exponentially distributed with rate  $\mu$ .

- Define states and set up the balance equations.

In terms of the solution of the balance equations, find

- the average rate at which customers enter the station.
- the average time that a customer spends in the station per visit.

10. Consider a single-server queue with Poisson arrivals and exponential service times having the following variation: Whenever a service is completed a departure occurs only with probability  $\alpha$ . With probability  $1 - \alpha$  the customer, instead of leaving, joins the end of the queue. Note that a customer may be serviced more than once.

- Set up the balance equations and solve for the steady-state probabilities, stating conditions for it to exist.

- (b) Find the expected waiting time of a customer from the time he arrives until he enters service for the first time.
- (c) What is the probability that a customer enters service exactly  $n$  times, for  $n = 1, 2, \dots$ ?
- (d) What is the expected amount of time that a customer spends in service (which does not include the time he spends waiting in line)?

**Hint:** Use (c).

- (e) What is the distribution of the total length of time a customer spends being served?

**Hint:** Is it memoryless?

\*11. A supermarket has two exponential checkout counters, each operating at rate  $\mu$ . Arrivals are Poisson at rate  $\lambda$ . The counters operate in the following way:

- (i) One queue feeds both counters.
- (ii) One counter is operated by a permanent checker and the other by a stock clerk who instantaneously begins checking whenever there are two or more customers in the system. The clerk returns to stocking whenever he completes a service, and there are fewer than two customers in the system.

- (a) Let  $P_n$  = proportion of time there are  $n$  in the system. Set up equations for  $P_n$  and solve.
- (b) At what rate does the number in the system go from 0 to 1? from 2 to 1?
- (c) What proportion of time is the stock clerk checking?

**Hint:** Be a little careful when there is one in the system.

12. Customers arrive at a single-service facility at a Poisson rate of 40 per hour. When two or fewer customers are present, a single attendant operates the facility, and the service time for each customer is exponentially distributed with a mean value of two minutes. However, when there are three or more customers at the facility, the attendant is joined by an assistant and, working together, they reduce the mean service time to one minute. Assuming a system capacity of four customers,

- (a) what proportion of time are both servers free?
- (b) each man is to receive a salary proportional to the amount of time he is actually at work servicing customers, the rate being the same for both. If together they earn \$100 per day, how should this money be split?

13. Consider a sequential-service system consisting of two servers,  $A$  and  $B$ . Arriving customers will enter this system only if server  $A$  is free. If a customer does enter, then he is immediately served by server  $A$ . When his service by  $A$  is completed, he then goes to  $B$  if  $B$  is free, or if  $B$  is busy, he leaves the system. Upon completion of service at server  $B$ , the customer departs. Assuming that the (Poisson) arrival rate is two customers an hour, and that  $A$  and  $B$  serve at respective (exponential) rates of four and two customers an hour,

- (a) what proportion of customers enter the system?
- (b) what proportion of entering customers receive service from  $B$ ?
- (c) what is the average number of customers in the system?
- (d) what is the average amount of time that an entering customer spends in the system?

14. Customers arrive at a two-server system according to a Poisson process having rate  $\lambda = 5$ . An arrival finding server 1 free will begin service with that server. An arrival finding server 1 busy and server 2 free will enter service with server 2. An arrival finding both servers busy goes away. Once a customer is served by either server, he departs the system. The service times at server  $i$  are exponential with rates  $\mu_i$ , where  $\mu_1 = 4$ ,  $\mu_2 = 2$ .

- (a) What is the average time an entering customer spends in the system?
- (b) What proportion of time is server 2 busy?

15. Customers arrive at a two-server station in accordance with a Poisson process with a rate of two per hour. Arrivals finding server 1 free begin service with that server. Arrivals finding server 1 busy and server 2 free begin service with server 2. Arrivals finding both servers busy are lost. When a customer is served by server 1, she then either enters service with server 2 if 2 is free or departs the system if 2 is busy. A customer completing service at server 2 departs the system. The service times at server 1 and server 2 are exponential random variables with respective rates of four and six per hour.

- (a) What fraction of customers do not enter the system?
- (b) What is the average amount of time that an entering customer spends in the system?
- (c) What fraction of entering customers receive service from server 1?

16. Customers arrive at a two-server system at a Poisson rate  $\lambda$ . An arrival finding the system empty is equally likely to enter service with either server. An arrival finding one customer in the system will enter service with the idle server. An arrival finding two others in the system will wait in line for the first free server. An arrival finding three in the system will not enter.

All service times are exponential with rate  $\mu$ , and once a customer is served (by either server), he departs the system.

- Define the states.
- Find the long-run probabilities.
- Suppose a customer arrives and finds two others in the system. What is the expected time he spends in the system?
- What proportion of customers enter the system?
- What is the average time an entering customer spends in the system?

17. There are two types of customers. Type  $i$  customers arrive in accordance with independent Poisson processes with respective rate  $\lambda_1$  and  $\lambda_2$ . There are two servers. A type 1 arrival will enter service with server 1 if that server is free; if server 1 is busy and server 2 is free, then the type 1 arrival will enter service with server 2. If both servers are busy, then the type 1 arrival will go away. A type 2 customer can only be served by server 2; if server 2 is free when a type 2 customer arrives, then the customer enters service with that server. If server 2 is busy when a type 2 arrives, then that customer goes away. Once a customer is served by either server, he departs the system. Service times at server  $i$  are exponential with rate  $\mu_i$ ,  $i = 1, 2$ .

Suppose we want to find the average number of customers in the system.

- Define states.
- Give the balance equations. Do not attempt to solve them.

In terms of the long-run probabilities, what is

- the average number of customers in the system?
- the average time a customer spends in the system?

\*18. Suppose in Exercise 17 we want to find out the proportion of time there is a type 1 customer with server 2. In terms of the long-run probabilities given in Exercise 17, what is

- the rate at which a type 1 customer enters service with server 2?
- the rate at which a type 2 customer enters service with server 2?
- the fraction of server 2's customers that are type 1?
- the proportion of time that a type 1 customer is with server 2?

19. Customers arrive at a single-server station in accordance with a Poisson process with rate  $\lambda$ . All arrivals that find the server free immediately enter service. All service times are exponentially distributed with rate  $\mu$ . An arrival that finds the server busy will leave the system and roam around "in orbit" for an exponential time with rate  $\theta$  at which time it will then return. If the server is busy when an orbiting customer returns, then that customer

returns to orbit for another exponential time with rate  $\theta$  before returning again. An arrival that finds the server busy and  $N$  other customers in orbit will depart and not return. That is,  $N$  is the maximum number of customers in orbit.

- Define states.
- Give the balance equations.

In terms of the solution of the balance equations, find.

- the proportion of all customers that are eventually served.
- the average time that a served customer spends waiting in orbit.

20. Consider the  $M/M/1$  system in which customers arrive at rate  $\lambda$  and the server serves at rate  $\mu$ . However, suppose that in any interval of length  $h$  in which the server is busy there is a probability  $\alpha h + o(h)$  that the server will experience a breakdown, which causes the system to shut down. All customers that are in the system depart, and no additional arrivals are allowed to enter until the breakdown is fixed. The time to fix a breakdown is exponentially distributed with rate  $\beta$ .

- Define appropriate states.
- Give the balance equations.

In terms of the long-run probabilities,

- what is the average amount of time that an entering customer spends in the system?
- what proportion of entering customers complete their service?
- what proportion of customers arrive during a breakdown?

\*21. Reconsider Exercise 20, but this time suppose that a customer that is in the system when a breakdown occurs remains there while the server is being fixed. In addition, suppose that new arrivals during a breakdown period are allowed to enter the system. What is the average time a customer spends in the system?

22. Poisson ( $\lambda$ ) arrivals join a queue in front of two parallel servers  $A$  and  $B$ , having exponential service rates  $\mu_A$  and  $\mu_B$ . When the system is empty, arrivals go into server  $A$  with probability  $\alpha$  and into  $B$  with probability  $1 - \alpha$ . Otherwise, the head of the queue takes the first free server.

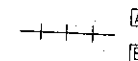


Figure 8.4.

- (a) Define states and set up the balance equations. Do not solve.  
 (b) In terms of the probabilities in part (a), what is the average number in the system? Average number of servers idle?  
 (c) In terms of the probabilities in part (a), what is the probability that an arbitrary arrival will get serviced in  $A$ ?

**23.** In a queue with unlimited waiting space, arrivals are Poisson (parameter  $\lambda$ ) and service times are exponentially distributed (parameter  $\mu$ ). However, the server waits until  $K$  people are present before beginning service on the first customer; thereafter, he services one at a time until all  $K$  units, and all subsequent arrivals, are serviced. The server is then "idle" until  $K$  new arrivals have occurred.

- (a) Define an appropriate state space, draw the transition diagram, and set up the balance equations.  
 (b) In terms of the limiting probabilities, what is the average time a customer spends in queue?  
 (c) What conditions on  $\lambda$  and  $\mu$  are necessary?

**24.** Consider a single-server exponential system in which ordinary customers arrive at a rate  $\lambda$  and have service rate  $\mu$ . In addition, there is a special customer who has a service rate  $\mu_1$ . Whenever this special customer arrives, it goes directly into service (if anyone else is in service, then this person is bumped back into queue). When the special customer is not being serviced, the customer spends an exponential amount of time (with mean  $1/\theta$ ) out of the system.

- (a) What is the average arrival rate of the special customer?  
 (b) Define an appropriate state space and set up balance equations.  
 (c) Find the probability that an ordinary customer is bumped  $n$  time.

**\*25.** Let  $D$  denote the time between successive departures in a stationary  $M/M/1$  queue with  $\lambda < \mu$ . Show, by conditioning on whether or not a departure has left the system empty, that  $D$  is exponential with rate  $\lambda$ .

**Hint:** By conditioning on whether or not the departure has left the system empty we see that

$$D = \begin{cases} \text{Exponential } (\mu), & \text{with probability } \lambda/\mu \\ \text{Exponential } (\lambda) * \text{Exponential } (\mu), & \text{with probability } 1 - \lambda/\mu \end{cases}$$

where  $\text{Exponential } (\lambda) * \text{Exponential } (\mu)$  represents the sum of two independent exponential random variables having rates  $\mu$  and  $\lambda$ . Now use moment-generating functions to show that  $D$  has the required distribution.

Note that the above does not prove that the departure process is Poisson. To prove this we need show not only that the interdeparture times are all exponential with rate  $\lambda$ , but also that they are independent.

**26.** For the tandem queue model verify that

$$P_{n,m} = (\lambda/\mu_1)^n (1 - \lambda/\mu_1) (\lambda/\mu_2)^m (1 - \lambda/\mu_2)$$

satisfies the balance equation (8.15).

**27.** Verify Equation (8.18) for a system of two servers by showing that it satisfies the balance equations for this model.

**28.** Consider a network of three stations. Customers arrive at stations 1, 2, 3 in accordance with Poisson processes having respective rates 5, 10, 15. The service times at the three stations are exponential with respective rates 10, 50, 100. A customer completing service at station 1 is equally likely to (a) go to station 2, (b) go to station 3, or (c) leave the system. A customer departing service at station 2 always goes to station 3. A departure from service at station 3 is equally likely to either go to station 2 or leave the system.

- (i) What is the average number of customers in the system (consisting of all three stations)?  
 (ii) What is the average time a customer spends in the system?

**29.** Consider a closed queueing network consisting of two customers moving among two servers, and suppose that after each service completion the customer is equally likely to go to either server—that is,  $P_{1,2} = P_{2,1} = \frac{1}{2}$ . Let  $\mu_i$  denote the exponential service rate at server  $i$ ,  $i = 1, 2$ .

- (a) Determine the average number of customers at each server.  
 (b) Determine the service completion rate for each server.

**30.** State and prove the equivalent of the arrival theorem for open queueing networks.

**31.** Customers arrive at a single-server station in accordance with a Poisson process having rate  $\lambda$ . Each customer has a value. The successive values of customers are independent and come from a uniform distribution on  $(0, 1)$ . The service time of a customer having value  $x$  is a random variable with mean  $3 + 4x$  and variance 5.

- (a) What is the average time a customer spends in the system?  
 (b) What is the average time a customer having value  $x$  spends in the system?

\*32. Compare the  $M/G/1$  system for first-come, first-served queue discipline with one of last-come, first-served (for instance, in which units for service are taken from the top of a stack). Would you think that the queue size, waiting time, and busy-period distribution differ? What about their means? What if the queue discipline was always to choose at random among those waiting? Intuitively which discipline would result in the smallest variance in the waiting time distribution?

33. In an  $M/G/1$  queue,

- (a) what proportion of departures leave behind 0 work?  
 (b) what is the average work in the system as seen by a departure?

34. For the  $M/G/1$  queue, let  $X_n$  denote the number in the system left behind by the  $n$ th departure.

(a) If

$$X_{n+1} = \begin{cases} X_n - 1 + Y_n, & \text{if } X_n \geq 1 \\ Y_n, & \text{if } X_n = 0 \end{cases}$$

what does  $Y_n$  represent?

(b) Rewrite the preceding as

$$X_{n+1} = X_n - 1 + Y_n + \delta_n \quad (8.58)$$

where

$$\delta_n = \begin{cases} 1, & \text{if } X_n = 0 \\ 0, & \text{if } X_n \geq 1 \end{cases}$$

Take expectations and let  $n \rightarrow \infty$  in Equation (8.58) to obtain

$$E[\delta_\infty] = 1 - \lambda E[S]$$

(c) Square both sides of Equation (8.58), take expectations, and then let  $n \rightarrow \infty$  to obtain

$$E[X_\infty] = \frac{\lambda^2 E[S^2]}{2(1 - \lambda E[S])} + \lambda E[S]$$

(d) Argue that  $E[X_\infty]$ , the average number as seen by a departure, is equal to  $L$ .

\*35. Consider an  $M/G/1$  system in which the first customer in a busy period has service distribution  $G_1$  and all others have distribution  $G_2$ . Let  $C$  denote the number of customers in a busy period, and let  $S$  denote the service time of a customer chosen at random.

Argue that

- (a)  $a_0 = P_0 = 1 - \lambda E[S]$ .  
 (b)  $E[S] = a_0 E[S_1] + (1 - a_0) E[S_2]$  where  $S_i$  has distribution  $G_i$ .  
 (c) Use (a) and (b) to show that  $E[B]$ , the expected length of a busy period, is given by

$$E[B] = \frac{E[S_1]}{1 - \lambda E[S_2]}$$

(d) Find  $E[C]$ .

36. Consider a  $M/G/1$  system with  $\lambda E[S] < 1$ .

- (a) Suppose that service is about to begin at a moment when there are  $n$  customers in the system.  
 (i) Argue that the additional time until there are only  $n - 1$  customers in the system has the same distribution as a busy period.  
 (ii) What is the expected additional time until the system is empty?  
 (b) Suppose that the work in the system at some moment is  $A$ . We are interested in the expected additional time until the system is empty—call it  $E[T]$ . Let  $N$  denote the number of arrivals during the first  $A$  units of time.  
 (i) Compute  $E[T|N]$ .  
 (ii) Compute  $E[T]$ .

37. Carloads of customers arrive at a single-server station in accordance to a Poisson process with rate 4 per hour. The service times are exponentially distributed with rate 20 per hour. If each carload contains either 1, 2, or 3 customers with respective probabilities  $\frac{1}{4}$ ,  $\frac{1}{2}$ ,  $\frac{1}{4}$ , compute the average customer delay in queue.

38. In the two-class priority queueing model of Section 8.6.2, what is  $W_Q$ ? Show that  $W_Q$  is less than it would be under FIFO if  $E[S_1] < E[S_2]$  and greater than under FIFO if  $E[S_1] > E[S_2]$ .

39. In a two-class priority queueing model suppose that a cost of  $C_i$  per unit time is incurred for each type  $i$  customer that waits in queue,  $i = 1, 2$ . Show that type 1 customers should be given priority over type 2 (as opposed to the reverse) if

$$\frac{E[S_1]}{C_1} < \frac{E[S_2]}{C_2}$$

40. Consider the priority queueing model of Section 8.6.2 but now suppose that if a type 2 customer is being served when a type 1 arrives then the type 2 customer is bumped out of service. This is called the preemptive case. Suppose that when a bumped type 2 customer goes back in service his service begins at the point where it left off when he was bumped.

(a) Argue that the work in the system at any time is the same as in the nonpreemptive case.

(b) Derive  $W_Q^1$ .

**Hint:** How do type 2 customers affect type 1's?

(c) Why is it not true that

$$V_Q^2 = \lambda_2 E[S_2] W_Q^2$$

(d) Argue that the work seen by a type 2 arrival is the same as in the nonpreemptive case, and so

$$W_Q^2 = W_Q^2(\text{nonpreemptive}) + E[\text{extra time}]$$

where the extra time is due to the fact that he may be bumped.

(e) Let  $N$  denote the number of times a type 2 customer is bumped. Why is

$$E[\text{extra time} | N] = \frac{NE[S_1]}{1 - \lambda_1 E[S_1]}$$

**Hint:** When a type 2 is bumped, relate the time until he gets back in service to a "busy period."

(f) Let  $S_2$  denote the service time of a type 2. What is  $E[N | S_2]$ ?

(g) Combine the preceding to obtain

$$W_Q^2 = W_Q^2(\text{nonpreemptive}) + \frac{\lambda_1 E[S_1] E[S_2]}{1 - \lambda_1 E[S_1]}$$

\*41. Calculate explicitly (not in terms of limiting probabilities) the average time a customer spends in the system in Exercise 21.

42. In the  $G/M/1$  model if  $G$  is exponential with rate  $\lambda$  show that  $\beta = \lambda/\mu$ .

43. Verify Erlang's loss formula, Equation (8.54), when  $k = 1$ .

44. Verify the formula given for the  $P_i$  of the  $M/M/k$ .

45. In the Erlang loss system suppose the Poisson arrival rate is  $\lambda = 2$ , and suppose there are three servers each of whom has a service distribution that is uniformly distributed over  $(0, 2)$ . What proportion of potential customers is lost?

46. In the  $M/M/k$  system,

- (a) what is the probability that a customer will have to wait in queue?  
 (b) determine  $L$  and  $W$ .

47. Verify the formula for the distribution of  $W_Q^*$  given for the  $G/M/k$  model.

\*48. Consider a system where the interarrival times have an arbitrary distribution  $F$ , and there is a single server whose service distribution is  $G$ . Let  $D_n$  denote the amount of time the  $n$ th customer spends waiting in queue. Interpret  $S_n, T_n$  so that

$$D_{n+1} = \begin{cases} D_n + S_n - T_n, & \text{if } D_n + S_n - T_n \geq 0 \\ 0, & \text{if } D_n + S_n - T_n < 0 \end{cases}$$

49. Consider a model in which the interarrival times have an arbitrary distribution  $F$ , and there are  $k$  servers each having service distribution  $G$ . What condition on  $F$  and  $G$  do you think would be necessary for there to exist limiting probabilities?

## References

1. J. Cohen, "The Single Server Queue," North-Holland, Amsterdam, 1969.
2. R. B. Cooper, "Introduction to Queueing Theory," Second Edition, Macmillan, New York, 1984.
3. D. R. Cox and W. L. Smith, "Queues," Wiley, New York, 1961.
4. F. Kelly, "Reversibility and Stochastic Networks," Wiley, New York, 1979.
5. L. Kleinrock, "Queueing Systems," Vol. I, Wiley, New York, 1975.
6. S. Nozaki and S. Ross, "Approximations in Finite Capacity Multiserver Queues with Poisson Arrivals," *J. Appl. Prob.* 13, 826-834 (1978).
7. L. Takacs, "Introduction to the Theory of Queues," Oxford University Press, London and New York, 1962.
8. H. Tijms, "Stochastic Models: An Algorithmic Approach," Wiley, New York, 1994.
9. P. Whittle, "Systems in Stochastic Equilibrium," Wiley, New York, 1986.
10. Wolff, "Stochastic Modeling and the Theory of Queues," Prentice Hall, New Jersey, 1989.