

Analysis of Queueing Systems with Sample Paths and Simulation

Nicky D. van Foreest

June 28, 2020

CONTENTS

Introduction	v	
0.1 Introduction	v	
0.2 Preliminaries	viii	
1 CONSTRUCTION AND SIMULATION OF QUEUEING SYSTEMS		1
1.1 Poisson Distribution	1	
1.2 Queueing Processes in Discrete-Time	4	
1.3 Exponential Distribution	9	
1.4 Single-server Queueing Process in Continuous Time	11	
2 FROM TRANSIENT TO STEADY-STATE ANALYSIS		15
2.1 Kendall's Notation	15	
2.2 Queueing Processes as Regulated Random Walks	16	
2.3 Rate, Stability and Load	18	
2.4 (Limits of) Empirical Performance Measures	20	
2.5 Graphical Summary	23	
3 APPROXIMATE QUEUEING MODELS		25
3.1 $G/G/c$ Queue: Sakasegawa's Formula	25	
3.2 Setups and Batch Processing	29	
3.3 Non-preemptive Interruptions, Server Adjustments	31	
3.4 Preemptive Interruptions, Server Failures	32	
3.5 Tandem Queues	34	
4 FUNDAMENTAL TOOLS		37
4.1 Renewal Reward Theorem	37	
4.2 Level Crossing and Balance Equations	38	
4.3 Poisson Arrivals See Time Averages	43	
4.4 Little's Law	44	
4.5 Graphical Summary	46	
5 EXACT QUEUEING MODELS		49
5.1 $M/M/1$ queue	49	
5.2 $M(n)/M(n)/1$ Queue	52	
5.3 $M^X/M/1$ Queue: Expected Waiting Time	54	
5.4 $M/G/1$ Queue: Expected Waiting Time	55	
5.5 $M^X/M/1$ Queue Length Distribution	57	
5.6 $M/G/1$ Queue Length Distribution	59	
6 QUEUEING CONTROL AND OPEN NETWORKS		63
6.1 N-policies for the $M/M/1$ queue	63	
6.2 N-policies for the $M/G/1$ queue	65	
6.3 Open Single-Class Product-Form Networks	67	
6.4 On $\lambda = \gamma + \lambda P$	70	
Bibliography	73	
Notation	75	
Formula Sheet	77	
Index	77	

INTRODUCTION

0.1 INTRODUCTION

Motivation and Examples

Queueing systems abound, and the analysis and control of queueing systems are major topics in the control, performance evaluation and optimization of production and service systems.

At my local supermarket, for instance, any customer that joins a queue of 4 or more customers gets his/her groceries for free. Of course, there are some constraints: at least one of the cashier facilities has to be unoccupied by a server and the customers in queue should be equally divided over the cashiers that are open (and perhaps there are some further rules, of which I am unaware). When $\pi(n)$ denotes fraction of customers that ‘see upon arrival’ the system with n customers, the manager that controls the occupation of the cashier positions is focused on keeping $\pi(4) + \pi(5) + \dots$, i.e., the fraction of customers that see upon arrival a queue length exceeding 3, very small. In a sense, this is easy enough: just hire many cashiers. However, the cost of personnel may then outweigh the yearly average cost of paying the customer penalties. Thus, the manager’s problem becomes to plan and control the service capacity in such a way that both the penalties and the personnel cost are small.

Fast food restaurants also deal with many interesting queueing situations. Consider, for instance, the preparation of hamburgers. Typically, hamburgers are made-to-stock, in other words, they are prepared before the actual demand has arrived. Thus, hamburgers in stock can be interpreted as customers in queue waiting for service, where the service time is the time between the arrival of two customers that buy hamburgers. The hamburgers have a typical lifetime, and they have to be scrapped if they remain on the shelf longer than a specified amount of time. Thus, the waiting time of hamburgers has to be closely monitored. Of course, it is easy to achieve zero scrap cost, simply by keeping no stock at all. However, to prevent lost-sales, it is very important to maintain a certain amount of hamburgers in stock. Thus, the manager has to balance the scrap cost against the cost of lost sales. In more formal terms, the problem is to choose a policy to prepare hamburgers such that the cost of excess waiting time (scrap) is balanced against the cost of an empty queue (lost sales).

Service systems, such as hospitals, call centers, courts, and so on, have a certain capacity available to serve customers. The performance of such systems is, in part, measured by the total number of jobs processed per year and the fraction of jobs processed within a certain time frame between receiving and closing the job. Here the problem is to organize the capacity such that the sojourn time, i.e., the typical time a job spends in the system, does not exceed some threshold, and such that the system achieves a certain throughput, i.e., jobs served per year.

Clearly, all of the above systems can be seen as queueing systems that have to be monitored and controlled to achieve a certain performance. The performance analysis of such systems can, typically, be characterized by the following performance measures:

1. The fraction of time $p(n)$ that the system contains n customers. In particular, $1 - p(0)$, i.e., the fraction of time the system contains jobs, is important, as this is a measure of the time-average occupancy of the servers, hence related to personnel cost.

2. The fraction of customers $\pi(n)$ that ‘see upon arrival’ the system with n customers. This measure relates to customer perception and lost sales, i.e., fractions of arriving customers that do not enter the system.
3. The average, variance, and/or distribution of the waiting time.
4. The average, variance, and/or distribution of the number of customers in the system.

Here the system can be anything that is capable of holding jobs, such as a queue, the server(s), an entire court, patients waiting for an MRI scan in a hospital, and so on.

It is important to realize that a queueing system can, typically, be decomposed into *two subsystems*: the queue itself and the service system. Thus, we are concerned with three types of waiting: waiting in queue, i.e., *queueing time*, waiting while being in service, i.e., the *service time*, and the total waiting time in the system, i.e., the *sojourn time*.

Organization

In these notes, we will be primarily concerned with making models of queueing systems such that we can compute or estimate the above-mentioned performance measures.

In Chapter 1 we construct queueing systems in discrete time and continuous time, and by implementing these models in code, we can simulate and analyze such systems. Simulation allows us to analyze many realistic queueing systems, while such systems are often (way) too hard to analyze by mathematical tools. Consider, for example, the service process at a check-in desk of an airline company. Business customers and economy customers are served by two separate queueing systems: the business customers are served by one server, say, while the economy class customers by three servers, say. What would happen to the sojourn time of the business customers if their server would be allowed to serve economy class customers when the business queue is empty? For the analysis of such complicated control policies, simulation appears to be the most natural approach.

Notwithstanding the power of simulation, it is often hard to obtain structural understanding into the behavior of queueing systems. Mathematical models, whether exact or approximate, are much more useful to help reason about and improve queueing systems, mainly because they offer insights into scaling laws, such as how the average waiting time depends on average service times or variability of such times. The aim of Chapter 3 is to apply Sakasegawa’s formula to understand how different production and service situations are affected by the system parameters such as service speed, batching rules, and outages. In passing we will use some general tools of probability theory, which proves useful for the sequel of the book.

In Chapter 5 we focus on exact models for single-station queueing systems, and provide the motivation that underlies Sakasegawa’s formula. The main idea here is to consider the *sample paths of a queueing process*, and assume that a typical sample path captures the ‘normal’ stochastic behavior of the system. This sample-path approach has two advantages. In the first place, many of the theoretical results follow from very concrete aspects of these sample paths. Second, the analysis of sample-paths carries over right away to simulation. In fact, simulation of a queueing system offers us one (or more) sample path(s), and based on such sample paths, we derive behavioral and statistical properties of the system. Thus, sample paths form a direct bridge between simulation on the one hand and mathematical analysis on the other.

Chapter 6 extends the models of the previous chapters to controlled queueing systems and queueing networks. The analysis of these systems requires a combination of material of

the previous chapters, but also other mathematical tools such as difference and differential equations, and non-negative matrices. As such, it also provides a stepping stone to the many-fold extensions of the theory to Markov decision theory, optimization, dynamic programming, and so on.

In our discussions we focus on obtaining an intuitive understanding of the analytical tools. For proofs and/or more extensive results we refer to the following books.

1. [Bolch et al. \[2006\]](#)
2. [El-Taha and Stidham Jr. \[1998\]](#)
3. [Tijms \[1994\]](#) and/or [Tijms \[2003\]](#)
4. [Capiński and Zastawniak \[2003\]](#)

Exercises

The main text contains hardly any examples or derivations. Instead, the exercises provide the material to *illustrate* the material and help the reader study the material. Also, a substantial part of exercises consists of consistency checks, to show how new results reduce to old results; like this these exercises also provide relations between various parts of the text. Typically, such checks are trivial; however, the algebra can be quite difficult at times. Another part of the exercises form a set of questions any student should ask while studying the material (even though asking good questions is difficult). For this reason, the reader is urged to try to make as many exercises as possible. Finally, the intention of the exercises is not to be easy.

The companion document gives hints and solutions to all problems. The solutions spell out nearly every intermediate step. For most of you all, this detail is not necessary, but over the years I got many questions like: "how do you go from 'here' to 'there'?" As service, I then added such intermediate steps. As a consequence, the companion is quite extensive. The companion document also contains many additional simple exercises and old exam questions.

Exercises marked as 'not obligatory' are interesting, but (too) hard. I will not use this in an exam.

Acknowledgements

I would like to acknowledge dr. J.W. Nieuwenhuis for our many discussions on queueing theory. To convince him about the more formal aspects, sample-path arguments proved very useful. Finally, I thank my students for submitting many improvements via github. It's very motivating to see a book like this turn into a joint piece of work.

0.2 PRELIMINARIES

Here is an overview of concepts you are supposed to have seen in earlier courses. We will use these concepts over and over in the rest of the book.

We use the notation:

$$\begin{aligned} [x]^+ &= \max\{x, 0\}, \\ f(x-) &= \lim_{y \uparrow x} f(y), \\ f(x+) &= \lim_{y \downarrow x} f(y), \\ \mathbb{1}_A &= \begin{cases} 1, & \text{if } A \text{ is true,} \\ 0, & \text{if } A \text{ is false.} \end{cases} \end{aligned}$$

The last equation defines the *indicator function*.

We write $f(h) = o(h)$ for a function f to say that f is such that $f(h)/h \rightarrow 0$ as $h \rightarrow 0$. If we write $f(h) = o(h)$ it is implicit that $|h| \ll 1$. We call this *small o notation*.

0.2.1. [0.2.1] Let c be a constant (in \mathbb{R}) and the functions f and g both of $o(h)$. Then show that (1) $f(h) \rightarrow 0$ when $h \rightarrow 0$, (2) $c \cdot f = o(h)$, (3) $f + g = o(h)$, and (4) $f \cdot g = o(h)$.

You should know that:

$$(a + b)^n = \sum_{i=0}^n \binom{n}{i} a^{n-i} b^i, \quad (0.2.1a)$$

$$e^x = \lim_{n \rightarrow \infty} (1 + x/n)^n, \quad (0.2.1b)$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots = \sum_{k=0}^{\infty} \frac{x^k}{k!}, \quad (0.2.1c)$$

$$\sum_{n=0}^N \alpha^n = \frac{1 - \alpha^{N+1}}{1 - \alpha}. \quad (0.2.1d)$$

You should know that for a non-negative, integer-valued random variable X with *probability mass function* $f(k) = P(X = k)$,

$$X = \sum_{n=0}^{\infty} X \mathbb{1}_{X=n} = \sum_{n=0}^{\infty} n \mathbb{1}_{X=n}, \quad (0.2.2a)$$

$$E[X] = E\left[\sum_{n=0}^{\infty} n \mathbb{1}_{X=n}\right] = \sum_{n=0}^{\infty} n E[\mathbb{1}_{X=n}] = \sum_{n=0}^{\infty} n f(n), \quad (0.2.2b)$$

$$E[g(X)] = \sum_{n=0}^{\infty} g(n) f(n), \quad (0.2.2c)$$

For general random variables X and Y :

$$E[\mathbb{1}_{X \leq x}] = P(X \leq x), \quad (0.2.3)$$

$$V[X] = E[X^2] - (E[X])^2. \quad (0.2.4)$$

$$E[X + Y] = E[X] + E[Y], \quad (0.2.5)$$

$$V[X + Y] = V[X] + V[Y], \quad \text{if } X \text{ and } Y \text{ are independent.} \quad (0.2.6)$$

0.2.2. [0.2.5] Define the survivor function of X as $G(k) = P(X > k)$. Show that

$$G(k) = \sum_{m=0}^{\infty} \mathbb{1}_{m>k} f(m).$$

As you will see below, this idea makes the computation of certain expressions quite a bit easier.

0.2.3. [0.2.6] Express the probability mass $f(k)$ and the survivor function $G(k)$ in terms of the distribution function $F(k) = P(X \leq k)$ of X .

0.2.4. [0.2.9] Use indicator functions to prove that $\sum_{i=0}^{\infty} iG(i) = E[X^2]/2 - E[X]/2$.

Let X be a continuous non-negative random variable with distribution function F . We write

$$E[X] = \int_0^{\infty} x dF(x)$$

for the expectation of X . Here $dF(x)$ acts as a (sort of) shorthand for $f(x)dx$ ¹. Recall that

$$E[g(X)] = \int_0^{\infty} g(x) dF(x).$$

0.2.5. [0.2.10] Use indicator functions to prove that $E[X] = \int_0^{\infty} x dF(x) = \int_0^{\infty} G(y) dy$, where $G(x) = 1 - F(x)$.

You should be able to use indicator functions and integration by parts to show that $E[X^2] = 2 \int_0^{\infty} yG(y) dy$, where $G(x) = 1 - F(x)$, provided the second moment exists.

0.2.6. [0.2.12] Show that $E[X^2]/2 = \int_0^{\infty} yG(y) dy$ for a continuous non-negative random variable X with survivor function G .

You should know that the *moment-generating function* $M_X(s)$ of a random variable X and $s \in \mathbb{R}$ sufficiently small is defined as:

$$M_X(s) = E[e^{sX}]. \quad (0.2.7a)$$

Moreover, $M_X(s)$ uniquely characterizes the distribution of X . From this definition it follows that:

$$E[X] = M'_X(0) = \left. \frac{dM_X(s)}{ds} \right|_{s=0}, \quad (0.2.7b)$$

$$E[X^2] = M''_X(0), \quad (0.2.7c)$$

$$M_{X+Y}(s) = M_X(s) \cdot M_Y(s), \quad \text{if } X \text{ and } Y \text{ are independent.} \quad (0.2.7d)$$

To help you recall the concept of *conditional probability* consider the following question.

0.2.7. [0.2.14] We have one gift to give to one out of three children. As we cannot divide the gift into parts, we decide to let 'fate decide'. That is, we choose a random number in the set $\{1, 2, 3\}$. The first child that guesses this number wins the gift. Show that the probability of winning the gift is the same for each child.

You should know that:

$$P(A|B) = \frac{P(AB)}{P(B)}, \quad \text{if } P(B) > 0, \quad (0.2.8a)$$

$$P(A) = \sum_{i=1}^n P(AB_i) = \sum_{i=1}^n P(A|B_i)P(B_i), \quad \text{if } A = \bigcup_{i=1}^n B_i \text{ and } P(B_i > 0) \text{ for all } i. \quad (0.2.8b)$$

¹ For the interested reader, $\int x dF(x)$ is a Lebesgue-Stieltjes integral with respect to the measure induced by the distribution function F , see the literature for further background.

CONSTRUCTION AND SIMULATION OF QUEUEING SYSTEMS

The first step to analyze a queueing system is to model it. And for this, there is often not a better start than to build a simulation model. Thus, the aim of this chapter is to teach you how to construct and simulate queueing processes.

In Section 1.2 we build discrete-time models of queueing systems, which means that we use the number of jobs that arrive and can be served in a period to construct the queueing process. Such a period can be an hour, or a day; in fact, any amount of time that makes sense in the context in which the model will be used. Typically we model the number of arrivals and potential services as random variables, and in many practical settings it is reasonable to take the number of arrivals in a period as Poisson distributed. This being the case, we consider the Poisson distribution in Section 1.1, and once we have an understanding of this, we can use random number generators to generate (Poisson distributed) random numbers of arrivals and services to drive the simulator.

In Section 1.4 we focus on constructing queueing processes in continuous time. In this setting, the inter-arrival times and service times of individual jobs become of importance, and then exponentially distributed random variables play a fundamental role. We therefore discuss the properties of the exponential distribution in Section 1.3. We also mention the interesting and close relationship between the exponential distribution and the Poisson distribution.

As will become apparent, both types of constructing queueing processes, the discrete-time or continuous-time models, are easy to implement as computer programs. We include a large number of exercises to show you the astonishing diversity of queueing systems that can be analyzed by simulation. In passing, we develop a number of performance measures to provide insight into the (transient and long-run average) behavior of queueing processes.

We expect you to *know all topics* summarized in Section 0.2; we use these extensively in Section 1.1 and Section 1.3, as well as in any other section that introduces some theory.

1.1 POISSON DISTRIBUTION

In this section, we provide motivation for the use of the Poisson process as an arrival process of customers or jobs at a shop, service station, or machine to receive service. In the exercises we derive numerous properties of this exceedingly important distribution; in the rest of the book we will use these results time and again.

Consider a stream of customers that enter a shop over time. Let us write $N(t)$ for the number of customers that enter during the time interval $[0, t]$ and $N(s, t) = N(t) - N(s)$ for the number that arrive in the time period $(s, t]$. Clearly, as we do not know in advance how many customers will enter, we model the set $\{N(t), t \geq 0\}$ as a family of random variables.

Our first assumption is that the rate at which customers enter stays constant over time. Then it is reasonable to assume that the expected number of arrivals is proportional to the length of the interval. Hence, it is reasonable to assume that there exists some constant λ such that

$$E[N(s, t)] = \lambda(t - s). \quad (1.1.1)$$

The constant λ is called the *arrival rate* of the arrival process.

The second assumption is that the process $N_\lambda = \{N(t), t \geq 0\}$ has *stationary and independent increments*. Stationarity means that the distributions of the number of arrivals are the same for all intervals of equal length, that is, $N(s, t]$ has the same distribution as $N(u, v]$ if $t - s = v - u$. Independence means, roughly speaking, that knowing that $N(s, t] = n$, does not help to make any predictions about the value of $N(u, v]$ if the intervals $(s, t]$ and $(u, v]$ do not overlap.

To find the distribution of $N(t)$ for some given t , let us split the interval $[0, t]$ into n sub-intervals, all of equal length, and ask: ‘What is the probability that a customer arrives in some given sub-interval?’ By our first assumption, the arrival rate is constant over time. Therefore, the probability p of an arrival in each interval should be constant. Moreover, if the time intervals are very small, we can safely neglect the probability that two or more customers arrive in one interval.

As a consequence, then, we can model the occurrence of an arrival in some period i as a Bernoulli distributed random variable B_i . We assume that the $B_i, i = 1, \dots$ are independent and identically distributed (*i.i.d.*) as the common random variable B such that $p = P(B = 1)$ and $1 - p = P(B = 0)$. The total number of arrivals $N_n(t)$ that occur in n intervals is then *binomially distributed*, i.e.,

$$P(N_n(t) = k) = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (1.1.2)$$

If we take $n \rightarrow \infty$, $p \rightarrow 0$ such that $np = \lambda t$, then $N_n(t)$ converges (in distribution) to a *Poisson distributed* random variable $N(t)$, i.e.,

$$P(N(t) = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}, \quad (1.1.3)$$

and then we write $N(t) \sim P(\lambda t)$.

We call the process $N_\lambda = \{N(t)\}$ a *Poisson process* with rate λ when N_λ is stationary, has independent increments, and its elements $N(t) \sim P(\lambda t)$ for all t . Observe that the process N_λ is a much more complicated object than a Poisson distributed random variable. The process is an uncountable set of random variables indexed by $t \in \mathbb{R}^+$, not just *one* random variable.

In the remainder of this section we derive a number of properties of the Poisson process that we will use time and again.

1.1.1. [1.1.4] Show that

$$P(N(t+h) = n | N(t) = n) = 1 - \lambda h + o(h)$$

when $N(t) \sim P(\lambda t)$ and h is small.

1.1.2. [1.1.5] Show that

$$P(N(t+h) = n+1 | N(t) = n) = \lambda h + o(h)$$

when $N(t) \sim P(\lambda t)$ and h is small.

1.1.3. [1.1.6] Show that if $N(t) \sim P(\lambda t)$, we have for small h ,

$$P(N(t+h) \geq n+2 | N(t) = n) = o(h).$$

1.1.4. [1.1.8] Show that if $N(t) \sim P(\lambda t)$, the expected number of arrivals during $[0, t]$ is

$$E[N(t)] = \lambda t.$$

1.1.5. [1.1.9] Show that $E[(N(t))^2] = \lambda^2 t^2 + \lambda t$ if $N(t) \sim P(\lambda t)$.

1.1.6. [1.1.12] Use the moment-generating function of $N(t) \sim P(\lambda t)$ to compute $E[N(t)]$ and $V[N(t)]$.

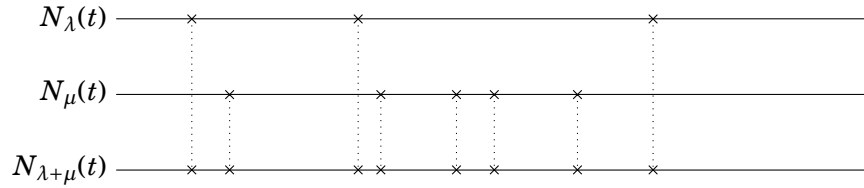
Define the *square coefficient of variation* (SCV) of a random variable X as

$$C^2 = \frac{V[X]}{(E[X])^2}. \quad (1.1.4)$$

As will become clear later, the SCV is a very important concept in queueing theory. Memorize it as a measure of *relative variability*.

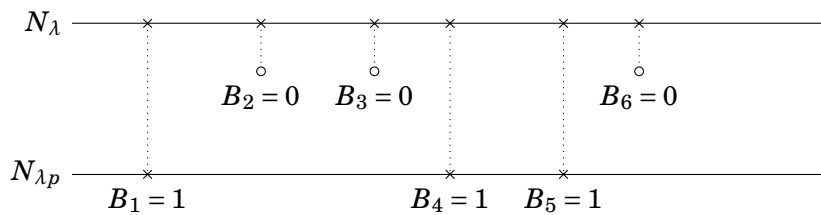
1.1.7. [1.1.13] Show that the SCV of $N(t) \sim P(\lambda t)$ is equal to $1/(\lambda t)$. What does this mean for t large?

Merging Poisson processes occurs often in practice. We have two Poisson processes, for instance, the arrival processes N_λ of men and N_μ of women at a shop. In the figure below, each cross represents an arrival; in the upper line it corresponds to a man, in the middle line to a woman, and in the lower line to an arrival of a general customer at the shop. Thus, the shop ‘sees’ the superposition of these two arrival processes. In fact, this merged process $N_{\lambda+\mu}$ is also a Poisson process with rate $\lambda + \mu$.



1.1.8. [1.1.14] If the Poisson arrival processes N_λ and N_μ are independent, show with a conditioning argument that $N_\lambda + N_\mu$ is a Poisson process with rate $\lambda + \mu$.

Besides merging Poisson streams, we can also consider the concept of *splitting*, or *thinning*, a stream into sub-streams, as follows. Model the stream of people passing by a shop as a Poisson process N_λ . In the figure below these arrivals are marked as crosses at the upper line. With probability p a person decides, independent of anything else, to enter the shop; the crosses at the lower line are the customers that enter the shop. In the figure, the Bernoulli random variable $B_1 = 1$ so that the first passerby enters the shop; the second passerby does not enter as $B_2 = 0$, and so on.



1.1.9. [1.1.18] Show with moment-generating functions that thinning the Poisson process N_λ by means of Bernoulli random variables with success probability p results in a Poisson process $N_{\lambda p}$.

The concepts of merging and thinning are useful to analyze queueing networks. Suppose the departure stream of a machine splits into two sub-streams, e.g., a fraction p of the jobs moves on to another machine and the rest $(1 - p)$ of the jobs leaves the system. Then we can model the arrival stream at the second machine as a thinned stream (with probability p) of the departures of the first machine. Merging occurs where the output streams of various stations arrive at another station.

1.1.10. [1.1.19] *Use moment-generating functions to prove that $N_n(t)$ converges to N , that is, the right-hand side in (1.1.2) converges to the Poisson distribution (1.1.3) when $n \rightarrow \infty, p \rightarrow 0$ such that $pn = \lambda t$ remains constant.*

1.2 QUEUEING PROCESSES IN DISCRETE-TIME

We start with a description of a case to provide motivation to study queueing systems. Then we develop a set of recursions of fundamental importance to construct and simulate queueing systems. With these recursions, we analyze the efficacy of several suggestions to improve the case system. To illustrate the power of this approach, we close the section with a large number of exercises in which you develop recursions for many different queueing situations.

Case

At a mental health department, five psychiatrists do intakes of future patients to determine the best treatment process for the patients. There are complaints about the time patients have to wait for their first intake; the desired waiting time is around two weeks, but the realized waiting time is sometimes more than three months. The organization considers this to be unacceptably long, but... what to do about it?

To reduce the waiting times, the five psychiatrists have various suggestions.

1. Not all psychiatrists have the same amount of time available per week to do intakes. This is not a problem during weeks when all psychiatrists are present; however, psychiatrists tend to take holidays, visit conferences, and so on. So, if the psychiatrist with the most intakes per week would go on leave, this might affect the behavior of the queue length considerably. This raises the question about the difference in the allocation of capacity allotted to the psychiatrists. What are the consequences on the distribution and average of the waiting times if they would all have the same weekly capacity?
2. The psychiatrists tend to plan their holidays consecutively, to reduce the variation in the service capacity. What if they would synchronize their holidays, to the extent possible, rather than spread their holidays?
3. Finally, suppose the psychiatrists would do 2 more intakes per week in busy times and 2 fewer in quiet weeks. Assuming that the system is stable, i.e., the average service capacity is larger than the average demand, then on average the psychiatrists would not do more intakes, i.e., their workload would not increase, but the queue length may be controlled better.

As this case is too hard to analyze by mathematical means, we need to develop a model to simulate the queueing system in discrete time. With this simulator we can evaluate the effect

of these suggestions on reducing the queueing dynamics. Interestingly, the structure of the simulation is very simple, so simple that it is also an exceedingly convincing tool to communicate the results of an analysis of a queueing system to managers (and the like).

Recursions

Let us start with discussing the essentials of the simulation of a queueing system. The easiest way to construct queueing processes is to ‘chop up’ time in periods and develop recursions for the behavior of the queue from period to period. Using fixed-sized periods has the advantage that we do not have to specify specific inter-arrival times or service times of individual customers; only the number of arrivals in a period and the number of potential services are relevant. Note that the length of such a period depends on the context for which the model is developed. For instance, to study queueing processes at a supermarket, a period can consist of 5 minutes, while for a production environment, e.g., a job shop, it can be a day, or even a week.

Let us define:

a_k = number of jobs that arrive *in* period k ,

c_k = the capacity, i.e., the maximal number of jobs that can be served, during period k ,

d_k = number of jobs that depart *in* period k ,

L_k = number of jobs in the system at the *end* of period k .

In the sequel we also call a_k the *size of the batch* arriving in period k . Note that the definition of a_k is a bit subtle: we may assume that the arriving jobs arrive either at the start or at the end of the period. In the first case, the jobs can be served in period k , in the latter case, they *cannot* be served in period k .

Since L_{k-1} is the queue length at the *end* of period $k-1$, it must also be the number of customers at the *start* of period k . Assuming that jobs arriving in period k cannot be served in period k , the number of customers that depart in period k is

$$d_k = \min\{L_{k-1}, c_k\}, \quad (1.2.1a)$$

since only the jobs that are present at the start of the period, i.e., L_{k-1} , can be served if the capacity exceeds the queue length. Now that we know the number of departures, the number at the end of period k is given by

$$L_k = L_{k-1} - d_k + a_k. \quad (1.2.1b)$$

Like this, if we are given L_0 , we can obtain L_1 , and from this L_2 , and so on.

Note that in this type of queueing system there is not a job in service, we only count the jobs in the system at the end of a period. Thus, the number of jobs in the system and in queue coincide in this model.

1.2.1. [1.2.3] Show that the scheme

$$\begin{aligned} L_k &= [L_{k-1} + a_k - c_k]^+, \\ d_k &= L_{k-1} + a_k - L_k. \end{aligned} \quad (1.2.2)$$

is equivalent to a modification of (1.2.1) in which we assume that jobs can be served in the period they arrive.

Of course we are not going to carry out these computations by hand. Typically we use company data to model the arrival process $\{a_k\}_{k=1,2,\dots}$ and the capacity $\{c_k\}_{k=1,2,\dots}$ and feed this data into a computer to carry out the recursions (1.2.1). If we do not have sufficient data, we make a probability model for these data and use the computer to generate random numbers with, hopefully, similar characteristics as the real data. At any rate, from this point on, we assume that it is easy, by means of computers, to obtain numbers a_1, \dots, a_n for $n \gg 1000$, and so on.

Case analysis

Here we continue with the case of the five psychiatrists and analyze the proposed rules to improve the performance of the system. We mainly want to reduce the long waiting times.

As a first step in the analysis, we model the arrival process of patients as a Poisson process, cf. Section 1.1. The duration of a period is taken to be a week. The average number of arrivals per period, based on data of the company, was slightly less than 12 per week; in the simulation we set it to $\lambda = 11.8$ per week. We model the capacity in the form of a matrix such that row i corresponds to the weekly capacity of psychiatrist i :

$$C = \begin{pmatrix} 1 & 1 & 1 & \dots \\ 1 & 1 & 1 & \dots \\ 1 & 1 & 1 & \dots \\ 3 & 3 & 3 & \dots \\ 9 & 9 & 9 & \dots \end{pmatrix}.$$

Thus, psychiatrists 1, 2, and 3 do just one intake per week, the fourth does 3, and the fifth does 9 intakes per week. The sum over column k is the total service capacity for week k of all psychiatrists together.

With the matrix C it is simple to make other capacity schemes. A more balanced scheme would be like this:

$$C = \begin{pmatrix} 2 & 2 & 2 & \dots \\ 2 & 2 & 2 & \dots \\ 3 & 3 & 3 & \dots \\ 4 & 4 & 4 & \dots \\ 4 & 4 & 4 & \dots \end{pmatrix}.$$

Next, we include the effects of holidays on the capacity. This is easily done by setting the capacity of a certain psychiatrist to 0 in a certain week. Let us assume that just one psychiatrist is on leave in a week, each psychiatrist has one week per five weeks off, and the psychiatrists' holiday schemes rotate. To model this, we set $C_{1,1} = C_{2,2} = \dots = C_{1,6} = C_{2,7} = \dots = 0$, i.e.,

$$C = \begin{pmatrix} 0 & 2 & 2 & 2 & 2 & 0 & \dots \\ 2 & 0 & 2 & 2 & 2 & 2 & \dots \\ 3 & 3 & 0 & 3 & 3 & 3 & \dots \\ 4 & 4 & 4 & 0 & 4 & 4 & \dots \\ 4 & 4 & 4 & 4 & 0 & 4 & \dots \end{pmatrix}.$$

Hence, the total average capacity must be $4/5 \cdot (2+2+3+4+4) = 12$ patients per week. The other holiday scheme—all psychiatrists take holiday in the same week—corresponds to setting entire

columns to zero, i.e., $C_{i,5} = C_{i,10} = \dots = 0$ for week 5, 10, and so on. Note that all these variations in holiday schemes result in the same average capacity.

Now that we have modeled the arrivals and the capacities, we can use the recursions (1.2.1) to simulate the queue length process for the four different scenarios proposed by the psychiatrists, unbalanced versus balanced capacity, and spread out holidays versus simultaneous holidays.

The results are shown in Fig. 1. It is apparent that Suggestions 1 and 2 above do not significantly affect the behavior of the queue length process.

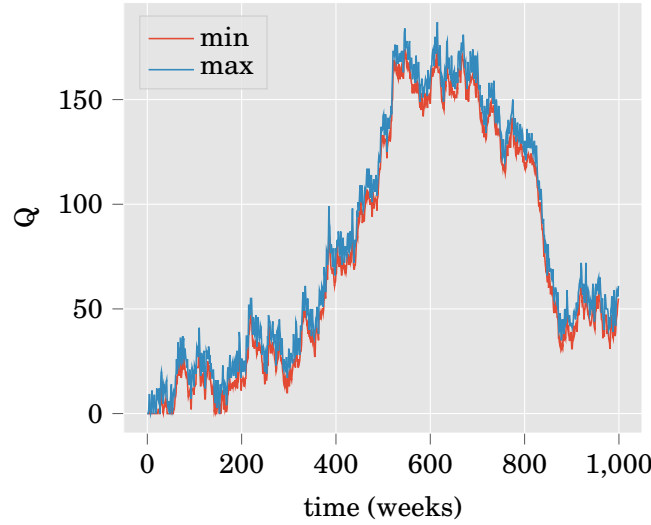


Figure 1: Effect of capacity and holiday plans. We plot for each time point the maximum and the minimum queue length for each of the policies. Apparently, the effect of each of these policies is, for all practical purposes, negligible.

Now we consider Suggestion 3, which comes down to doing more intakes when it is busy, and fewer when it is quiet. A simple rule would be to use week's queue L_{n-1} : if $L_{n-1} < 12$, serve e intakes less (in other words, when the service capacity of week k is larger than the queue length serve e less), when $L_{n-1} > 24$, serve e more. Here, $e = 1$ or 2, or perhaps a larger number. The larger e , the larger the control we exercise.

Let's consider three different control levels, $e = 1$, $e = 2$, and $e = 5$; thus in the last case all psychiatrists do five extra intakes. The previous simulation shows that it is safe to disregard the holiday plans, so just assume a flat service capacity of 12 intakes a week. The results, see Fig. 2, show a striking difference indeed. The queue does not explode anymore, and already taking $e = 1$ has a large influence.

From this simulation experiment, we learn that changing holiday plans or spreading the work over multiple servers, i.e., psychiatrists, does not significantly affect the queueing behavior. However, controlling the service rate as a function of the queue length improves the situation quite dramatically.

Conclusion

From the above case we learn that even with these (deceitfully) simple recursions, we can obtain considerable insight into the behavior and performance of this quite complicated controlled

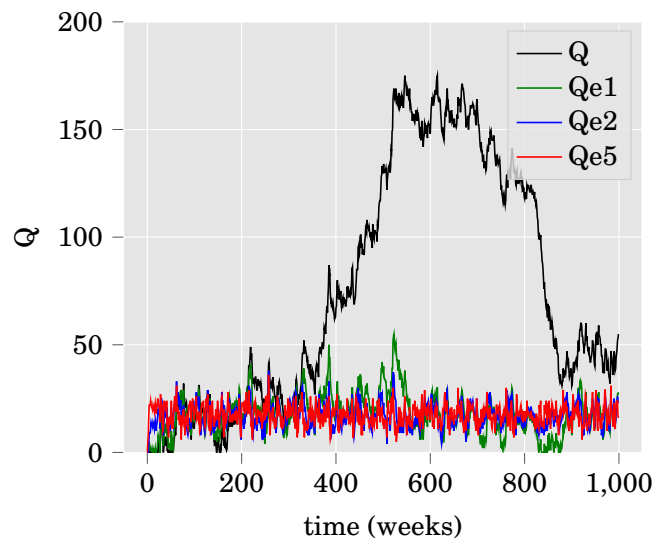


Figure 2: Controlling the number of intakes. Clearly, adapting the service rate ‘does wonders’ to control the queue length.

queueing system¹. In general, the study of queueing systems is focused on studying the probabilistic properties of the queueing length process and related concepts such as waiting time, server occupancy, fraction of customers lost, and so on. Once we have constructed the queueing process, we can compute all relevant performance measures, such as the average waiting time. If it turns out that the performance of the system is not according to what we desire, we can change parts of the system with the aim to improve the situation and assess the effect of this change. For instance, if the average waiting time is too long, we might want to add service capacity. With simulation it is easy to study, hence evaluate, the effect of such decisions.

Exercises

The reader should understand from the above case that, once we have the recursions, we can analyze the system and make, for instance, plots to evaluate suggestions for improvement. Thus, formulating the recursions is crucial to construct queueing processes. For this reason, many of the exercises below concentrate on obtaining recursions for specific queueing systems.

It may be that the recursions you find are not identical to the recursions in the solution. The reason is that the assumptions you make might not be equal to the ones I make. I don’t quite know how to get out of this paradoxical situation. In a sense, to completely specify the model, we need the recursions. However, if the problem statement would contain the recursions, there would be nothing left to practice anymore. Another way is to make the problem description five times as long, but this is also undesirable. So, let’s be pragmatic: the aim is that you practice with modeling, and that you learn from the solutions. If you obtain *reasonable* recursions, but they are different from mine, then your answer is just as good.

1.2.2 (Queue with Blocking). [1.2.4] *A queueing system under daily review, i.e., at the end of the day the queue length is measured. We assume that at the end of the day no jobs are still in service.*

¹ If you doubt the value of simulation, try to develop a mathematical method to analyze multi-server queueing systems with vacations, of which this case is an example.

We assume that jobs that arrive at day k cannot be served in day k . The queue length cannot exceed level K . Formulate a set of recursions to cover this case. What is the loss per period? What is the fraction of jobs lost?

1.2.3 (Estimating the lead time distribution). **[1.2.5]** Take $d_k = \min\{L_{k-1} + a_k, c_k\}$, and assume that jobs are served in FIFO sequence. Find an expression for the shortest possible waiting time $W_{-,k}$ of a job that arrives at time k , and an expression for the largest possible waiting time $W_{+,k}$.

1.2.4 (Cost models). **[1.2.8]** A single-server queueing station processes customers. At the start of a period the server capacity is chosen, so that for period k the capacity is c_k . Demand that arrives in a period can be served in that period. It costs β per unit time per unit processing capacity to operate the machine, i.e., to have it switched on. There is also a cost h per unit time per job in the system. Make a cost model to analyze the long-run average cost for this case.

1.2.5 (N-policies). **[1.2.9]** A machine can switch on and off. If the queue length hits N , the machine switches on, and if the system becomes empty, the machine switches off. It costs K to switch on the machine. There is also a cost β per unit time while the machine is switched on, and it costs h per unit time per customer in the system. Make a cost model.

1.2.6 (Priority queueing). **[1.2.12]** An interesting situation is a system with two queues served by one server, but such that one queue, queue A, gets priority over the other queue. Again find a set of recursions to describe this case.

1.2.7 (Queue with protected service capacity and lost capacity). **[1.2.14]** Consider a single-server that serves two parallel queues A and B. Each queue receives a minimal service capacity every period. Reserved capacity unused for one queue cannot be used to serve the other queue. Any extra capacity beyond the reserved capacity is given to queue A with priority. Formulate a set of recursions to analyze this situation.

Let r_A be the reserved capacity for queue A, and likewise for r_B . We assume of course that $c_k \geq r_A + r_B$, for all k .

1.2.8 (Tandem networks). **[1.2.15]** Consider a production network with two production stations in tandem, that is, the jobs processed by station A are in the next period to the downstream Station B. Extend the recursions of (1.2.1) to simulate this situation.

1.2.9 (Merging departure streams). **[1.2.17]** Consider another production situation with two machines, A and B say, that send their products to Station C. Derive a set of recursion relations to simulate this system.

1.3 EXPONENTIAL DISTRIBUTION

In Section 1.1 we introduced the Poisson process as a natural model of the (random) number of jobs arriving during intervals of time. As we will see in the sections to come, we can model a single-server queueing system in continuous time if we specify the (probability) distribution of the inter-arrival times, i.e., the time between consecutive arrival epochs of jobs. A particularly fruitful model for the distribution of the inter-arrival times is the exponential distribution because, as it turns out, it is intimately related to the Poisson distribution. Besides explaining this relation, we derive many useful properties of the exponential distribution, in particular, that it is *memoryless*.

We say that X is an *exponentially distributed* random variable with mean $1/\lambda$ if

$$P(X \leq t) = 1 - e^{-\lambda t},$$

and then we write $X \sim \text{Exp}(\lambda)$.

The Poisson process N and exponentially distributed inter-arrival times are intimately related: A counting process $\{N(t)\}$ is a *Poisson process* with rate λ if and only if the inter-arrival times $\{X_i\}$ are i.i.d. with common random variable $X \sim \text{Exp}(\lambda)$, in short, $X_i \sim \text{Exp}(\lambda) \Leftrightarrow N(t) \sim P(\lambda t)$. We next provide further relations between the Poisson distribution and the exponential distribution.

1.3.1. [1.3.2] *If the random variable $X \sim \text{Exp}(\lambda)$, show that its mean $E[X] = \frac{1}{\lambda}$.*

1.3.2. [1.3.3] *If $X \sim \text{Exp}(\lambda)$, show that its second moment $E[X^2] = \frac{2}{\lambda^2}$.*

1.3.3. [1.3.7] *Use the moment-generating function of $X \sim \text{Exp}(\lambda)$ to show that*

$$E[X] = \frac{1}{\lambda}, \quad E[X^2] = \frac{2}{\lambda^2}.$$

We now provide a number of relations between the Poisson distribution and the exponential distribution to conclude that a process N_λ is a Poisson process with rate λ iff the inter-arrival times $\{X_i\}$ between individual jobs are i.i.d. $\sim \text{Exp}(\lambda)$.

1.3.4. [1.3.9] *If N_λ is a Poisson process with rate λ , show that the time X_1 to the first arriving job is $\text{Exp}(\lambda)$.*

1.3.5. [1.3.11] *Let A_i be the arrival time of customer i and set $A_0 = 0$. Assume that the inter-arrival times $\{X_i\}$ are i.i.d. with exponential distribution with mean $1/\lambda$ for some $\lambda > 0$. Prove that A_i has density*

$$f_{A_i}(t) = \lambda e^{-\lambda t} \frac{(\lambda t)^{i-1}}{(i-1)!}.$$

1.3.6. [1.3.13] *If the inter-arrival times $\{X_i\}$ are i.i.d. $\sim \text{Exp}(\lambda)$, prove that the number $N(t)$ of arrivals during the interval $[0, t]$ is Poisson distributed.*

We now introduce another fundamental concept. A random variable X is called *memoryless* when it satisfies

$$P(X > t + h | X > t) = P(X > h).$$

In words, the probability that X is larger than some time $t + h$, conditional on it being larger than a time t , is equal to the probability that X is larger than h .

1.3.7. [1.3.14] *Show that $X \sim \text{Exp}(\lambda)$ is memoryless.*

In fact, it can be shown that only exponential random variables have the memoryless property. The proof of this fact requires quite some work; we refer the reader to the literature if s/he wants to check this, see e.g. [Yushkevich and Dynkin \[1969, Appendix 3\]](#).

1.3.8. [1.3.15] *If $X \sim \text{Exp}(\lambda)$ and $S \sim \text{Exp}(\mu)$, and X and S are independent, show that*

$$Z = \min\{X, S\} \sim \text{Exp}(\lambda + \mu),$$

hence $E[Z] = (\lambda + \mu)^{-1}$.

1.3.9. [1.3.16] *If $X \sim \text{Exp}(\lambda)$, $S \sim \text{Exp}(\mu)$, and X and S are independent, show that*

$$P(X \leq S) = \frac{\lambda}{\lambda + \mu}.$$

1.4 CONSTRUCTION OF THE SINGLE-SERVER QUEUEING PROCESS IN CONTINUOUS TIME

In Section 1.2 we modeled time as progressing in discrete ‘chunks’: minutes, hours, days, and so on. For given numbers of arrivals and capacity per period, we use the recursions (1.2.1) to compute the departures and queue length per period. However, we can also model time in a continuous way, so that jobs can arrive at any moment and have arbitrary service times. In this section, we consider a single-server FIFO queueing process in continuous time.

Assume we are given the *arrival process* $\{A(t); t \geq 0\}$, i.e., the number of jobs that arrived during $[0, t]$. Thus, $\{A(t); t \geq 0\}$ is a *counting process*.

From this arrival process, we can derive various other interesting concepts, such as the *arrival times* of individual jobs. Especially, if we know that $A(s) = k - 1$ and $A(t) = k$, then the arrival time A_k of the k th job must lie somewhere in $(s, t]$. Thus, from $\{A(t)\}$, we can define²

$$A_k = \min\{t : A(t) \geq k\}, \quad (1.4.1)$$

and set $A_0 = 0$. Once we have the set of arrival times $\{A_k\}$, the set of *inter-arrival times* $\{X_k, k = 1, 2, \dots\}$ between consecutive customers can be constructed as

$$X_k = A_k - A_{k-1}. \quad (1.4.2)$$

However, often the basic data consists of the inter-arrival times $\{X_k; k = 1, 2, \dots\}$ rather than the arrival times $\{A_k\}$ or the arrival process $\{A(t)\}$. Then we construct the arrival times as

$$A_k = A_{k-1} + X_k,$$

with $A_0 = 0$. From the arrival times $\{A_k\}$ we can, in turn, construct the arrival process $\{A(t)\}$ as

$$A(t) = \max\{k : A_k \leq t\}. \quad (1.4.3)$$

Thus, from the inter-arrival times $\{X_k\}$ it is possible to construct $\{A_k\}$ and $\{A(t)\}$, and from $\{A(t)\}$ we can obtain $\{A_k\}$ and $\{X_k\}$.

1.4.1. Show that we can also define $A(t)$ as

$$A(t) = \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq t}. \quad (1.4.4)$$

Thus, we just count all arrivals that occur up to and including time t .

The *service times* of the jobs are given by the sequence $\{S_k\}$. Typically we assume that the service times are i.i.d., and also independent of the arrival process $\{A(t)\}$. With the arrival times and service times we can construct the *waiting time in queue* $\{W_{Q,k}\}$ as seen by the jobs at the moment they arrive. From Fig. 3 it is evident that

$$W_{Q,k} = [W_{Q,k-1} + S_{k-1} - X_k]^+. \quad (1.4.5)$$

With this it is easy to compute the waiting times: from the initial condition $W_{Q,0} = 0$ we can obtain $W_{Q,1}$, and then $W_{Q,2}$, and so on.

² If we want to be mathematically precise, we must take \inf rather than \min . However, in this book we do not want to distinguish between subtleties.



Figure 3: Construction of the single-server queue in continuous time. The sojourn time W_k of the k th job is the sum of the work in queue $W_{Q,k}$ at its arrival epoch A_k and its service time S_k ; its departure time is then $D_k = A_k + W_k$. The waiting time of job k is clearly equal to $W_{k-1} - X_k$. We also see that job $k+1$ arrives at an empty system, hence its sojourn time $W_{k+1} = S_{k+1}$. Finally, the virtual waiting time process is shown by the lines with slope -1 .

1.4.2. [1.4.6] If $S \sim U[0, 7]$ and $X \sim U[0, 10]$, where $U[I]$ stands for the uniform distribution concentrated on the interval I , compute $P(S - X \leq u)$, for S and X independent.

The time job k leaves the queue and moves to the server is given by

$$\tilde{A}_k = A_k + W_{Q,k},$$

because a job can only move to the server after its arrival plus the time it needs to wait in queue. Note that we here explicitly use the FIFO assumption. Right after the job moves from the queue to the server, its service starts. Thus, \tilde{A}_k is also the epoch at which the service of job k starts.

After completing its service, the job leaves the system. Hence, the *departure time of the system* is

$$D_k = \tilde{A}_k + S_k.$$

This in turn specifies the departure process $\{D(t)\}$ as

$$D(t) = \max\{k : D_k \leq t\} = \sum_{k=1}^{\infty} \mathbb{1}_{D_k \leq t}.$$

The *sojourn time*, or *waiting time in the system*, W_k , is the time a job spends in the entire system. With the above relations we see that

$$W_k = D_k - A_k = \tilde{A}_k + S_k - A_k = W_{Q,k} + S_k, \quad (1.4.6)$$

where each of these equations has its own interpretation.

1.4.3. [1.4.7] Explain that the following recursion for W_k is valid:

$$W_{Q,k} = [W_{k-1} - X_k]^+, \quad W_k = W_{Q,k} + S_k = [W_{k-1} - X_k]^+ + S_k, \quad (1.4.7)$$

from which follows a recursion for D_k :

$$D_k = A_k + W_k. \quad (1.4.8)$$

1.4.4. [1.4.11] Explain the following recursions for a single server queue:

$$\begin{aligned} A_k &= A_{k-1} + X_k, \\ D_k &= \max\{A_k, D_{k-1}\} + S_k, \\ W_k &= D_k - A_k. \end{aligned} \tag{1.4.9}$$

The *virtual waiting time process* $\{V(t)\}$ is the amount of waiting that an arrival would see if it would arrive at time t . To construct $\{V(t)\}$, we simply draw lines that start at points (A_k, W_k) and have slope -1, unless the line hits the x -axis, in which case the virtual waiting time remains zero until the next arrival occurs.

1.4.5. [1.4.12] Provide a specification of the virtual waiting time process $\{V(t)\}$ for all t .

Once we have the arrival and departure processes it is easy to compute the *number of jobs in the system* at time t as, cf. Fig. 4,

$$L(t) = A(t) - D(t) + L(0), \tag{1.4.10}$$

where $L(0)$ is the number of jobs in the system at time $t = 0$; typically we assume that $L(0) = 0$. As in a queueing system, jobs can be in queue or in service, we distinguish between the number in the system $L(t)$, the number in queue $L_Q(t)$, and the number of jobs in service $L_s(t)$.

In summary, starting from a sequence of inter-arrival times $\{X_k\}$ and service times $\{S_k\}$ we can obtain a set of recursions by which we can simulate a queueing process in continuous time. A bit of experimentation with Python (or R perhaps) will reveal that this is easy.

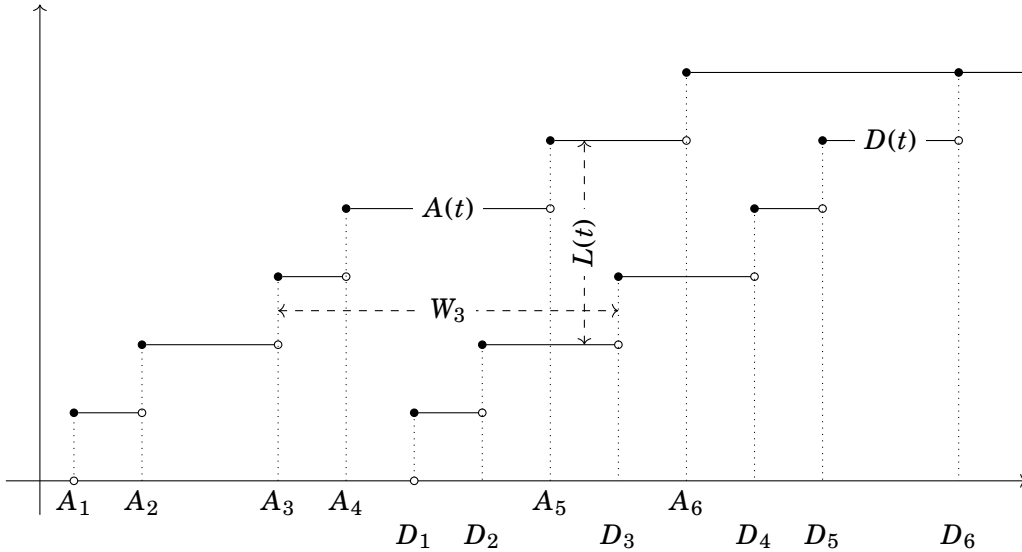


Figure 4: Relation between the arrival process $\{A(t)\}$, the departure process $\{D(t)\}$, the number in the system $\{L(t)\}$ and the waiting times $\{W_k\}$. In particular, $L(t)$ is the difference between the graphs of $A(t)$ and $D(t)$.

1.4.6. [1.4.18] Show that $L(t) = \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq t < D_k}$ when the system starts empty.

Define the number of jobs in the system as seen by the k th arrival as

$$L(A_k-). \tag{1.4.11}$$

Observe that we write A_k- , and not A_k . As we are concerned with a process with jumps, we need to be quite particular about left and right limits at jump epochs.

1.4.7. [1.4.20] *With the recursions (1.4.7) it is apparently easy to compute the waiting time (in queue), but it is less simple to compute the number of jobs in queue or in the system. In this exercise we develop an algorithm to compute the number of jobs in the system as seen by arrivals. Explain why the following (algorithmic efficient) procedure works:*

$$L(A_k-) = L(A_{k-1}-) + 1 - \sum_{i=k-1-L(A_{k-1}-)}^{k-1} \mathbb{1}_{D_i < A_k}.$$

FROM TRANSIENT TO STEADY-STATE ANALYSIS

With the tools developed in Chapter 1 we can simulate queueing processes. The next step is to develop mathematical models of queueing systems. The aim of this chapter is to make start with this. However, as we will see in Section 2.2, the mathematical analysis of the time-dependent behavior of queueing systems is beyond our capabilities; the transient behavior of even the simplest queueing system is already extremely complicated. Thus, we have to lower our goals, and for this reason we will focus the steady-state behavior of queueing systems. Intuitively speaking, this requires the system to be stable, for otherwise the queue length process grows to infinity.

We introduce concepts of stability and load in Section 2.3 and express these in terms of the arrival, service and departure rates. The notions of arrival and service rate are crucial because they capture our intuition that when jobs arrive faster on average than they can leave, the queue must ‘explode’. As we will see, when the arrival rate is smaller than the service rate, the system is stable. Once stability is ensured, we can properly define in Section 2.4 a number of measures to characterize the performance of the queueing system, such as the average waiting time. In Chapter 3 we will see how to use the load and expected service time to approximate expected waiting times in queue.

Before introducing these definitions, however, we need to introduce some notational shortcuts to characterize the type of queueing process, which is the topic of Section 2.1. We provide in Section 2.5 an overview of the relations we introduce in this chapter.

2.1 KENDALL’S NOTATION

As became apparent in Sections 1.2 and 1.4, the construction of any single-station queueing process involves three main elements: the distribution of the inter-arrival times between consecutive jobs, the distribution of the service times of the individual jobs, and the number of servers present to process jobs. In this characterization, it is implicit that the inter-arrival times form a set of i.i.d. (independent and identically distributed) random variables, the service times are also i.i.d., and finally, the inter-arrival times and service times are mutually independent.

To characterize the type of queueing process it is common to use *Kendall’s abbreviation* $A/B/c/K$, where A is the distribution of the inter-arrival times, B the distribution of the service times, c the number of servers, and K the system size, i.e., the total number of customers that can be simultaneously present, whether in queue or in service.¹ In this notation it is assumed that jobs are served in first-in-first-out (FIFO) order; FIFO scheduling is also often called first-come-first-served (FCFS).

When at an arrival a number of jobs arrive simultaneously (like a bus at a restaurant), we say that a batch arrives. Likewise, the server can work in batches, for instance, when an oven processes multiple jobs at the same time. We write $A^X/B^Y/c$ to denote that X is the distribution of the arrival batch size and Y is the distribution of the service batch sizes. When $X \equiv Y \equiv 1$,

¹ The meaning of K differs among authors. Sometimes it stands for the capacity of the queue, not the entire system. In this book K corresponds to the system’s size.

i.e., single batch arrivals and single batch services, we suppress the X and Y in the queueing formula.

Two inter-arrival and service distributions are the most important in queueing theory: the exponential distribution denoted with the shorthand M , as it is memoryless, and a general distribution (with the implicit assumption that its first moment is finite) denoted with G . We write D for a deterministic (constant) random variable.

Familiarize yourself with this notation as it is used continuously in the rest of the book. Here are some exercises to illustrate the notation.

2.1.1. [2.1.7] *What is the meaning of $M(n)/M(n)/1$?*

2.1.2. [2.1.8] *What is the meaning of $M^X/M/1$?*

2.1.3. [2.1.13] *Is the $M/D/1$ queue a specific type of $M/G/c$ queue?*

You should also understand the differences between different scheduling rules. The next exercise should help with this.

2.1.4. [2.1.14] *What are some advantages and disadvantages of using the Shortest Processing Time First (SPTF) rule to serve jobs?*

When a customer finds a large queue in front of it, s/he can use the normal distribution to estimate the distribution of the time s/he will spend in queue. The next exercise shows how.

2.1.5. [2.1.15] *Suppose for the $G/G/1$ that a job sees n jobs in the system upon arrival. Use the central limit theorem to estimate the distribution of the waiting time in queue for this job.*

2.2 QUEUEING PROCESSES AS REGULATED RANDOM WALKS

In this section, we provide an elegant construction of a queueing process $\{L_k\}$ based on a *random walk* $\{Z_k\}$. This serves two aims. The first is to show that a characterization of the transient behavior of a queueing process is typically extremely hard. Thus, in the rest of the book we will only study queueing systems in steady-state. The second is to show that queueing theory is essentially based on concepts of fundamental interest in probability theory (the random walk), hence is strongly related to many other applications of random walks, such as in finance, inventory theory, and insurance theory. Moreover, we can use tools of random walk to analyze queueing systems, such as a characterization of the distribution of the time until an especially large queue length is reached; we refer the reader to literature for this.

In the construction of queueing processes as set out in Section 1.2 we are given two sequences of i.i.d. random variables: the number of arrivals $\{a_k\}$ per period and the service capacities $\{c_k\}$, cf., (1.2.2). Observe that in (1.2.2) the process $\{L_k\}$ shares a resemblance to a random walk $\{Z_k, k = 0, 1, \dots\}$ with Z_k given by

$$Z_k = Z_{k-1} + a_k - c_k. \quad (2.2.1)$$

To see that $\{Z_k\}$ is indeed a random walk, observe that Z makes jumps of size $a_k - c_k, k = 1, \dots$, and $\{a_k - c_k\}$ is a sequence of i.i.d. random variables since, by assumption, $\{a_k\}$ and $\{c_k\}$ are i.i.d. Clearly, $\{Z_k\}$ is ‘free’, i.e., it can take positive and negative values, but $\{L_k\}$ is restricted to the non-negative integers.

2.2.1. [2.2.1] Show that L_k satisfies the relation

$$L_k = Z_k - \min_{1 \leq i \leq k} Z_i \wedge 0, \quad (2.2.2)$$

where Z_k is defined by the above random walk and we write $a \wedge b$ for $\min\{a, b\}$.

This recursion for L_k leads to really interesting graphs. In Fig. 5 we take $a_k \sim B(0.3)$, i.e., a_k is Bernoulli distributed with success parameter $p = 0.3$, i.e., $P(a_k = 1) = 0.3 = 1 - P(a_k = 0)$, and $c_k \sim B(0.4)$. In Fig. 6, $a_k \sim B(0.49)$ and the random walk is constructed as

$$Z_k = Z_{k-1} + 2a_k - 1. \quad (2.2.3)$$

Thus, if $a_k = 1$, the random walk increases by one step, while if $a_k = 0$, the random walk decreases by one step, so that $Z_k \neq Z_{k-1}$ always. Observe that this is slightly different from a random walk that satisfies (2.2.1); there, $Z_k = Z_{k-1}$, if $a_k = c_k$.

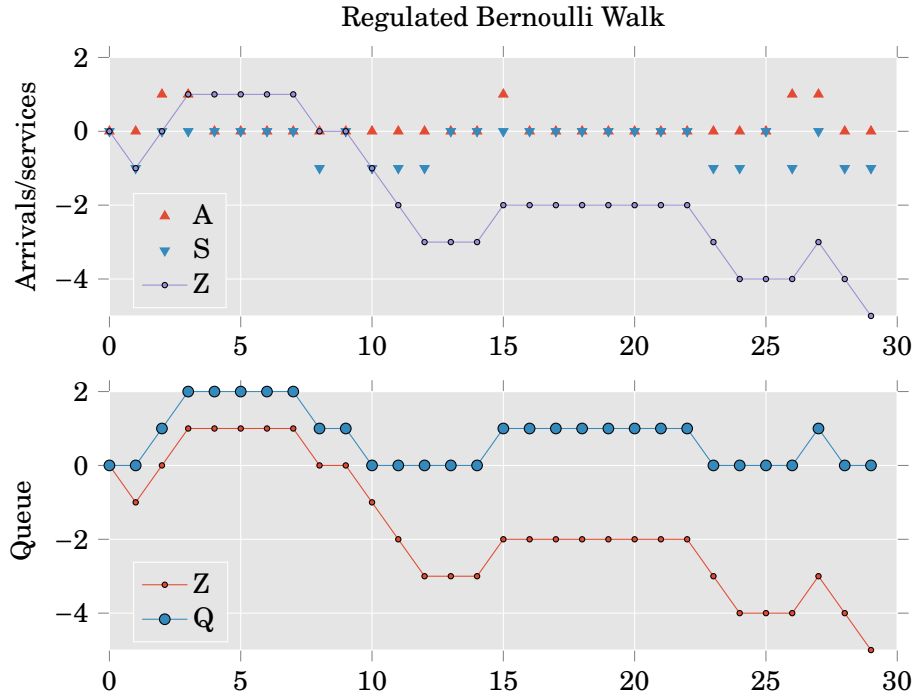


Figure 5: The upper panel shows a graph of the random walk Z . An upward pointing triangle corresponds to an arrival, a downward triangle to a potential service. The lower panel shows the queueing process $\{L_k\}$ as a random walk with reflection.

With (2.2.2), we see that a random walk $\{Z_k\}$ can be converted into a queueing process $\{L_k\}$, and we might try to understand the transient behavior of the latter by investigating the transient behavior of the former. Suppose that $a_k \sim P(\lambda)$ and $c_k \sim P(\mu)$.

2.2.2. [2.2.2] Show that if $\{a_k\}$ forms an i.i.d. sequence of random variables all Poisson distributed $P(\lambda)$ then, $\sum_{j=1}^k a_j = P(\lambda k)$.

With the above exercise,

$$Z_k = Z_0 + N_{\lambda k} - N_{\mu k},$$

and we call $Z = \{Z_k\}$ the *free* (discrete-time) $M/M/1$ queue as, contrary to the real $M/M/1$ queue, Z can take negative values.



Figure 6: Another example of a reflected random walk.

2.2.3. [2.2.3] Show that when $n > m$ and $Z_0 = m$,

$$P(Z_k = n) = e^{-(\lambda+\mu)k} (\lambda k)^{n-m} \sum_{j=0}^{\infty} \frac{(\lambda \mu k^2)^j}{j!(n-m+j)!}.$$

The solution of the above exercise shows that there is no simple function by which we can compute the transient distribution of this simple random walk Z . Since a queueing process is typically a more complicated object (as we need to obtain L from Z via (2.2.2)), our hopes of finding anything simple for the transient analysis of the $M/M/1$ queue should not be too high. But the $M/M/1$ queue is about the simplest queueing system; other queueing systems are (much) more complicated. We therefore give up the analysis of the transient behavior of queueing systems and henceforth contend ourselves with the analysis of queueing systems in the limit as $t \rightarrow \infty$.

2.3 RATE, STABILITY AND LOAD

In this section, we develop a number of essential concepts to analyze queueing systems: the arrival, service and departure rate. With these we define the load, which is, arguably, the most important performance measure of a queueing system to check.

We first formalize the arrival rate and departure rate in terms of the arrival and departure processes $\{A(t)\}$ and $\{D(t)\}$; recall that these are *counting processes*. The *arrival rate* is the long-run average number of jobs that arrive per unit time along a sample path, i.e.,

$$\lambda = \lim_{t \rightarrow \infty} \frac{A(t)}{t}. \quad (2.3.1)$$

We remark in passing that this limit does not necessarily exist if $A(t)$ is some pathological function. If, however, the inter-arrival times $\{X_k\}$ are the basic data, and $\{X_k\}$ are i.i.d. and distributed as a generic random variable X with finite mean $E[X]$, we can construct $\{A_k\}$ and $\{A(t)\}$ as described in Section 1.4, and then the strong law of large numbers guarantees that the above limit exists.²

Likewise, define the *departure rate* as

$$\delta = \lim_{t \rightarrow \infty} \frac{D(t)}{t}. \quad (2.3.2)$$

² In fact, $A(t)/t \rightarrow \lambda$ with probability one.

Observe that, if the system is empty at time 0, the number of departures must be smaller than or equal to the number of arrivals, i.e., $D(t) \leq A(t)$ for all t . Therefore,

$$\delta = \lim_{t \rightarrow \infty} \frac{D(t)}{t} \leq \lim_{t \rightarrow \infty} \frac{A(t)}{t} = \lambda. \quad (2.3.3)$$

It is evident that when $\lambda > \delta$, the system length process $L(t) \rightarrow \infty$ as $t \rightarrow \infty$. We therefore call a system *rate-stable* if

$$\lambda = \delta.$$

In words: the system is rate-stable whenever jobs leave the system just as fast as they arrive in the long run.

2.3.1. [2.3.2] *If the system starts empty, then we know that the number $L(t)$ in the system at time t is equal to $A(t) - D(t)$. Show that the system is rate-stable if $L(t)$ remains finite, or, more generally, $L(t)/t \rightarrow 0$ as $t \rightarrow \infty$.*

We next relate the arrival rate λ to the expected inter-arrival time $E[X]$. Observe that at time $t = A_n$ precisely n arrivals occurred. Thus, we see that $A(A_n) = n$, and therefore

$$\frac{1}{n} \sum_{k=1}^n X_k = \frac{A_n}{n} = \frac{A_n}{A(A_n)}.$$

But since $A_n \rightarrow \infty$ if $n \rightarrow \infty$, it follows from (2.3.1) that the average inter-arrival time between two consecutive jobs is

$$E[X] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n X_k = \lim_{n \rightarrow \infty} \frac{A_n}{A(A_n)} = \lim_{t \rightarrow \infty} \frac{t}{A(t)} = \frac{1}{\lambda}, \quad (2.3.4)$$

where we take $t = A_n$ in the limit for $t \rightarrow \infty$. In words, the arrival rate λ is the *inverse* of the expected inter-arrival time $E[X]$.

Assume now that there is a single server. Let S_k be the required service time of the k th job to be served, so that $U_n = \sum_{k=1}^n S_k$ becomes the total service time required by the first n jobs. With this, let $U(t) = \max\{n : U_n \leq t\}$, so that we can define the *service rate* or *processing rate* as

$$\mu = \lim_{t \rightarrow \infty} \frac{U(t)}{t}.$$

Similar to the relation $E[X] = 1/\lambda$, we have the relation

$$E[S] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n S_k = \lim_{n \rightarrow \infty} \frac{U_n}{n} = \lim_{n \rightarrow \infty} \frac{U_n}{U(U_n)} = \lim_{t \rightarrow \infty} \frac{t}{U(t)} = \frac{1}{\mu}.$$

Thus, the service rate μ is the *inverse* of the expected service time $E[S]$.

It turns out that, when $\mu = \lambda$ and $V[S] > 0$ or $V[X] > 0$, the queue length process behaves in a very peculiar way. This follows from the fact that the random walk without drift, i.e., $\mu = \lambda$, has some unexpected behavior³, and we know from Section 2.2 that queueing systems and random walks are intimately related. To avoid such problems, we henceforth (and implicitly) require that $\mu > \lambda$. Observe also the evident fact that jobs cannot depart faster than they can be served, hence, $D(t) \leq U(t)$ for all t , hence $\delta \leq \mu$. Combining this with $\lambda < \mu$ and rate-stability, we see that $\delta < \mu$.

³ See any good book on probability theory.

2.3.2. [2.3.4] Define $\tilde{X}_k = S_{k-1} - X_k$. Show that $E[\tilde{X}_k] < 0$ implies that $\lambda < \mu$.

Perhaps the most important performance measure is the concept of *load* ρ , which is defined as the rate at which jobs arrive times the average amount of work per job: $\rho = \lambda E[S]$. From the identities $\lambda^{-1} = E[X]$ and $\mu^{-1} = E[S]$, we get some further relations:

$$\rho = \lambda E[S] = \frac{\lambda}{\mu} = \frac{E[S]}{E[X]}. \quad (2.3.5)$$

Observe that by assumption $\mu > \lambda$, hence $\rho < 1$. The relation $\rho = E[S]/E[X] < 1$ then tells us that the average time it takes to serve a job must be less than the average time between two consecutive arrivals, i.e., $E[S] < E[X]$. Also, when $\lambda < \mu$, it is easy to check with simulation that when $L(0)$ is very large, $L(t) \approx L(0) - (\mu - \lambda)t$ until the system is empty, while if $\lambda > \mu$, we have that $L(t) \approx L(0) + (\lambda - \mu)t$.

2.3.3. [2.3.5] Consider a queueing system with c servers with identical production rates μ . What would be a reasonable stability criterion for this system?

2.4 (LIMITS OF) EMPIRICAL PERFORMANCE MEASURES

If the queueing system is rate-stable, we can sensibly define a number of long-run average performance measures such as the load, the average waiting time in queue, and so on. This we do here and refer to Fig. 8 for an overview of the relations between these performance measures.

Recall that we constructed a single-server queueing process $\{L(t)\}$ in cf. Section 1.4 based on the arrival process $\{A(t)\}$ and service process $\{S(t)\}$. Once we have the queueing process, we can compute the waiting time process $\{W_{Q,k}\}$, the sojourn time process $\{W_k\}$.

Define the *expected sojourn time as seen by arrivals* as

$$E[W] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n W_k. \quad (2.4.1)$$

Note that this performance measure is the limit of an *empirical* performance measure as observed by arriving jobs: the first job experiences a sojourn time W_1 when it arrives, the second a sojourn time W_2 , and so on. For this reason, we colloquially say such a performance measure is as ‘seen by arrivals’. Similarly, the *expected waiting time in queue* is defined as

$$E[W_Q] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n W_{Q,k}. \quad (2.4.2)$$

The *distribution of the sojourn times* seen by arrivals can be found by counting:

$$P(W \leq x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{W_k \leq x}. \quad (2.4.3)$$

The (sample) *average number of jobs* in the system as seen by arrivals is given by

$$E[L] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n L(A_k -), \quad (2.4.4)$$

since $L(A_k -)$ is the number of jobs in the system at the arrival epoch of the k th job. Finally, the *distribution of $\{L(t)\}$* as seen by arrivals, is given by

$$P(L \leq m) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{L(A_k -) \leq m}. \quad (2.4.5)$$

A related set of performance measures follows by tracking the system's behavior over time and taking the *time-average*, rather than the average at arrival epochs. Assuming the limit exists we use (1.4.10) to define the *time-average number of jobs* as

$$E[L] = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t L(s) ds. \quad (2.4.6)$$

Observe that, notwithstanding that the symbols are the same, this expectation is not necessarily the same as (2.4.4). In a loose sense we can say that $E[L]$ is the average number in the system as perceived by the *server*. Next, define the *time-average fraction of time the system contains at most m jobs* as

$$P(L \leq m) = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t \mathbb{1}_{L(s) \leq m} ds. \quad (2.4.7)$$

Again, this probability need not be the same as what customers see upon arrival.

2.4.1. [2.4.1] *Design a queueing system to show that the average number of jobs in the system as seen by the server can be very different from what customers see upon arrival.*

Taking the above limits requires considerable care, and at least two questions are important. What type of limit do we actually mean here? And, once such limits are given proper meaning, what is the rate of convergence of, for instance, the random variables $\{L_k\}$ to the limiting random variable L ? Here we sidestep all such fundamental issues, but see Remark 2.4.1. Assuming that the first of these questions is answered, the limiting random variables are known as the *steady-state* or *stationary* limits, and the related distributions are often called *limiting* or *stationary distributions*.

Related to the second question about the convergence rate, we provide some intuition by means of an example. We consider the sequence of waiting times $\{W_{Q,k}\}$ to a limiting random variable W_Q , where $W_{Q,k}$ is constructed according to the recursion (1.4.5). Suppose that $X_k \sim U\{1, 2, 4\}$ and $S_k \sim U\{1, 2, 3\}$. Starting with $W_{Q,0} = 5$ we use (1.4.5) to compute the *exact* distribution of $W_{Q,k}$ for $k = 1, 2, \dots, 20$, cf., the left panel in Fig. 7. We see that when $k = 5$, the ‘hump’ of $P(W_{Q,5} = x)$ around $x = 5$ is due the starting value of $W_{Q,0} = 5$. However, for $k > 10$ the distribution of $W_{Q,k}$ hardly changes, at least not visually. Apparently, the convergence of the sequence of distributions of $W_{Q,k}$ is rather fast.

In the middle panel we show the results of a set of *simulations* for increasing simulation length, up to $N = 1000$ samples, where we use the *empirical distribution*

$$P(W_Q \leq x) = \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{W_{Q,k} \leq x}$$

to estimate the waiting time distribution. As should be clear from the figure, the simulated distribution also converges quite fast to some limiting function. Finally, in the right panel we compare the densities as obtained by the exact method and simulation with $n = 1000$. Clearly, for all practical purposes, these densities can be treated as the same.

The combination of the fast convergence to the steady-state situation and the difficulties with the transient analysis validates, to some extent, that most queueing theory is concerned with the analysis of the system in *stationary* or *steady state*.

2.4.2. [2.4.3] *Suppose that $X_k \in \{1, 3\}$ such that $P(X_k = 1) = P(X_k = 3)$ and $S_k \in \{1, 2\}$ with $P(S_k = 1) = P(S_k = 2)$. Write a computer program to see how fast the distributions of $W_{Q,k}$ converge to a limiting distribution function.*

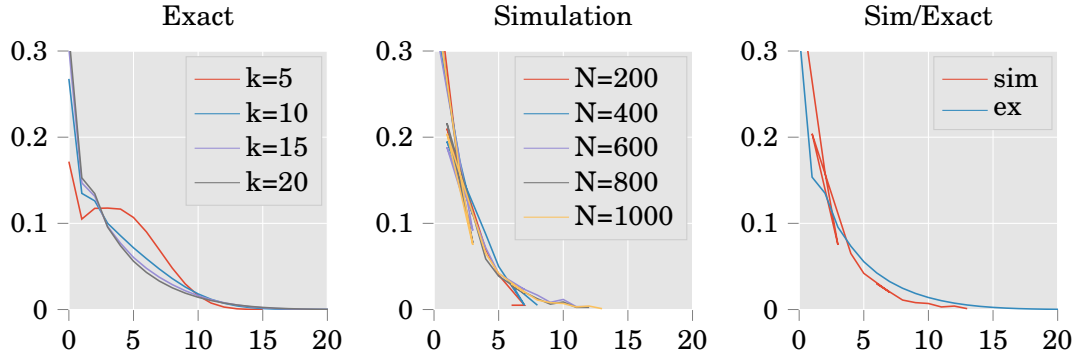


Figure 7: The density of $W_{Q,k}$ for $k = 5, 10, 15, 20$ computed by an exact method as compared the density obtained by simulation of different run lengths $N = 200, 400, \dots, 1000$. The right panel compares the exact density of $W_{Q,20}$ to the density obtained by simulation for $N = 1000$.

Up to now, we considered queueing systems in continuous time. However, the performance measures for discrete-time queueing models needs some modification, mainly for the reason that multiple jobs can arrive in a single period. What statistics do these jobs assemble? The next exercise focuses on this problem.

2.4.3. [2.4.4] Consider the discrete-time model of the queueing system specified by (1.2.1). In such queueing systems, jobs arrive in batches, for instance, when $a_k = 3$, three jobs arrive in slot k . Assuming that $L_{k-1} = 5$, what queue length have these 3 arrivals seen?

Provide a definition similar to (2.4.5) for the case in which we say that all arrivals of a batch see the same number in the system. For this, assume that the batch sees, upon arrival, $L_{k-1} - d_k$ jobs in the system. In other words, we assume that in a period the served jobs first depart before new jobs arrive.

Provide another definition to express that the first of a batch of arrivals sees less jobs in the system than the last arrival of a batch. (To get the details right is harder than you might think; I did it wrong several times.)

Remark 2.4.1. The long-run limiting behavior of a queueing system (i.e., the first question) is a subtle topic. For instance, does there exist a random variable L such that $L_k \rightarrow L$ in some sense? The answer to this question requires a considerable amount of mathematics. To sketch what has to be done, realize that we first need to define $\{L_k\}$ as a stochastic process (up to now we just considered each L_k as a number, i.e., a measurement or simulation of the queue length time of the k th period.) The construction of L_k as a proper random variable is not as simple as the definition of the service times $\{S_k\}$, for instance. In the latter case, we just assume these random variables to be i.i.d. However, in the former case, it is apparent from (1.4.7) that the queue lengths $\{L_k\}$ are constructed in terms of recursions, hence they are certainly not i.i.d. Next, based on these recursions, we need to show that the sequence of distribution functions $\{G_k\}$ associated with the random variables $\{L_k\}$ converges to some limiting distribution function G , say. Finally, it is necessary to show that it is possible to construct a random variable L that has G as its distribution function. In this sense, finally, we say that $L_k \rightarrow L$. All this can be done, but we refer to the specialized literature for this.

2.5 GRAPHICAL SUMMARY

Here is, in graphical form, an overview to show the relation between the concepts developed in this chapter.

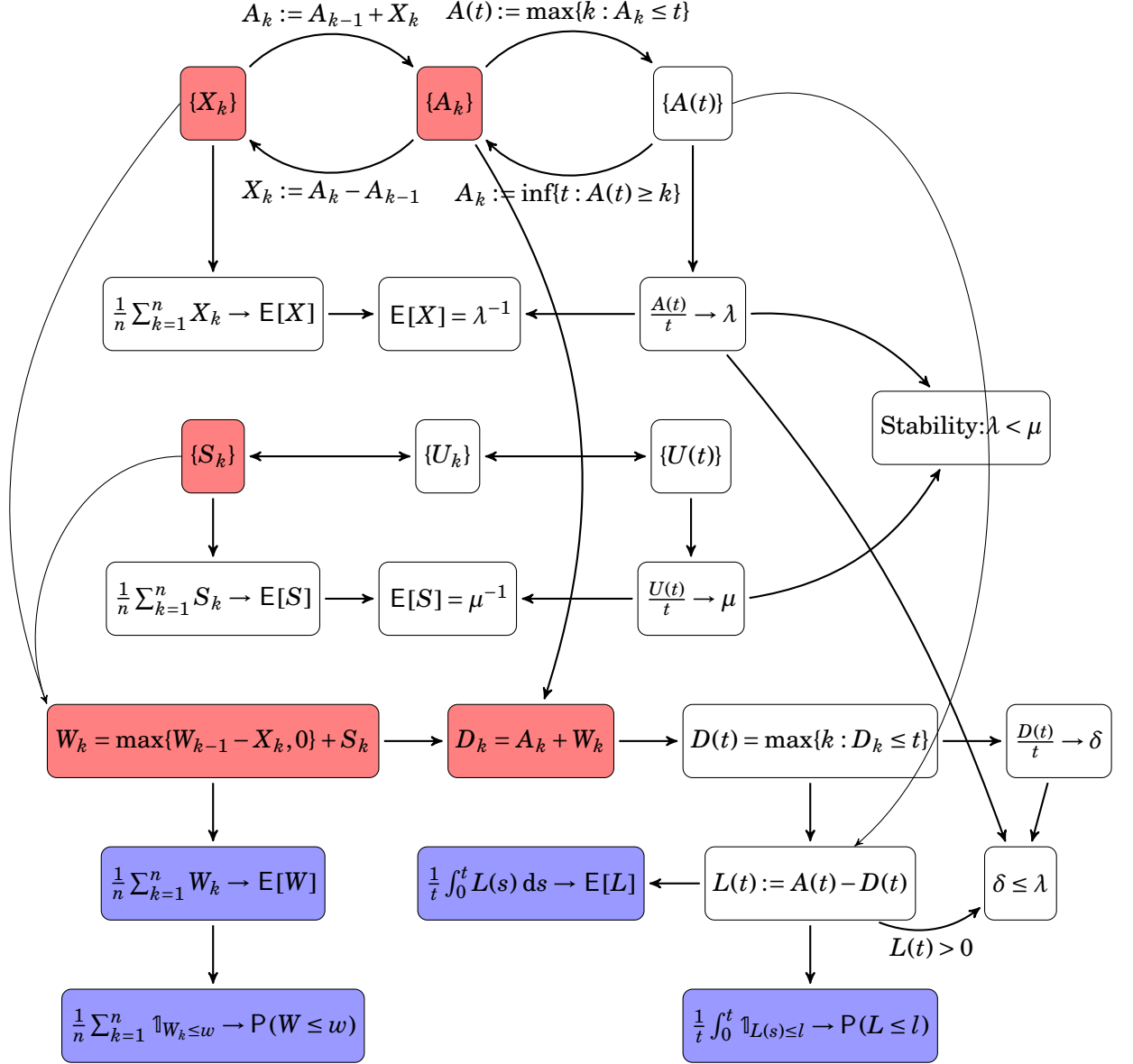


Figure 8: Here we sketch the relations between the construction of the $G/G/1$ queue from the primary data, i.e., the inter-arrival times $\{X_k; k \geq 0\}$ and the service times $\{S_k; k \geq 0\}$, and different performance measures.

APPROXIMATE QUEUEING MODELS

In this chapter we familiarize the reader with Sakasegawa's formula by which the expected queueing time in a $G/G/c$ can be approximated, and another formula by which it is possible to characterize how variability propagates through a tandem network of $G/G/c$ queues. We might say that these two formulas are the most important formulas to understand the behavior of queueing systems. With a bit of exaggeration, it is justified to say that the entire philosophy behind lean manufacturing and the world-famous Toyota production system are based on the principles that can be derived from these two formulas.

Here we take these formulas for granted and focus on the insights they provide and how to use them to guide improvement procedures for practical queueing problems that are, in some way or another, often encountered in production and service systems. In later sections, most notably Section 5.4, we will provide the theoretical background of Sakasegawa's formula.

In Section 3.1 we introduce and discuss the main insights of Sakasegawa's formula. Then, in the subsequent three sections, we illustrate how to use this formula to estimate waiting times in three queueing settings in which the service process is interrupted. In the first case, Section 3.2, the server has to produce jobs from different families, and there is a change-over time required to switch from one production family to another. As such setups reduce the time the server is available to serve jobs, the load must increase. In fact, to reduce the load, the server produces in batches of fixed sizes. In the second case, in Section 3.3, the server sometimes requires small adjustments, for instance, to prevent the production quality to degrade below a certain level. Clearly, such adjustments are typically not required during a job's service; however, they can occur at arbitrary moments in time. Thus, this is different from batch production in which the batch sizes are constant. In the third example, in Section 3.4, quality problems or break downs can occur during a job's service. In the final Section 3.5 we concentrate on tandem queues and study the consequences of two different scenarios to improve a simple network of queues. With the tools we develop here it becomes possible to analyze various scenarios to organize a process and quantitatively compare the scenarios.

In passing we use some interesting results of probability theory and the Poisson process; we will use these again, for instance in Chapter 6.

3.1 $G/G/c$ QUEUE: SAKASEGAWA'S FORMULA

In manufacturing settings it is quite often the case that the arrival process at a station is not a Poisson process. For instance, if processing times at a station are nearly constant, and the jobs of this station are sent to a second station for further processing, the inter-arrival times at the second station must be more or less equal too. Hence, in this case, the arrival process at this second station is most probably more regular than at the first station. As a second, trivial, case if the inter-arrival times of jobs are 1 hour always and service times 59 minutes always, there simply cannot be a queue.

In this section, we discuss Sakasegawa's formula by which we can estimate the expected waiting time in queue for the $G/G/c$ queue. The aim here is to show how to use this formula. In later sections we will provide some theoretical background.

While there is no closed-form expression available to compute the expected waiting time in queue for the $G/G/c$ queue, Sakasegawa's formula provides a reasonable approximation. This takes the form

$$E[W_Q] = \frac{C_a^2 + C_s^2}{2} \frac{\rho^{\sqrt{2(c+1)}-1}}{c(1-\rho)} E[S], \quad (3.1.1)$$

where λ is the rate at which jobs arrive at the system, $E[S]$ is the expected service time,

$$\rho = \frac{\lambda E[S]}{c}$$

is the *load* of the station (but not of the individual machines), $C_a^2 = V[X]/(E[X])^2$ is the SCV of the inter-arrival times, $C_s^2 = V[S]/(E[S])^2$ of the service times, and c the number of servers.

It is crucial to memorize the *scaling* relations that can be obtained from this formula. Even though the above results are only approximations, they prove to be exceedingly useful when designing queueing systems and analyzing the effect of certain changes, in particular changes in processing speed (i.e., service times) and service variability. Roughly:

1. $E[W_Q] \sim (1-\rho)^{-1}$. The consequence is that the waiting time increases *very steeply* when ρ is large. Hence, the waiting time is extremely sensitive to the actual value of ρ when ρ is large.
2. $E[W_Q] \sim C_a^2$ and $E[W_Q] \sim C_s^2$. Hence, reductions of the variation of the inter-arrival and service times do affect the waiting time, but only linearly.
3. $E[W_Q] \sim E[S]/c$. This means that, from the perspective of a job in queue, the average services are c times as short.

These insights prove very useful when trying to reduce waiting times in any practical situation. First try to reduce the load (by blocking demand or increasing the capacity), then try to reduce the variability (e.g., by planning the arrival times of jobs), and finally, attempt to split jobs into multiple smaller jobs and use the resulting freedom to reschedule jobs in the queue.

3.1.1. [3.1.1] Show for the $M/G/1$ and $M/M/1$ queue that the approximation (3.1.1) reduces to

$$E[W_Q(M/G/1)] = \frac{1+C_s^2}{2} \frac{\rho}{1-\rho} E[S], \quad E[W_Q(M/M/1)] = \frac{\rho}{1-\rho} E[S].$$

The above results explain part of the form of Sakasegawa's formula. In case $C_a^2 = C_s^2 = 1$, we obtain the expression for $E[W_Q(M/M/1)]$. As we will see in Section 5.1, this is an exact result. Then, by generalizing to the $M/G/1$ queue, we obtain $(1+C_s^2)/2$

Section 4.4, and Section 5.4, these results are exact.

3.1.2. [3.1.2] (Hall 5.19) When a bus reaches the end of its line, it undergoes a series of inspections. The entire inspection takes 5 minutes on average, with a standard deviation of 2 minutes. Buses arrive with inter-arrival times uniformly distributed on [3,9] minutes, hence, 10 buses arrive on average per hour.

As a first case, assuming a single server, estimate $E[W_Q]$ with the $G/G/1$ waiting time formula. As a second case, compare this result to an $M/G/1$ system with arrival rate 10 per hour and the same service time distribution.

Explain why the queueing time in the first setting is smaller.

The next exercise shows us that in a queue with finite capacity, there is a relation between the loss and the average number of servers occupied. Thus, if we want to reduce the loss, we need to reduce the load on the servers.

3.1.3. [3.1.3] Consider a queue with c servers, with generally distributed inter-arrival times, generally distributed service times, and the system can contain at most K customers, i.e., the $G/G/c/K$ queue. Let λ be the arrival rate, μ the service rate, β the long-run fraction of customers lost, and ρ the average number of busy/occupied servers. Show that

$$\beta = 1 - \rho \frac{\mu}{\lambda}.$$

Clearly, Sakasegawa's equation requires an estimate of C_a^2 and C_s^2 . Now it is not always easy in practice to determine the actual service time distribution, one reason being that service times are often only estimated by a planner, but not actually measured. Similarly, the actual arrival moments of jobs are often not registered, only just the date, or perhaps the hour, that a customer arrived. Hence, it is often not possible to estimate C_a^2 and C_s^2 from the information that is available. However, when for instance the number of arrivals per period has been logged for some time so that we know the arrivals per period $\{a_n, n = 1, \dots, N\}$ for some N , we can use this information instead of the inter-arrival times $\{X_k\}$ to obtain insight into C_a^2 . The relation we present here to compute C_a^2 from $\{a_n\}$ can of course also be applied to estimate C_s^2 .

Theorem 3.1.1. The SCV of the inter-arrival times can be estimated with the formula

$$C_a^2 \approx \frac{\tilde{\sigma}^2}{\tilde{\lambda}},$$

where

$$\tilde{\lambda} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n a_i, \quad \tilde{\sigma}^2 = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n a_i^2 - \tilde{\lambda}^2.$$

In words, $\tilde{\lambda}$ is the average number of arrivals per period, e.g., per day, and $\tilde{\sigma}^2$ is the variance of the number of arrivals per period.

Proof. The proof is based on an argument in Cox [1962]. We use quite a bit of the notation developed in Section 2.3. Let $\{A(t), t \geq 0\}$ be the number of arrivals that occur up to (and including) time t . We assume that $\{A(t)\}$ is a renewal process such that the inter-arrival times $\{X_k, k = 1, 2, \dots\}$ with $X_k = A_k - A_{k-1}$, are i.i.d. with mean $1/\lambda$ and standard deviation σ . (Observe that σ is not the same as $\tilde{\sigma}$ above.) Note that C_a^2 is defined in terms of λ and σ as:

$$C_a^2 = \frac{V[X_i]}{(E[X_i])^2} = \frac{\sigma^2}{1/\lambda^2} = \lambda^2 \sigma^2.$$

Next, let A_k be the arrival time of the k th arrival. The following useful relation between $A(t)$ and A_k enables us to prove our result (recall that we used a similar relation in the derivation of the Poisson process):

$$P(A(t) < k) = P(A_k > t).$$

Since the inter-arrival times have finite mean and second moment by assumption, we can apply the central limit law to obtain that, as $k \rightarrow \infty$,

$$\frac{A_k - k/\lambda}{\sigma\sqrt{k}} \rightarrow N(0, 1),$$

where $N(0, 1)$ is a standard normal random variable with distribution $\Phi(\cdot)$. Similarly,

$$\frac{A(t) - \lambda t}{\alpha \sqrt{t}} \rightarrow N(0, 1)$$

for an α that is yet to be determined. Thus, $E[A(t)] = \lambda t$ and $V[A(t)] = \alpha^2 t$.

Using that $P(N(0, 1) \leq y) = P(N(0, 1) > -y)$ (and $P(N(0, 1) = y) = 0$) we have that

$$\begin{aligned} \Phi(y) &\approx P\left(\frac{A_k - k/\lambda}{\sigma \sqrt{k}} \leq y\right) \\ &= P\left(\frac{A_k - k/\lambda}{\sigma \sqrt{k}} > -y\right) \\ &= P\left(A_k > \frac{k}{\lambda} - y\sigma \sqrt{k}\right). \end{aligned}$$

Define for ease

$$t_k = \frac{k}{\lambda} - y\sigma \sqrt{k}.$$

We can use the above relation between the distributions of $A(t)$ and A_k to see that $P(A_k > t_k) = P(A(t_k) < k)$. With this we get,

$$\begin{aligned} \Phi(y) &\approx P(A_k > t_k) \\ &= P(A(t_k) < k) \\ &= P\left(\frac{A(t_k) - \lambda t_k}{\alpha \sqrt{t_k}} < \frac{k - \lambda t_k}{\alpha \sqrt{t_k}}\right). \end{aligned}$$

Since $(A(t_k) - \lambda t_k)/\alpha \sqrt{t_k} \rightarrow N(0, 1)$ as $t_k \rightarrow \infty$, the above implies that

$$\frac{k - \lambda t_k}{\alpha \sqrt{t_k}} \rightarrow y,$$

as $t_k \rightarrow \infty$. Using the above definition of t_k , the left-hand side of this equation can be written as

$$\frac{k - \lambda t_k}{\alpha \sqrt{t_k}} = \frac{\lambda \sigma \sqrt{k}}{\alpha \sqrt{k/\lambda + \sigma \sqrt{k}}} y.$$

Since $t_k \rightarrow \infty$ is implied by (and implies) $k \rightarrow \infty$, we therefore want that α is such that

$$\frac{\lambda \sigma \sqrt{k}}{\alpha \sqrt{k/\lambda + \sigma \sqrt{k}}} y \rightarrow y,$$

as $k \rightarrow \infty$. This is precisely the case when

$$\alpha = \lambda^{3/2} \sigma.$$

Finally, for t large (or, by the same token k large),

$$\frac{\sigma_k^2}{\lambda_k} = \frac{V[A(t)]}{E[A(t)]} \approx \frac{\alpha^2 t}{\lambda t} = \frac{\alpha^2}{\lambda} = \frac{\lambda^3 \sigma^2}{\lambda} = \lambda^2 \sigma^2 = C_a^2,$$

where the last equation follows from the above definition of C_a^2 . The proof is complete. \square

3.2 SETUPS AND BATCH PROCESSING

In some cases machines have to be setup before they can start producing items. Consider, for instance, a machine that paints red and blue bikes. When the machine requires a color change, it may be necessary to clean up the machine. Another example is an oven that needs a temperature change when different item types require different production temperatures. Service operations form another setting with setup times: when servers (personnel) have to move from one part of a building to another, the time spent moving cannot be spent on serving customers. In all such cases, the setups consume a significant amount of time; in fact, setup times of an hour or longer are not uncommon. Clearly, in such situations, it is necessary to produce in batches: a server processes a batch of jobs (or customers) of one type or at one location, then the server changes from type or location, starts serving a batch of another type or at another location. Once done with one type, the server is setup again, and so on. Here we focus on the effect of change-over, or setup, times on the average sojourn time of jobs.

First we make a model and provide a list of elements required to compute the expected sojourn time of an item, then we illustrate how to use these elements in a concrete case.

Specifically, we analyze the following batch queueing situation. There are two job families, e.g., red and blue, each served by the same single server. Jobs arrive at rate λ_r and λ_b , respectively, so that the arrival rate of jobs is $\lambda = \lambda_b + \lambda_r$. Jobs of both types require an average *net processing time* of $E[S_0]$, provided the server is already setup for the correct job color. The setup times $\{R_i\}$ are assumed to form an i.i.d. sequence with common random variable R and independent of S_0 . The sojourn time comprises the following steps. First, jobs of each color are assembled into batches of size B , which we assume to be the same for both colors. Once a batch is complete, the batch enters a queue (of batches). After some time the batch reaches the head of the queue. Then the machine performs a setup, and starts processing each job individually until the batch is complete. Finally, once a job is finished, it can leave the system; as a consequence, it does not have to wait for other jobs in the same batch to form a new batch.

3.2.1. [3.2.1] *Show that the average time a job has to wait to fill the batch (to which this job belongs) is given by*

$$E[W_r] = \frac{B-1}{2\lambda_r}. \quad (3.2.1)$$

Now that we have a batch of jobs, we need to estimate the average time a batch spends in queue. For this we can use the $G/G/1$ waiting time formula, but we have to convert the effects of the setup times into job service times. Define, to this end, the *effective processing time* $E[S]$ as the average time the server is occupied with processing a job from a batch including the setup time required to setup the batch.

3.2.2. [3.2.2] *Motivate that the effective processing time of an item should be defined as*

$$E[S] = E[S_0] + \frac{E[R]}{B}. \quad (3.2.2)$$

With the previous exercise, the load becomes

$$\rho = \lambda \left(E[S_0] + \frac{E[R]}{B} \right).$$

There is another important way to look at the load, see the next exercise.

3.2.3. [3.2.3] Explain that the load can also be written as

$$\rho = \lambda_B (B E[S_0] + E[R]),$$

where $\lambda_B = \lambda/B$ is the arrival rate of batches and $E[S_B] = B E[S_0] + E[R]$ is the service time an entire batch.

3.2.4. [3.2.4] Show that the requirement $\rho < 1$ leads to the following constraint on the minimal batch size B

$$B > \frac{\lambda E[R]}{1 - \lambda E[S_0]}.$$

Observe that we nearly have all elements to apply Sakasegawa's formula to the batch: the service time of a batch is

$$E[S_B] = E[R] + B E[S_0], \quad (3.2.3)$$

and the load ρ is given above. Therefore it remains to find $C_{a,B}^2$ and $C_{s,B}^2$, for the batches, not the items.

3.2.5. [3.2.5] Explain that the SCV of the batch inter-arrival times is given by

$$C_{a,B}^2 = \frac{C_a^2}{B}. \quad (3.2.4)$$

3.2.6. [3.2.6] Show that the SCV $C_{s,B}^2$ of the service times of the batches takes the form

$$C_{s,B}^2 = \frac{B V[S_0] + V[R]}{(B E[S_0] + E[R])^2}. \quad (3.2.5)$$

Observe that we now have all elements to fill in Sakasegawa's formula, which becomes

$$E[W_{Q,B}] = \frac{C_{a,B}^2 + C_{s,B}^2}{2} \frac{\rho}{1 - \rho} E[S_B].$$

It is left to find a rule to determine what happens to an item once the batch to which the item belongs enters service. If the job has to wait until all jobs in the batch are served, the expected time an item spends at the server is $E[R] + B E[S_0]$.

3.2.7. [3.2.7] Show that, when items can leave right after being served, the time at the server is given by

$$E[R] + \frac{B+1}{2} E[S_0]. \quad (3.2.6)$$

3.2.8. [3.2.8] Jobs arrive at $\lambda = 3$ per hour at a machine with $C_a^2 = 1$; service times are exponential with an average of 15 minutes. Assume $\lambda_r = 0.5$ per hour, hence $\lambda_b = 3 - 0.5 = 2.5$ per hour. Between any two batches, the machine requires a cleanup of 2 hours, with a standard deviation of 1 hour, during which it is unavailable for service. What is the smallest batch size that can be allowed?

What is the average time a red job spends in the system in case $B = 30$ jobs? Finally, observe that there is B that minimizes the average sojourn time.

3.2.9. [3.2.9] What important insights can you learn from the above about setting proper batch sizes?

A much more interesting and realistic problem is to assume that we have many families such that items of family i arrive at rate λ_i and require service time S_i . Typically, the setup time depends on the sequence in which the families are produced; let this be given by R_{ij} when the machine switched from producing family i to family j . Often, in practice, $R_{ij} \neq R_{ji}$, for example, a switch in color from white to black takes less cleaning time and cost than from black to white. Then the problem becomes to determine a good schedule in which to produce the families and the batch size B_i , which may depend on the family. Here is an exercise to show how to handle a simple case.

3.2.10. [3.2.10] Consider a paint factory that contains a paint mixing machine that serves two classes of jobs, A and B. The processing times of jobs of types A and B are constant and require t_A and t_B hours. The job arrival rate is λ_A for type A and λ_B for type B jobs. It takes a setup time of S hours to clean the mixing station when changing from paint type A to type B, and there is no time required to change from type B to A.

To keep the system stable, it is necessary to produce the jobs in batches, for otherwise the server, i.e., the mixing machine, spends a too large fraction of time on setups. Motivate that the following linear program can be used to determine the minimal batch sizes:

$$\text{minimize } T$$

$$\text{such that } T = k_A t_A + S + k_B t_B, \lambda_A T < k_A \text{ and } \lambda_B T < k_B.$$

3.3 NON-PREEMPTIVE INTERRUPTIONS, SERVER ADJUSTMENTS

In Section 3.2 we studied the effect of setup times between job batches with a fixed size B . However, other types of interruptions can occur, such as a machine requiring random adjustments that can occur between any two jobs. This type of outages is *non-preemptive* as the outages do not interrupt the processing of a job in service. In this section we develop a simple model to understand the impact of such outages on job sojourn times; we use the same notation as in Section 3.2 and follow the same line of reasoning. With this model, we can analyze the effects of reducing adjustment times or the variability of these adjustments times. For instance, we might decide to do fewer adjustments, but the average outage times become longer.

We assume that adjustments $\{R_i\}$ occur geometrically distributed between any two jobs with a mean of B jobs between any two adjustments. Consequently, the probability of an outage between any two jobs is $p = 1/B$. Observe that geometrically distributed random variables satisfy the memoryless property in discrete time. Hence, our assumption implies that the occurrence of an adjustment between jobs i and $i + 1$ has no effect on the probability that an adjustment is necessary between jobs $i + 1$ and $i + 2$.

Contrary to the batch processing case of Section 3.2, we now only need to find the mean and variance of the effective processing times.

3.3.1. [3.3.1] Show that the average effective processing time is given by

$$E[S] = E[S_0] + \frac{E[R]}{B} \quad (3.3.1)$$

Clearly, the effective server load including down-times is $\rho = \lambda E[S]$. Since the adjustments do not affect the job arrival process, we only have to find an expression how it affects C_s^2 .

3.3.2. [3.3.3] Derive

$$V[S] = V[S_0] + \frac{V[R]}{B} + (B-1) \left(\frac{E[R]}{B} \right)^2. \quad (3.3.2)$$

With this, we have all elements to fill in the $G/G/1$ waiting time formula!

3.3.3. [3.3.4] Jobs arrive as a Poisson process with rate $\lambda = 9$ per working day. The machine works two 8 hour shifts a day. Work not processed on a day is carried over to the next day. Job service times are 1.5 hours, on average, with standard deviation of 0.5 hours. Outages occur on average between 30 jobs. The average duration of an outage is 5 hours and has a standard deviation of 2 hours.

Compute the average waiting time in queue.

Remark 3.3.1. *Let us compare the results of this section to those of Section 3.3. In both cases the expected service time is the same: an amount $E[R]/B$ gets added to the net processing time $E[S_0]$. However, the impact on the SCV is different. Unexpected interruptions must have a larger effect on variability than expected interruptions.*

In more general terms, we can use part of the model of Section 3.2 to analyze the effect of planned maintenance—do an adjustment after every B jobs—while with the model of this section we consider adjustments that arrive completely unexpectedly. With these two models we can compare the influence on waiting times of doing planned adjustments more often (take B a bit small), but remove the unexpected interruptions. Hence, we have another set of tools to compare different scenarios to organize production systems.

3.4 PREEMPTIVE INTERRUPTIONS, SERVER FAILURES

In Sections 3.2 and 3.3 we assumed that servers are never interrupted while serving a job. However, in many situations this assumption is not satisfied: a person might receive a short phone call while working on a job, a machine may fail in the midst of processing, and so on. In this section, we develop a model to compute the influence on the mean waiting time of such *preemptive outages*, i.e., interruptions that occur *during* a service. As in Section 3.3, we can use Sakasegawa's formula to estimate expected queueing time for the $G/G/1$ queue. For this, it suffices to find expressions for $E[S]$ and $V[S]$.

We will first present the main formulas and show how to apply them. The exercises at the end of the section ask you to provide the derivations. Specifically, here we compute the expectation and variance of a random number of failures, each with a random duration. Similar problems occur in, for instance, insurance theory (or inventory theory) where we are interested in (the distribution of) the total claim (demand) size that occurs in a period. Noting that the total claim (demand) is given by a random number of claims (demands), each of which is a random variable by itself, that problem is exactly the same as the problem here.

The techniques we develop here are of general interest, and they provide relations between the topics discussed in Section 0.2, Section 1.1, Section 5.6, and Section 6.2.

As in the previous sections, we derive expressions for the expectation and variance of the effective processing time S .

Supposing that N interruptions occur during the net service time S_0 of a job, and the repair times $\{R_i\}$ form an i.i.d. sequence distributed as a common random variable R , write

$$S_N = \sum_{i=1}^N R_i \quad (3.4.1)$$

for the total duration of the interruptions that occur during the service of a job. Then, the effective service time becomes

$$S = S_0 + S_N = S_0 + \sum_{i=1}^N R_i.$$

It is important to realize that N is a random number.

A common assumption is that the time between two interruptions is memoryless; let the interruptions occur at a given *failure rate* λ_f . Then, if the net service time S_0 is a constant, the expected number of failures $E[N]$ is $\lambda_f S_0$. More generally, we show below that $E[N] = \lambda_f E[S_0]$.

Define the *availability* as

$$A = \frac{m_f}{m_f + m_r},$$

where m_f is the mean time to fail and m_r the mean time to repair. From 3.4.2 it follows that

$$A = \frac{1}{1 + \lambda_f E[R]}. \quad (3.4.2)$$

Then, from 3.4.4,

$$E[S] = \frac{E[S_0]}{A}. \quad (3.4.3)$$

Now the load is easy to compute:

$$\rho = \lambda E[S] = \lambda \frac{E[S_0]}{A}.$$

Clearly, the server load increases due to failures.

By assuming that repair times are exponentially distributed with mean $E[R]$, we can find with the exercises below that

$$C_s^2 = C_0^2 + 2A(1 - A) \frac{E[R]}{E[S_0]}, \quad (3.4.4)$$

where C_0^2 is the SCV of S_0 , i.e., the service time without interruptions.

Before deriving the above expressions, let us see how to apply them.

3.4.1. [3.4.1] Suppose we have a machine with memoryless failure behavior, with a mean-time-to-fail of 3 hours. Regular service times are deterministic with an average of 10 minutes, jobs arrive as a Poisson process with rate of 4 per hour. Repair times are exponential with a mean duration of 30 minutes. What is the average sojourn time?

The rest of the section is concerned with the derivations of the formulas.

3.4.2. [3.4.4] Derive (3.4.2) for our model of interruptions.

3.4.3. [3.4.7] Show that $E[N] = \lambda_f E[S_0]$ and $E[N^2] = \lambda_f^2 E[S_0^2] + \lambda_f E[S_0]$.

3.4.4. [3.4.8] Show that $E[S] = E[S_0] + \lambda_r E[S_0] E[R]$, and conclude that $E[S] = E[S_0]/A$.

3.4.5. [3.4.9] The derivation of C_s^2 is a bit more involved. To understand why, explain first that $V[S] \neq V[S_0] + V[\sum_{i=1}^N R_i]$.

3.4.6. [3.4.10] Show that

$$E[S^2] = E[S_0^2] + 2E\left[S_0 \sum_{i=1}^N R_i\right] + E\left[\sum_{i=1}^N R_i^2\right] + E\left[\sum_{i=1}^N \sum_{j \neq i} R_i R_j\right].$$

Principally we are done now. With $E[S^2]$ and $E[S]$ we can compute the SCV via $(E[S^2] - (E[S])^2)/(E[S])^2$, so that we have all components for Sakasegawa's formula. The rest is polishing.

3.4.7. [3.4.16] *Show that*

$$C_s^2 = \frac{V[S]}{(E[S])^2} = C_0^2 + \frac{\lambda_f E[R^2] A^2}{E[S_0]},$$

3.4.8. [3.4.17] *With the above assumption that R is exponentially distributed, show that*

$$C_s^2 = C_0^2 + 2A(1-A) \frac{E[R]}{E[S_0]}.$$

3.5 TANDEM QUEUES

Consider two $G/G/1$ stations in tandem. Suppose we have the financial means to reduce the variability of the processing times at one of the servers, but not at both. Then we like to improve the one that has the most impact on the total waiting time in the system. To obtain insight into this problem, we present a formula that approximates the SCV of the inter-departure times of a $G/G/c$ queue. Noting that the output of the first machine forms the input of the second machine, the SCV of the inter-departure times of first station must then be the SCV of the inter-arrival times at the second station. Thus, with this formula for the SCV of the inter-departure times, we can model the propagation of variability through a simple network of $G/G/c$ stations in tandem.

We remark that the literature contains algorithms that can deal with more complicated networks of $G/G/c$ queues in which the output streams of several stations merge into the input stream of another station and rework is allowed. However, when the machines act like $M/M/c$ queues, it is possible to fully analyze the system, cf. Section 6.3.

Theory and Exercises

With simulation it has been tested that the SCV of the inter-departure times of a $G/G/c$ queue can be reasonable well approximated by

$$C_d^2 \approx 1 + (1 - \rho^2)(C_a^2 - 1) + \frac{\rho^2}{\sqrt{c}}(C_s^2 - 1). \quad (3.5.1)$$

3.5.1. [3.5.1] *Show that (3.5.1) reduces to the following for the $G/G/1$ queue:*

$$C_d^2 = (1 - \rho^2)C_a^2 + \rho^2 C_s^2. \quad (3.5.2)$$

Clearly, for single-server queues the expression becomes quite simple for which we provide the following intuitive interpretation. Suppose that the load ρ is very high. Then the server will seldom be idle, so that most of inter-departure times are the same as service times. If, however, the load is low, the server will be idle most of the time, and the inter-departure times are approximately equal to the inter-arrival times. The approximation then consists of an interpolation between these two extremes.

3.5.2. [3.5.2] *What is C_d^2 for the $D/D/1$ queue according to (3.5.2)?*

3.5.3. [3.5.3] *Show that $C_d^2 = 1$ for the $M/M/1$ queue according to (3.5.2). In fact, in 6.3.1 we provide a proof that this is an exact result.*

The next exercise proves a very useful insight: if we make service times more regular, the departure process also becomes more regular.

3.5.4. [3.5.4] Use (3.5.2) to show for the $G/D/1$ that $C_d^2 < C_a^2$.

3.5.5. [3.5.5] Consider two $G/G/1$ stations in tandem. Suppose $\lambda = 2$ per hour, $C_{a,1}^2 = 2$ at station 1, $C_s^2 = 0.5$ at both stations, and $E[S_1] = 20$ minutes and $E[S_2] = 25$ minutes. What is the total time jobs spend on average in the system?

Let us next apply the above to analyze the performance of a tandem network of two $M/M/1$ queues. For this, assume that jobs arrive at the first station at rate λ , and are served at rate μ_i at station i . Thus, $\rho_i = \lambda/\mu_i$ and $E[S_i] = 1/\mu_i$, for $i = 1, 2$. First we consider the base case, i.e., the system without any improvement at either of the servers. Then we assume we can remove all variability at the second station; this is the best that we can possibly do at station 2. Next, we assume we can remove all variability at the first station. The final step is to compare the different improvement scenarios.

To analyze this case, observe that the departures of the first station form the arrivals at the second station. Thus, the SCV of the inter-arrival times at the second station is the SCV of the inter-departure times of the first station, that is $C_{d,1}^2 = C_{a,2}^2$. This idea can of course be used for longer tandem networks.

3.5.6. [3.5.6] Use Sakagawa's formula to show that the average queueing time for the tandem of two $M/M/1$ queues is given by

$$E[W_Q] = \frac{\rho_1}{1 - \rho_1} \frac{1}{\mu_1} + \frac{\rho_2}{1 - \rho_2} \frac{1}{\mu_2}, \quad (3.5.3)$$

Now suppose we can remove all variability of the service process at the second station.

3.5.7. [3.5.7] Show that in this case the total time in queue is equal to

$$E[W_Q] = \frac{\rho_1}{1 - \rho_1} \frac{1}{\mu_1} + \frac{1}{2} \frac{\rho_2}{1 - \rho_2} \frac{1}{\mu_2}.$$

Suppose now that we reduce the variability of the service process of the first station.

3.5.8. [3.5.8] Motivate that

$$E[W_Q] = \frac{1}{2} \frac{\rho_1}{1 - \rho_1} \frac{1}{\mu_1} + \frac{2 - \rho_1^2}{2} \frac{\rho_2}{1 - \rho_2} \frac{1}{\mu_2}$$

is a reasonable approximation of the queueing time in the network.

3.5.9. [3.5.9] Comparing these three scenarios, what do you conclude?

FUNDAMENTAL TOOLS

To develop mathematical models of queueing systems we need a few concepts that are fundamentally important and have a general interest beyond queueing. All of these concepts rely on *sample-path constructions* of queueing, or more general stochastic, systems. Thus, sample paths, which are nothing but the output of a simulation, form an elegantly unifying principle.

Here we keep the discussion in these notes mostly at an intuitive level; we refer to [El-Taha and Stidham Jr. \[1998\]](#) for proofs and further background.

4.1 RENEWAL REWARD THEOREM

We state the *renewal reward theorem* and provide graphical motivation for its validity. This results proves to be extremely useful: we will apply it regularly in the sequel of the book and here we use it to provide another definition of the load ρ .

In essence the renewal reward theorem is very simple: it states that when customers arrive at rate λ and each customer pays an average amount X , the system earns money at rate $Y = \lambda X$. Figure 9 provides graphical motivation about why this theorem is true; [El-Taha and Stidham Jr. \[1998\]](#) gives a (simple) proof.

Theorem 4.1.1 (Renewal Reward Theorem, $Y = \lambda X$). *Consider epochs $\{T_k, k = 0, 1, \dots\}$ such that $0 = T_0 < T_1 < \dots$. Let $N = \{N(t), t \geq 0\}$ be the associated counting process with $N(t) = \max\{k : T_k \leq t\}$. Let $\{Y(t), t \geq 0\}$ be a non-decreasing right-continuous (deterministic) process. Define $X_k = Y(T_k) - Y(T_{k-1})$. Suppose that $N(t)/t \rightarrow \lambda$ as $t \rightarrow \infty$, where $0 < \lambda < \infty$. Then $Y(t)/t$ has a limit iff $n^{-1} \sum_{k=1}^n X_k$ has a limit, and then $Y = \lambda X$. In other words,*

$$\lim_{t \rightarrow \infty} \frac{Y(t)}{t} = Y \iff \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n X_k = X,$$

and then $Y = \lambda X$.

Let us use this theorem to understand the idea of ‘load’ in a different way. Define the *load* or *utilization* as the limiting fraction of time the server is busy, i.e.,

$$\rho = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t \mathbb{1}_{L(s) > 0} ds.$$

4.1.1. [4.1.1] Use the renewal reward theorem to prove that $\rho = \lambda E[S]$ for the rate-stable G/G/1 queue.

4.1.2. [4.1.2] We can derive the relation $\rho = \lambda E[S]$ in a somewhat more direct way by considering the fact that

$$\sum_{k=1}^{A(t)} S_k \geq \int_0^t \mathbb{1}_{L(s) > 0} ds \geq \sum_{k=1}^{D(t)} S_k.$$

Explain this, and conclude that

$$\lambda E[S] \geq \rho \geq \delta E[S]. \quad (4.1.1)$$

Hence, if the system is rate-stable, $\delta = \lambda$, and then $\rho = \lambda E[S]$.



Figure 9: A graphical ‘proof’ of $Y = \lambda X$. Here $Y(t)/t \rightarrow Y$, $n/T_n \rightarrow \lambda$ and $n^{-1} \sum_i^n X_i \rightarrow X$. Observe that in the figure X_k does not represent an inter-arrival time; instead it corresponds to the increment of (the graph of) $Y(t)$ between two consecutive epochs T_{k-1} and T_k at which $Y(t)$ is observed.

4.2 LEVEL CROSSING AND BALANCE EQUATIONS

Let us say that the system is in *state* n at time t when it contains n jobs at that moment, i.e., when $L(t) = n$. The system *up-crosses level* n at time t when its state changes from n to $n + 1$ due to an arrival, and it *down-crosses level* n when its state changes from $n + 1$ to n due to a departure. Clearly, the number of up-crossings and down-crossings cannot differ by more than 1, because it is only possible to up-cross level n after a down-crossing (or the other way around). This simple idea will prove a key stepping stone in the analysis of queueing systems.

To establish the section’s main result (4.2.6), we need a few definitions that are quite subtle and might seem a bit abstract, but below we will provide intuitive interpretations in terms of system KPIs. After this, we will generalize the principle of level-crossing to *balance equations* which allow us to deal with more general types of transitions. We refer to Fig. 13 in which all concepts developed here are summarized.

LEVEL CROSSING Define

$$A(n, t) = \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq t} \mathbb{1}_{L(A_k-) = n} \quad (4.2.1a)$$

as the number of arrivals up to time t that saw n customers in the system upon their arrival, cf. Fig. 10.

Next, let

$$Y(n, t) = \int_0^t \mathbb{1}_{L(s) = n} ds \quad (4.2.1b)$$

be the total time the system contains n jobs during $[0, t]$, and

$$p(n, t) = \frac{1}{t} \int_0^t \mathbb{1}_{L(s) = n} ds = \frac{Y(n, t)}{t}, \quad (4.2.1c)$$

be the fraction of time that $L(s) = n$ in $[0, t]$. Figure 11 illustrates the relation between $Y(n, t)$ and $A(n, t)$.

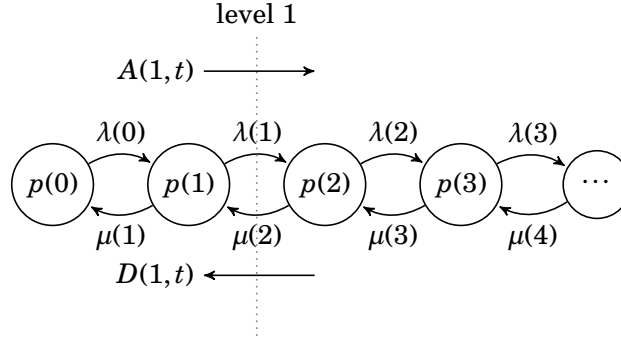


Figure 10: $A(1, t)$ counts the number of jobs up to time t that saw 1 job in the system upon arrival, and right after such arrivals the system contains 2 jobs. Thus, each time $A(1, t)$ increases by one, level 1 (the dotted line separating states 1 and 2) is crossed from below. Similarly, $D(1, t)$ counts the number of departures that leave 1 job behind, and just before such departures the system contains 2 jobs. Hence, level 1 is crossed from above. It is evident that the number of times this level is crossed from below must be the same (plus or minus 1) as the number of times it is crossed from above. (We introduce $\lambda(n)$, $\mu(n)$ and $p(n)$ below.)



Figure 11: Plots of $Y(1, t)$ and $A(1, t)$. (For visual clarity, we subtracted $1/2$ from $A(1, t)$, for otherwise its graph would partly overlap with the graph of Y .)

4.2.1. [4.2.5] Consider the following (silly) queueing process. At times $0, 2, 4, \dots$ customers arrive, each customer requires 1 unit of service, and there is one server. Find an expression for $A(n, t)$. (What acronym would describe this queueing situation?)

4.2.2 (Continuation of 4.2.1). [4.2.6] Find an expression for $Y(n, t)$.

Define also the limits:

$$\lambda(n) = \lim_{t \rightarrow \infty} \frac{A(n, t)}{Y(n, t)}, \quad p(n) = \lim_{t \rightarrow \infty} p(n, t), \quad (4.2.2)$$

as the *arrival rate in state n* and the *long-run fraction of time the system spends in state n* . To clarify the former definition, observe that $A(n, t)$ counts the number of arrivals that see n jobs in the system upon arrival, while $Y(n, t)$ tracks the amount of time the system contains n jobs. Suppose that at time T a job arrives that sees n jobs in the system. Then $A(n, T) = A(n, T-) + 1$, and this job finishes an interval that is tracked by $Y(n, t)$, precisely because this job sees n jobs in the system just prior to its arrival. Thus, just as $A(t)/t$ is the total number of arrivals during $[0, t]$ divided by t , $A(n, t)/Y(n, t)$ is the number of arrivals that see n jobs divided by the time the system contains n jobs.

4.2.3 (Continuation of 4.2.2). [4.2.7] Compute $p(n)$ and $\lambda(n)$.

Similar to the definition for $A(n, t)$, let

$$D(n, t) = \sum_{k=1}^{\infty} \mathbb{1}_{D_k \leq t} \mathbb{1}_{L(D_k)=n}$$

denote the number of departures up to time t that leave n customers behind. Then, define

$$\mu(n+1) = \lim_{t \rightarrow \infty} \frac{D(n, t)}{Y(n+1, t)},$$

as the departure rate from state $n+1$. (It is easy to get confused here: to leave n jobs behind, the system must contain $n+1$ jobs just prior to the departure.) Figure 10 shows how $A(n, t)$ and $\lambda(n)$ relate to $D(n, t)$ and $\mu(n+1)$.

4.2.4 (Continuation of 4.2.3). [4.2.9] Compute $D(n, t)$ and $\mu(n+1)$ for $n \geq 0$.

Observe that customers arrive and depart as single units. Thus, if $\{T_k\}$ is the ordered set of arrival and departure times of the customers, then $L(T_k) = L(T_k-) \pm 1$. But then we must also have that

$$|A(n, t) - D(n, t)| \leq 1, \quad (4.2.3)$$

(think about this). From this observation it follows immediately that

$$\lim_{t \rightarrow \infty} \frac{A(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{D(n, t)}{t}. \quad (4.2.4)$$

With this equation we can obtain two nice and fundamental identities. The first we develop now; the second follows in Section 4.3.

The rate of jobs that ‘see the system with n jobs’ can be defined as $A(n, t)/t$. Taking limits we get

$$\lim_{t \rightarrow \infty} \frac{A(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{A(n, t)}{Y(n, t)} \frac{Y(n, t)}{t} = \lambda(n)p(n), \quad (4.2.5a)$$

where we use the above definitions for $\lambda(n)$ and $p(n)$. Similarly, the departure rate of jobs that leave n jobs behind is

$$\lim_{t \rightarrow \infty} \frac{D(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{D(n, t)}{Y(n+1, t)} \frac{Y(n+1, t)}{t} = \mu(n+1)p(n+1). \quad (4.2.5b)$$

Combining this with (4.2.4) we arrive at the *level-crossing equations*

$$\lambda(n)p(n) = \mu(n+1)p(n+1). \quad (4.2.6)$$

4.2.5 (Continuation of 4.2.4). [4.2.10] Compute $\lambda(n)p(n)$ for $n \geq 0$, and check $\lambda(n)p(n) = \mu(n+1)p(n+1)$.

Result (4.2.6) turns out to be exceedingly useful, as will become evident from Section 5.1 onward. More specifically, by specifying (i.e., modeling) $\lambda(n)$ and $\mu(n)$, we can compute the long-run fraction of time $p(n)$ that the system contains n jobs. To see this, rewrite the above into

$$p(n+1) = \frac{\lambda(n)}{\mu(n+1)} p(n). \quad (4.2.7)$$

Thus, this equation fixes the ratios between the probabilities. In other words, if we know $p(n)$ we can compute $p(n+1)$, and so on. Hence, if $p(0)$ is known, then $p(1)$ follows, from which $p(2)$ follows, and so on. A straightaway iteration then leads to

$$p(n+1) = \frac{\lambda(n)\lambda(n-1)\cdots\lambda(0)}{\mu(n+1)\mu(n)\cdots\mu(1)}p(0). \quad (4.2.8)$$

To determine $p(0)$ we can use the fact that the numbers $p(n)$ represent probabilities. Hence, from the normalizing condition $\sum_{n=0}^{\infty} p(n) = 1$, we get $p(0) = G^{-1}$ with G being the *normalization constant*

$$G = 1 + \sum_{n=0}^{\infty} \frac{\lambda(n)\lambda(n-1)\cdots\lambda(0)}{\mu(n+1)\mu(n)\cdots\mu(1)}. \quad (4.2.9)$$

In the next few sections we will make suitable choices for $\lambda(n)$ and $\mu(n)$ to model many different queueing situations so that, based on (4.2.6), we can obtain simple expressions for $p(n)$ in terms of the arrival and service rates.

With $p(n)$ we define two easy, but important performance measures. The time-average number of items in the system becomes

$$E[L] = \sum_{n=0}^{\infty} np(n),$$

and the long-run fraction of time the system contains at least n jobs is

$$P(L \geq n) = \sum_{i=n}^{\infty} p(i).$$

4.2.6. [4.2.11] Derive $E[L] = \sum_{n=0}^{\infty} np(n)$ from (2.4.6).

Finally, the following two exercises show that level-crossing arguments extend well beyond the queueing systems modeled by Fig. 10.

4.2.7. [4.2.12] Consider a single server that serves one queue and serves only in batches of 2 jobs at a time (so never 1 job or more than 2 jobs), i.e., the $M/M^2/1/3$ queue. Single jobs arrive at rate λ and the inter-arrival times are exponentially distributed so that we can assume that $\lambda(n) = \lambda$. The batch service times are exponentially distributed with mean $1/\mu$. Then, by the memoryless property, $\mu(n) = \mu$. At most 3 jobs fit in the system. Make a graph of the state-space and show, with arrows, the transitions that can occur.

4.2.8. [4.2.13] Use the graph of 4.2.7 and a level-crossing argument to express the steady-state probabilities $p(n)$, $n = 0, \dots, 3$ in terms of λ and μ .

INTERPRETATION The definitions in (4.2.1) may seem a bit abstract, but they obtain an immediate interpretation when relating them to applications. To see this, we discuss two examples.

Consider the sorting process of post parcels at a distribution center of a post-delivery company. Each day tens of thousands of incoming parcels have to be sorted to their final destination. In the first stage of the process, parcels are sorted to a region in the Netherlands. Incoming parcels are deposited on a conveyor belt. From the belt, they are carried to outlets (chutes), each chute corresponding to a specific region. Employees take out the parcels from the chutes and put the parcels in containers. The arrival rate of parcels for a certain chute may temporarily exceed the working capacity of the employees, as such the chute serves as a queue.

When the chute overflows, parcels are directed to an overflow container and are sorted the next day. The target of the sorting center is to deliver at least a certain percentage of the parcels within one day. Thus, the fraction of parcels rejected at the chute should remain small.

Suppose a chute can contain at most 20 parcels, say. Then, each parcel on the belt that ‘sees’ 20 parcels in its chute will be blocked. Let $L(t)$ be the number of parcels in the chute at time t . Then, $A(20, t)$ as defined in (4.2.1a) is the number of *blocked parcels* up to time t , and $A(20, t)/A(t)$ is the fraction of rejected parcels. In fact, $A(20, t)$ and $A(t)$ are continuously tracked by the sorting center and used to adapt employee capacity to control the fraction of rejected parcels. Thus, in simulations, if one wants to estimate loss fractions, $A(n, t)/A(t)$ is the most natural concept to consider.

For the second example, suppose there is a cost associated with keeping jobs in queue. Let w be the cost per job in queue per unit time so that the cost rate is nw when n jobs are in queue. But then $wnY(n, t)$ is the total cost up to time t to have n jobs in queue, hence the total cost up to time t is

$$C(t) = w \sum_{n=0}^{\infty} nY(n, t),$$

and the average cost is

$$\frac{C(t)}{t} = w \sum_{n=0}^{\infty} n \frac{Y(n, t)}{t} = w \sum_{n=0}^{\infty} np(n, t).$$

All in all, the concepts developed above have natural interpretations in practical queueing situations; they are useful in theory and in simulation, as they relate the theoretical concepts to actual measurements.

BALANCE EQUATIONS It is important to realize that the level-crossing argument cannot always be used as we do here. The reason is that sometimes there does not exist a line between two states such that the state space splits into two disjoint parts. For a more general approach, we focus on a single state and count how often this state is entered and left, cf. Fig. 12. Specifically, define

$$I(n, t) = A(n-1, t) + D(n, t),$$

as the number of times the queueing process enters state n either due to an arrival from state $n-1$ or due to a departure leaving n jobs behind. Similarly,

$$O(n, t) = A(n, t) + D(n-1, t),$$

counts how often state n is left either by an arrival (to state $n+1$) or a departure (to state $n-1$).

Of course, $|I(n, t) - O(n, t)| \leq 1$. Thus, from the fact that

$$\lim_{t \rightarrow \infty} \frac{I(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{A(n-1, t)}{t} + \lim_{t \rightarrow \infty} \frac{D(n, t)}{t} = \lambda(n-1)p(n-1) + \mu(n+1)p(n+1)$$

and

$$\lim_{t \rightarrow \infty} \frac{O(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{A(n, t)}{t} + \lim_{t \rightarrow \infty} \frac{D(n-1, t)}{t} = \lambda(n)p(n) + \mu(n)p(n)$$

we get that

$$\lambda(n-1)p(n-1) + \mu(n+1)p(n+1) = (\lambda(n) + \mu(n))p(n).$$

These equations hold for any $n \geq 0$ and are known as the *balance equations*. We will use these equations when studying queueing systems in which level-crossing cannot be used, for instance for queueing networks.

Again, just by using properties, i.e., counting differences, that hold along any sensible sample path we obtain very useful statistical and probabilistic results.

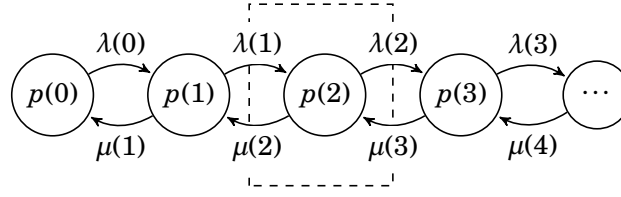


Figure 12: For the balance equations we count how often a box around a state is crossed from inside and outside. In the long run, the entering and leaving rates should be equal. For the example here, the rate out is $p(2)\lambda(2) + p(2)\mu(2)$ while the rate in is $p(1)\lambda(1) + p(3)\mu(3)$.

4.3 POISSON ARRIVALS SEE TIME AVERAGES

Suppose the following limit exists:

$$\pi(n) = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{k=1}^m \mathbb{1}_{L(A_k-) = n}, \quad (4.3.1)$$

then $\pi(n)$ is the long-run fraction of jobs that, upon arrival, observe n customers in the system. It is natural to ask whether $\pi(n)$ and $p(n)$, as defined by (4.2.2), are related, that is, whether what customers see upon arrival is related to the time-average behavior of the system. In this section we will derive the famous *Poisson arrivals see time averages (PASTA)* condition that ensures that $\pi(n) = p(n)$ if jobs arrive in accordance with a Poisson process.

Since $A(t) \rightarrow \infty$ as $t \rightarrow \infty$, it is reasonable that (see 4.3.5 for a proof)

$$\begin{aligned} \pi(n) &= \lim_{t \rightarrow \infty} \frac{1}{A(t)} \sum_{k=1}^{A(t)} \mathbb{1}_{L(A_k-) = n} = \lim_{t \rightarrow \infty} \frac{1}{A(t)} \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq t, L(A_k-) = n} \\ &= \lim_{t \rightarrow \infty} \frac{A(n, t)}{A(t)}, \end{aligned} \quad (4.3.2)$$

where we use (4.2.1a) in the last row. But, with (2.3.1),

$$\frac{A(n, t)}{t} = \frac{A(t)}{t} \frac{A(n, t)}{A(t)} \rightarrow \lambda \pi(n), \quad \text{as } t \rightarrow \infty, \quad (4.3.3)$$

while by (4.2.5),

$$\frac{A(n, t)}{t} = \frac{A(n, t)}{Y(n, t)} \frac{Y(n, t)}{t} \rightarrow \lambda(n) p(n), \quad \text{as } t \rightarrow \infty.$$

Thus

$$\lambda \pi(n) = \lambda(n) p(n). \quad (4.3.4)$$

This leads to our final result:

$$\lambda(n) = \lambda \iff \pi(n) = p(n).$$

This means that if the arrival rate does not depend on the state of the system, i.e., $\lambda(n) = \lambda$, the sample probabilities $\{\pi(n)\}$ are equal to the time-average probabilities $\{p(n)\}$. In other words, the customer perception at arrival moments is the same as the server perception.

As the next exercises show, this property is not satisfied in general. However, when the arrival process is Poisson we have that $\lambda(n) = \lambda$. This fact is called *PASTA: Poisson Arrivals See Time Averages*.

4.3.1. [4.3.1] Show for the case of 4.2.1 that $\pi(0) = 1$ and $\pi(n) = 0$, for $n > 0$.

4.3.2. [4.3.2] Check that (4.3.4) holds for the system of 4.3.1.

With the above reasoning, we can also establish a relation between $\pi(n)$ and the statistics of the system as obtained by the departures. Define, analogous to (4.3.2),

$$\delta(n) = \lim_{t \rightarrow \infty} \frac{D(n, t)}{D(t)} \quad (4.3.5)$$

as the long-run fraction of jobs that leave n jobs behind. From (4.2.4),

$$\frac{A(t)}{t} \frac{A(n, t)}{A(t)} = \frac{A(n, t)}{t} \approx \frac{D(n, t)}{t} = \frac{D(t)}{t} \frac{D(n, t)}{D(t)}.$$

Taking limits at the left and right, and using (2.3.2), we obtain for (queueing) systems in which customers arrive and leave as single units that

$$\lambda \pi(n) = \delta \delta(n). \quad (4.3.6)$$

Thus, if the system is rate-stable and transitions occur one-by-one, the statistics obtained by arrivals is the same as statistics obtained by departures, i.e.,

$$\lambda = \delta \iff \pi(n) = \delta(n). \quad (4.3.7)$$

4.3.3. [4.3.3] For the $G/G/1$ queue, prove that the fraction of jobs that see n jobs in the system upon arrival is the same as the fraction of departures that leave n jobs behind. What condition have you used to prove this?

4.3.4. [4.3.4] When $\lambda \neq \delta$, is $\pi(n) \geq \delta(n)$?

4.3.5. [4.3.7] There is a subtle problem in the transition from (4.3.1) to (4.3.2) and the derivation of (4.3.3): $\pi(n)$ is defined as a limit over arrival epochs while in $A(n, t)/t$ we take the limit over time. Now the observant reader might ask why these limits should relate at all. Use the renewal reward theorem to show that (4.3.2) is valid.

4.4 LITTLE'S LAW

There is an important relation between the average time $E[W]$ a job spends in the system and the long-run time-average number $E[L]$ of jobs that is contained in the system, which is called *Little's law*:

$$E[L] = \lambda E[W]. \quad (4.4.1)$$

Part of the usefulness of Little's law is that it applies under very general conditions to all input-output systems, whether the system is a queueing system or an inventory system or some much more general system. Hence, we will apply Little's law often in the forthcoming sections. The aim of this section is to prove this law.

We start by defining a few intuitively useful concepts. From (1.4.10), we see that

$$\frac{1}{t} \int_0^t L(s) ds = \frac{1}{t} \int_0^t (A(s) - D(s)) ds$$

is the time-average of the number of jobs in the system during $[0, t]$. Next, the waiting time of the k th job is the time between the moment the job arrives and departs, that is,

$$W_k = \int_0^\infty \mathbb{1}_{A_k \leq s < D_k} ds.$$

Fig. 4 relates W_k to $L(t)$.

Consider a departure time T at which the system is empty so that $A(T) = D(T)$. Then, for $k \leq A(T)$,

$$W_k = \int_0^T \mathbb{1}_{A_k \leq s < D_k} ds,$$

and for $s \leq T$,

$$L(s) = \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq s < D_k} = \sum_{k=1}^{A(T)} \mathbb{1}_{A_k \leq s < D_k}.$$

4.4.1. [4.4.3] Prove Little's law under the assumptions that $A(T_i) = D(T_i)$ for an infinite number of times $\{T_i\}$ such $T_i \rightarrow \infty$ and that all limits exist.

4.4.2. [4.4.7] For a given single-server queueing system the average number of customers in the system is $E[L] = 10$, customers arrive at rate $\lambda = 5$ per hour and are served at rate $\mu = 6$ per hour. Suppose that at the moment you join the system, the number of customers in the system is 10. What is your expected time in the system?

With the PASTA property and Little's law it becomes quite easy to derive expressions for the average queue length and waiting times for the $M/M/1$ queue. The average waiting time $E[W]$ in the entire system is the expected time in queue plus the expected time in service, i.e.,

$$E[W] = E[W_Q] + E[S]. \quad (4.4.2)$$

By the PASTA property we have for the $M/M/1$ queue that

$$E[W_Q] = E[L] E[S]. \quad (4.4.3)$$

4.4.3. [4.4.8] Use Little's law to show for the $M/M/1$ queue that

$$\begin{aligned} E[W] &= \frac{E[S]}{1 - \rho}, & E[L] &= \frac{\rho}{1 - \rho}, \\ E[L_Q] &= \frac{\rho^2}{1 - \rho}, & E[L_s] &= \rho. \end{aligned}$$

4.4.4. [4.4.9] Why is (4.4.3) not true in general for the $M/G/1$ queue?

In the above proof of Little's law we assumed that there is a sequence of moments $\{T_k, k = 0, 1, \dots\}$ at which the system is empty and such that $T_k < T_{k+1}$ and $T_k \rightarrow \infty$. However, in many practical queueing situations the system is never empty. Thus, to be able to apply Little's law to such more general situations we should slacken the assumption that such a sequence exists. The aim of this set of questions is to find an educated guess for a more general assumption under which Little's law can hold.

4.4.5. [4.4.10] Motivate in words (or with a derivation, if you prefer this) why the following is true:

$$\sum_{k=1}^{A(t)} W_k \geq \int_0^t L(s) ds \geq \sum_{k=1}^{D(t)} W_k.$$

4.4.6. [4.4.11] Take suitable limits (and assume all these limits exist) to show that

$$\lambda E[W] \geq E[L] \geq \delta E[W].$$

Make explicit all the points where you use the strong law of large numbers.

4.4.7. [4.4.12] Suppose that $A(t) = \lambda t$ and $D(t) = [A(t) - 10]^+$. Explain that for this system the above assumption on $\{T_k\}$ is violated. Show that Little's law is still true.

4.4.8. [4.4.13] Based on the above formulate an educated guess for more general conditions under which Little's law holds. (You don't have to prove Little's law under your condition; postpone that to after the exam.)

4.5 GRAPHICAL SUMMARY

We finish this chapter with providing two summaries in graphical form to clarify how all concepts developed in this chapter relate.

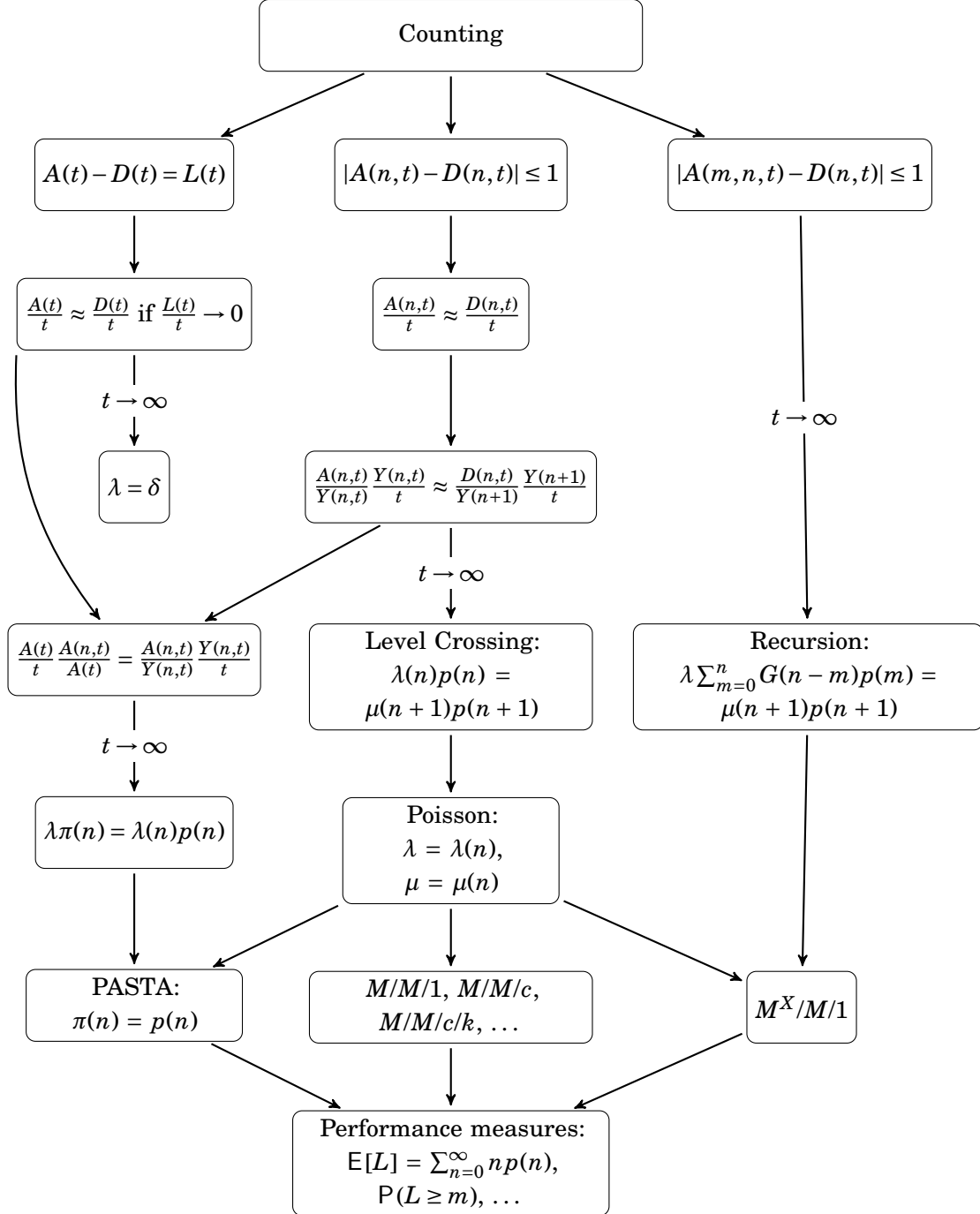


Figure 13: With level-crossing arguments we can derive a number of useful relations. This figure presents an overview of these relations that we derive in this and the next sections.

EXACT QUEUEING MODELS

In this chapter we use the concepts of Chapter 4 to model and analyze a large number of queueing systems in steady state. The simplest, non-trivial, case is the $M/M/1$ queue, which is the topic of Section 5.1. As the main ideas are based on sample-paths, it turns out to be nearly trivial to extend the analysis of the $M/M/1$ to the $M(n)/M(n)/1$ queue in Section 5.2. The $M(n)/M(n)/1$ queue is a very generic model with which we can find closed-form expressions for the queue length distributions for many other queueing models such as the $M/M/c$ queue or the $M/M/1/K$, i.e., a queueing system with loss,

We then focus on finding the expected queueing time for batch queues, in Section 5.3, and the $M/G/1$ queue, in Section 5.4. In the last two sections, Section 5.5 and Section 5.6, we derive expressions for the queue length distributions of the batch queue and the $M/G/1$ queue.

Many of the queueing systems we analyze here are either generalizations of some other model, for instance, the $M(n)/M(n)/1$ generalizes to the $M/M/c$ queue, or reduce to special cases in certain parameter settings, such as that the $M/G/1$ becomes the $M/M/1$ queue when the service times are exponentially distributed. Quite a number of exercises in this chapter are targeted on *checking* that the general results reduce to those of the special cases. The reader should understand the importance of such checks. These exercises are simple in a sense—it is perfectly clear what to do, there is no model to make for instance—, but the algebra can be quite tough at times, and hence it is good practice.

5.1 $M/M/1$ QUEUE

In the $M/M/1$ queue, one server serves jobs arriving with exponentially distributed inter-arrival times and each job requires an exponentially distributed processing time. With the level-crossing equations (4.2.7) we derive a number of important results for this queueing process.

Recall from Section 2.2 that we can construct the $M/M/1$ queue as a reflected random walk where the arrivals are generated by a Poisson process N_λ and the departures are generated according to the Poisson process N_μ (but only during periods at which the system is not empty!). Since the rates of these processes do not depend on the state of the random walk nor on the queue process, it follows that $\lambda(n) = \lambda$ for all $n \geq 0$ and $\mu(n) = \mu$ for all $n \geq 1$. Thus, (4.2.7) reduces to

$$p(n+1) = \frac{\lambda(n)}{\mu(n+1)} p(n) = \frac{\lambda}{\mu} p(n) = \rho p(n),$$

where we use the definition of the load $\rho = \lambda/\mu$. Since this holds for any $n \geq 0$, it follows with recursion that

$$p(n+1) = \rho^{n+1} p(0).$$

Then, by using normalization, it follows from (4.2.9) and (0.2.1d) that

$$p(0) = 1 - \rho, \quad p(n) = (1 - \rho)\rho^n. \quad (5.1.1)$$

It is now easy to compute the most important performance measures. The utilization of the server is $\rho = \lambda/\mu$, as observed above. Then, with a bit of algebra,

$$E[L] = \frac{\rho}{1-\rho}, \quad V[L] = \frac{\rho}{(1-\rho)^2}, \quad P(L > n) = \rho^{n+1}. \quad (5.1.2)$$

5.1.1. [5.1.10] *Derive (5.1.2) in three ways: by using indicator variables, by means of the expression $(1-\rho)^{-1} = \sum_{n=0}^{\infty} \rho^n$, and with moment-generating functions. (This is a good exercise to train your calculus and algebra skills.)*

Let us interpret (5.1.2). The fact that $E[L] \sim (1-\rho)^{-1}$ for $\rho \rightarrow 1$ implies that the average waiting time increases very fast when $\rho \rightarrow 1$. If we want to avoid long waiting times, this formula tells us that situations with $\rho \approx 1$ should be avoided. As a practical guideline, it is typically best to keep ρ quite a bit below 1, and accept that servers are not fully utilized.

Clearly, the probability that the queue length exceeds some threshold decreases geometrically fast (for $\rho < 1$). If we make the simple assumption that customers decide to leave (or rather, not join) the system when the queue is longer than 9 say, then $P(L \geq 10) = \rho^{10}$ is an estimator for the fraction of customers lost.

SUPERMARKET PLANNING Let us consider the example of cashier planning of a supermarket to demonstrate how to use the tools we developed up to now. Out of necessity, our approach is a bit heavy-handed—Turning the example into a practically useful scheme requires more sophisticated queueing models and data assembly—but the present example contains the essential analytic steps to solve the planning problem.

The *service objective* is to determine the minimal service capacity c (i.e., the number of cashiers) such that the fraction of the time that more than 10 people are in queue is less than 1%. (If the supermarket has 3 cashiers open, 10 people in queue means about 3 people per queue.)

The next step is to find the *relevant data*: the arrival process and the service time distribution. For the arrival process, it is reasonable to model it as a Poisson process. There are many potential customers, each choosing with a small probability to go to the supermarket at a certain moment in time. Thus, we only have to characterize the arrival rate. Estimating this for a supermarket is relatively easy: the cash registers track all customers payments. Thus, we know the number of customers that left the shop, hence entered the shop. (We neglect the time customers spend in the shop.) Based on these data we make a *demand profile*: the average number of customers arriving per hour, cf. Fig. 14. Then we model the arrival process as Poisson with an arrival rate that is constant during a certain hour as specified by the demand profile.

It is also easy to find the service distribution from the cash registers. The first item scanned after a payment determines the start of a new service, and the payment closes the service. (As there is always a bit of time between the payment and the start of a new service we might add 15 seconds, say, to any service.) To keep things simple here, we just model the service time distribution as exponential with a mean of 1.5 minutes.

We also *model* the behavior of all the cashiers together (a multi-server queue) as a single fast server. Thus, we neglect any differences between a station with, for instance, 3 cashiers and a single server that works 3 times as fast as a normal cashier. (We analyze in 5.2.3 the quality of this approximation.) As yet another simplification, we change the objective somewhat such that the number of jobs in the system, rather than the number in queue, should not exceed 10.

We now find a formula to convert the demand profile into the *load profile*, which is the minimal number of servers per hour needed to meet the service objective. We already know for



Figure 14: A demand profile of the arrival rate λ modeled as constant over each hour.

the $M/M/1$ that $P(L > 10) = \rho^{11}$. Combining this with the objective $P(L > 10) \leq 1\%$, we get that $\rho^{11} \leq 0.01$, which translates into $\rho \leq 0.67$. Using that $\rho = \lambda E[S]/c$ and our estimate $E[S] = 1.5$ minutes, we get the following rough bound on c :

$$c \geq \frac{\lambda E[S]}{0.67} \approx \frac{3}{2} \cdot \lambda \cdot 1.5 = 2.25\lambda,$$

where λ is the arrival rate (per minute, *not* per hour). For instance, for the hour from 12 to 13, we read in the demand profile in Fig. 14 that $\lambda = 120$ customers per hour, hence $c = 2.25 \cdot 120/60 = 4.5$. With this formula, the conversion of the demand profile to the load profile becomes trivial: divide the hourly arrival rate by 60 and multiply by 2.25.

The last step is to *cover the load profile with service shifts*. This is typically not easy since shifts have to satisfy all kinds of rules, such as: after 2 hours of work a cashier should take a break of at least 10 minutes; a shift length must be at least four hours, and no longer than 9 hours including breaks; when the shift is longer than 4 hours it needs to contain at least one break of 30 minutes; and so on. These shifts also have different costs: shifts with hours after 18h are more expensive per hour; when the supermarket covers traveling costs, short shifts have higher marginal traveling costs; and so on.

The usual way to solve such covering problems is by means of an integer problem. First, generate all (or a subset of the) allowed shift types with associated starting times. For instance, suppose only 4 shift plans are available

1. ++-++
2. +++-+
3. ++-+++
4. +++-++,

where a + indicates a working hour and – a break of an hour. Then generate shift types for each of these plans with starting times 8 am, 9 am, and so on, until the end of the day. Thus, a shift type is a shift plan that starts at a certain hour. Let x_i be the number of shifts of type i and c_i the cost of this type. Write $t \in s_i$ if hour t is covered by shift type i . Then the problem is to solve

$$\min \sum_i c_i x_i,$$

such that

$$\sum_i x_i \mathbb{1}_{t \in s_i} \geq 2.25 \frac{\lambda_t}{60}$$

for all hours t the shop is open and λ_t is the demand for hour t .

5.2 $M(n)/M(n)/1$ QUEUE

As it turns out, many more single-server queueing situations than the $M/M/1$ queue can be analyzed by making a judicious choice of $\lambda(n)$ and $\mu(n)$ in the level-crossing equations (4.2.7). For these queueing systems, we just present the results. The first set of exercises are somewhat technical; we ask you to derive the formulas—the main challenge is not to make computational errors. The second set of exercises shows how the combination of PASTA and Little's law allows us to analyze an astonishingly large number of non-trivial practical queueing situations.

It is important to realize that the inter-arrival times and service times need to be memoryless for the analysis below; the rates, however, may depend on the number of jobs in the system. Specifically, we require that for all s and t ,

$$P(A_{A(t)+1} \leq t + s \mid L(t) = n) = 1 - e^{-\lambda(n)s},$$

where we use that $A_{A(t)+1}$ is the arrival time of the next job after time t . Similarly, we assume for all t and s ,

$$P(D_{D(t)+1} \leq t + s \mid L(t) = n) = 1 - e^{-\mu(n)s}.$$

5.2.1. [5.2.1] Model the $M/M/1/K$ queue in terms of an $M(n)/M(n)/1$ queue and compute $p(K)$, i.e., the fraction of time that the system is full.

5.2.2. [5.2.3] Model the $M/M/c$ queue in terms of an $M(n)/M(n)/1$ queue and compute $E[L_Q]$.

5.2.3. [5.2.5] It should be clear that the $M/M/c$ queue is a bit harder to analyze than the $M/M/1$ queue, at least the expressions are more extensive. It is tempting to approximate the $M/M/c$ queue by an $M/M/1$ queue with a server that works c times as fast. As we now have the formulas for the $M/M/c$ queue and the $M/M/1$ queue we can use these to obtain some basic understanding of the difference.

Let us therefore consider a numerical example. Suppose that we have an $M/M/3$ queue, with arrival rate $\lambda = 5$ per day and $\mu = 2$ per server, and we compare it to an $M/M/1$ with the same arrival rate but with a service rate of $\mu = 3 \cdot 2 = 6$. Make a graph of the ratios of $E[L]$ and $E[L_Q]$ of both models as a function of ρ . Explain why these ratios become 1 as $\rho \uparrow 1$.

5.2.4. [5.2.6] Model the $M/M/c/c$ queue in terms of an $M(n)/M(n)/1$ queue and determine the performance measures. This model is also known as the Erlang B-formula and is often used to determine the number of beds at hospitals, where the beds act as servers and the patients as jobs.

5.2.5. [5.2.7] Take the limit $c \rightarrow \infty$ in the $M/M/c$ queue (or the $M/M/c/c$ queue) and obtain the performance measures for the $M/M/\infty$ queue, i.e., a queueing system with ample servers.

5.2.6. [5.2.12] Derive the steady state probabilities $p(n)$ for a single-server queue with a finite calling population with N jobs, i.e., jobs that are in service cannot arrive to the system. Check the answer you obtained for the cases $N = 1$ and $N = 2$. What happens if $N \rightarrow \infty$? Interpret the results.

5.2.7. [5.2.13] Give an example of a system with a finite calling population.

Finally, we consider queues with *balking*, that is, queues in which customers leave when they find the queue too long at the moment they arrive. A simple example model with customer balking is given by

$$\lambda(n) = \begin{cases} \lambda, & \text{if } n = 0, \\ \lambda/2, & \text{if } n = 1, \\ \lambda/4, & \text{if } n = 2, \\ 0, & \text{if } n > 2, \end{cases}$$

and $\mu(n) = \mu$.

Observe that here we make a subtle implicit assumption; in Section 4.3 we elaborate on this assumption. To make the problem clear, note that balking customers *decide at the moment they arrive* to either join or leave; in other words, they decide based on what they ‘see upon arrival’. In yet other words, they make decisions based on the state of the system at arrival moments, not on time-averages. However, the notion of $p(n)$ is a long-run *time-average*, and is typically not the same as what customers ‘see upon arrival’. As a consequence, the performance measure $P(L \leq n)$ is not necessarily in accordance with the perception of customers. To relate these two ‘views’, i.e., time-average versus observer-average, we need a new concept, *PASTA*, to be developed in Section 4.3.

5.2.8. [5.2.15] In what way is a queueing system with balking, at level b say, different from a queueing system with finite calling population of size b ?**5.2.9 (Hall 5.3). [5.2.19]**

After observing a queue with two servers for several days, the following steady-state probabilities have been determined: $p(0) = 0.4$, $p(1) = 0.3$, $p(2) = 0.2$, $p(3) = 0.05$ and $p(4) = 0.05$. The arrival rate is 10 customers per hour.

1. Determine $E[L]$ and $E[L_Q]$.
2. Using Little’s formula, determine $E[W]$ and $E[W_Q]$.
3. Determine $V[L]$ and $V[L_Q]$.
4. Determine the service time and the utilization.

5.2.10 (Hall 5.22). [5.2.22] At a large hotel, taxi cabs arrive at a rate of 15 per hour, and parties of riders arrive at the rate of 12 per hour. Whenever taxicabs are waiting, riders are served immediately upon arrival. Whenever riders are waiting, taxicabs are loaded immediately upon arrival. A maximum of three cabs can wait at a time (other cabs must go elsewhere).

1. Let p_{ij} be the steady-state probability of there being i parties of riders and j taxicabs waiting at the hotel. Write the state transition equation for the system.
2. Calculate the expected number of cabs waiting and the expected number of parties waiting.
3. Calculate the expected waiting time for cabs and the expected waiting time for parties. (For cabs, compute the average among those that do not go elsewhere.)
4. In words, what would be the impact of allowing four cabs to wait at a time?

5.2.11 (Continuation of **5.2.10**). **[5.2.24]** Did you have to use the PASTA property to solve **5.2.10**? If so, how did you use it? If not, why not?

5.2.12 (Continuation of **5.2.10**). **[5.2.25]** Suppose cabs can contain at most 4 riders, and the size of a party (i.e., a batch) has distribution B_k with $P(B_k = i) = 1/7$ for $i = 1, \dots, 7$. Parties of riders have the same destination, so riders of different parties cannot be served by one taxi. Provide a set of recursions to simulate this system. (This is a real hard exercise, but doable. I asked it at an exam to see who would deserve the highest grade. I was lenient with the grading...)

5.3 $M^X/M/1$ QUEUE: EXPECTED WAITING TIME

Sometimes jobs arrive in batches, rather than as single units. For instance, when a car or a bus arrives at a fast-food restaurant, a batch consists of the number of people in the vehicle. When the batches arrive as a Poisson process and the individual items within a batch have exponential service times we denote such queueing systems by the shorthand $M^X/M/1$. We derive expressions for the load and the expected waiting time and queue length for this queueing model.

Assume that jobs arrive as a Poisson process with rate λ and each job contains multiple items. Let A_k be the arrival time of job k and $A(t)$ the number of (job) arrivals up to time t . Denote by B_k the batch size of the k th job, i.e., the number of items of job k . We assume that $\{B_k\}$ is a sequence of independent discrete random variables each distributed as the generic random variable B . Let $f(k) = P(B = k)$ be given; let the survivor function be $G(k) = P(B > k)$. The service time of each item is $E[S]$.

5.3.1. [5.3.1] Explain that the average time to serve an entire batch is $E[B] E[S]$, so that the load must given by $\rho = \lambda E[B] E[S]$.

The aim of the remainder of the section is to derive a cornerstone of queueing theory, which is the following formula for the expected time an item spends in queue:

$$E[W_Q] = \frac{1 + C_s^2}{2} \frac{\rho}{1 - \rho} E[B] E[S] + \frac{1}{2} \frac{\rho}{1 - \rho} E[S], \quad (5.3.1)$$

where C_s^2 is the SCV of the batch size distribution. By applying (4.4.3), it follows right away that the expected number of items in the system takes the form

$$E[L] = \frac{E[W_Q]}{E[S]} = \frac{1 + C_s^2}{2} \frac{\rho}{1 - \rho} E[B] + \frac{1}{2} \frac{\rho}{1 - \rho}. \quad (5.3.2)$$

Note that $\rho < 1$ is required, as usual.

Before deriving the above, let us try to use it.

5.3.2. [5.3.3] If the batch size is geometrically distributed with success probability p , what is $E[L]$?

5.3.3. [5.3.4] A common operational problem is a machine that receives batches of various sizes. Management likes to know how a reduction of the variability of the batch sizes would affect the average queueing time. Suppose, for the sake of an example, that the batch size

$$P(B = 1) = P(B = 2) = P(B = 3) = \frac{1}{3}.$$

Batches arrive at rate 1 per hour. The average processing time for an item is 25 minutes. Compute by how much the number of items in the system would decrease if batch sizes were constant and equal to 2; hence the load is the same in both cases.

5.3.4. [5.3.6] Show that $E[W_Q(M^X/M/1)] \geq E[W_Q(M/M/1)]$ when the loads are the same. What do you conclude? (This solution of this exercise is more useful than you might think.)

Let us now focus on deriving (5.3.1). Assume that an arriving batch joins the end of the queue (if present), and once the queue in front of it has been cleared, it moves in its entirety to the server. Thus, all items in one batch spend the same time in queue. Once the batch moves to the server, the server processes the items one after another until the batch is empty. Write $E[L_Q^B]$ for the number of batches in queue and $E[L_S^B]$ for the number of items of the job (if any) at the server. Observe first that the average time an item spends in queue is

$$E[W_Q] = E[L] E[S] = \left(E[L_Q^B] E[B] + E[L_S^B] \right) E[S].$$

We also see that the average time a batch spends in queue is

$$E[W_Q^B] = E[L_Q^B] E[B] E[S] + E[L_S^B] E[S].$$

Hence, $E[W_Q] = E[W_Q^B]$.

5.3.5. [5.3.7] Use Little's law to show that

$$E[W_Q^B] = \frac{E[L_S^B]}{1-\rho} E[S], \text{ hence } E[L] = \frac{E[L_S^B]}{1-\rho}. \quad (5.3.3)$$

Clearly, we are done if we can find an expression for $E[L_S^B]$. For this we can use the renewal reward theorem; in fact, we can use 4.1.1 as inspiration. (Solve this exercise if you have not done yet.) Define $Y(t) = \int_0^t L_S^B(s) ds$ to see that $E[L_S^B] = Y = \lim_{t \rightarrow \infty} Y(t)/t$.

5.3.6. [5.3.9] Use the renewal reward theorem with sampling epochs $T_k = D_k$ to prove that

$$E[L_S^B] = \lambda \frac{E[B^2]}{2} E[S] + \frac{\rho}{2}.$$

5.3.7. [5.3.11] Use the above to derive (5.3.1).

5.4 M/G/1 QUEUE: EXPECTED WAITING TIME

In many practical single-server queueing systems the service times are not really well approximated by the exponential distribution. The M/G/1 queue then becomes a better model than the M/M/1 queue. In this section we first present a formula to compute the average waiting time in queue for the M/G/1 queue, and then we derive it by means of sample path arguments. The derivation is also of general interest as it develops some general results of renewal theory.

The fundamentally important *Pollaczek-Khinchine formula*, or *PK formula*, for the average waiting time in queue for the M/G/1 queue has the form

$$E[W_Q] = \frac{1+C_s^2}{2} \frac{\rho}{1-\rho} E[S]. \quad (5.4.1)$$

Before deriving this formula, let us apply it.

5.4.1. [5.4.4] A queueing system receives Poisson arrivals at the rate of 5 per hour. The single server has a uniform service time distribution, with a range of 4 minutes to 6 minutes. Determine $E[L_Q]$, $E[L]$, $E[W_Q]$, $E[W]$.

5.4.2. [5.4.5] Consider a workstation with just one machine. We model the job arrival process as a Poisson process with rate $\lambda = 3$ per day. The average service time $E[S] = 2$ hours, $C_s^2 = 1/2$, and the shop is open for 8 hours. What is $E[W_Q]$?

Suppose the expected waiting time has to be reduced to 1h. How to achieve this?

5.4.3 (Hall 5.16). [5.4.6] The manager of a small firm would like to determine which of two people to hire. One employee is fast, on average, but somewhat inconsistent. The other is a bit slower, but very consistent. The first has a mean service time of 2 minutes, with a standard deviation of 1 minute. The second has a mean service time of 2.1 minutes, with a standard deviation of 0.1 minutes. If the arrival rate is Poisson with rate 20 per hour, which employee would minimize $E[L_Q]$? Which would minimize $E[L]$?

To derive the PK-formula, suppose at first that we know the expected remaining service time $E[S_r]$, i.e., the expected time it takes to complete the job in service, if present, at the time a job arrives.

5.4.4. [5.4.12] Show that, given $E[S_r]$,

$$E[W_Q] = \frac{E[S_r]}{1 - \rho}. \quad (5.4.2)$$

It remains to compute the average remaining service time $E[S_r]$ for generally distributed service times. Just like in Section 5.3 we use the renewal reward theorem. Consider the k th job of some sample path of the $M/G/1$ queueing process. Let its service time start at time \tilde{A}_k so that it departs at time $D_k = \tilde{A}_k + S_k$.

5.4.5. [5.4.13] Use Fig. 15 to explain that the remaining service time of job k at time s is given by $(D_k - s) \mathbb{1}_{\tilde{A}_k \leq s < D_k}$. With this, explain that

$$Y(t) = \int_0^t (D_{D(s)+1} - s) \mathbb{1}_{L(s) > 0} ds$$

is the total remaining service time as seen by the server up to t .

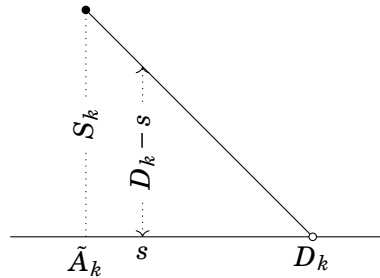


Figure 15: Remaining service time.

5.4.6. [5.4.14] Apply the renewal reward theorem to the result of 5.4.5 to prove that

$$E[S_r] = \frac{\lambda}{2} E[S^2]. \quad (5.4.3)$$

Then simplify to get (5.4.1).

5.4.7. [5.4.18] Show from (5.4.3) that

$$E[S_r | S_r > 0] = \frac{E[S^2]}{2E[S]}. \quad (5.4.4)$$

5.5 $M^X/M/1$ QUEUE LENGTH DISTRIBUTION

In Sections 5.3 and 5.4 we established the Pollaczek-Khinchine formula for the waiting times of the $M^X/M/1$ queue and the $M/G/1$ queue, respectively. To compute more difficult performance measures such as the loss probability $P(L > n)$, we need expressions for the stationary distribution $\pi(n)$. Here we present a numerical, recursive, scheme to compute these probabilities. Recall, that by the PASTA property, the long-run fraction of time $p(n)$ the system contains n jobs is the same as the fraction $\pi(n)$ of arriving jobs that observe n jobs in the system. Hence, by finding $\pi(n)$ we also obtain $p(n)$.

To find $\pi(n)$, $n = 0, 1, \dots$, we turn again to level-crossing arguments. However, the reasoning that led to the level-crossing equation (4.2.4) needs to be generalized. To see this, we consider an example. If $L(t) = 3$, the system contains 3 items. (This is not necessarily the same as 3 batches.) Since the server serves single items, down-crossings of level $n = 3$ occur in single units. However, due to the batch arrivals, when a job arrives it typically brings multiple items to the queue. For instance, suppose that $L(A_k-) = 3$, i.e., job k sees 3 items in the system at its arrival epoch. If its size $B_k = 20$, then right after the k th arrival the system contains 23 items, that is, $L(A_k) = 3 + 20 = 23$. Thus, upon the arrival of job k , all levels between states 3 and 23 are crossed.

The left panel in Fig. 16 shows all up- and down-crossings of some level n . The down-crossing rate is easy: just as in Fig. 10 there is just one arrow from right to left. However, level n can be up-crossed from below from many states, in fact from any level $m \in \{0, 1, \dots, n-1\}$. More formally, to count the number of up-crossings define

$$A(m, n, t) = \sum_{k=1}^{A(t)} \mathbb{1}_{L(A_k-) = m} \mathbb{1}_{B_k > n-m}$$

as the number of jobs up to time t that see m in the system upon arrival and have batch size larger than $n - m$. Recall from Section 5.3 that $f(k) = P(B = k)$ and $G(k) = P(B > k)$.

5.5.1. [5.5.4] Assuming that the limits exist, show that

$$\lim_{t \rightarrow \infty} \frac{A(m, n, t)}{t} = \lambda \pi(m) G(n - m).$$

Equating the number of up- and down-crossing gives $\sum_{m=0}^n A(m, n, t) \approx D(n, t)$. Then, dividing by t , taking the limit $t \rightarrow \infty$, and using 5.5.1 results in the level-crossing equation for the $M^X/M/1$ queue:

$$\lambda \sum_{m=0}^n \pi(m) G(n - m) = \mu \pi(n + 1). \quad (5.5.1)$$

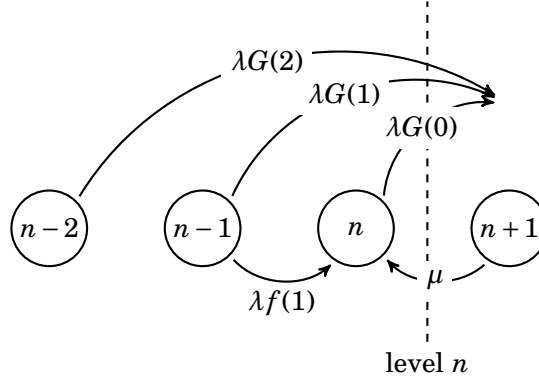


Figure 16: Level crossing of level n . Observe that when the system is in state $n-2$, the arrival of any batch larger than 2 ensures that level n is crossed from below. The rate at which such events happen is $\lambda\pi(n-2)G(2)$. Similarly, in state $n-1$, the arrival of any batch larger than one item ensures that level n is crossed, and this occurs with rate $\lambda\pi(n-1)G(1)$, and so on.

5.5.2. [5.5.5] Provide an interpretation of (5.5.1) in terms of a thinned Poisson arrival process.

It is left to find the normalization constant. As the recursion (5.5.1) does not lead to a closed form expression for $\pi(n)$, such as (5.1.1), we need to use a criterion to stop this iterative procedure. It is not so easy to find general conditions when to stop, but we can use a pragmatic approach. When the demand is finite, the numbers $\{\pi(k)\}$ should decrease geometrically fast for all $k \geq N$ where N is some large number¹. Stop when $\pi(N) \ll \pi(0)$, and take $\sum_{i=0}^N \pi(i)$ as the normalization constant.

Once we have $\pi(n)$, we can compute the influence on the batch size distribution, λ , and μ on the system's performance.

5.5.3. [5.5.7] Why is (4.3.7), i.e., $\pi(n) = \delta(n)$, not true for the $M^X/M/1$ batch queue? Provide an example.

5.5.4. [5.5.10] Substitute recursion (5.5.1) for $\pi(n)$ into the expression $E[L] = \sum_{n=0}^{\infty} n\pi(n)$ and derive (5.3.2).

5.5.5. [5.5.11] Implement the recursion (5.5.1) in a computer program for the case $f(1) = f(2) = f(3) = 1/3$. Take $\lambda = 1$ and $\mu = 3$.

We consider the $M^X/M/1/K$ queue, i.e., a batch queue in which at most K jobs fit into the system. When customers can be blocked in a batch queue it is necessary to specify a policy that decides which items in a batch to accept. Three common rules are

1. Complete rejection: if a batch does not fit entirely into the system, it will be rejected completely.
2. Partial acceptance: accept whatever fits of a batch, and reject the rest.
3. Complete acceptance: accept all batches that arrive when the system contains K or less jobs, and reject the entire batch otherwise.

5.5.6. [5.5.12] Derive a set of recursions, analogous to (5.5.1), to compute $\pi(n)$ for the $M^X/M/1/K$ queue with complete rejection.

¹ An interesting question, why should it decrease monotonically after some, large, N ?

5.6 M/G/1 QUEUE LENGTH DISTRIBUTION

In Section 5.5 we used level-crossing arguments to find a recursive method to compute the stationary distribution $p(n)$ of the number of items in an $M^X/M/1$ queue. Here we apply similar arguments to find $p(n) = P(L = n)$ for the $M/G/1$ queue. However, we cannot simply copy the derivation of the $M^X/M/1$ queue to the $M/G/1$ queue, because in the $M^X/M/1$ queue the service times of the items are exponential, hence memoryless, while in the $M/G/1$ this is not the case.

When job service times are not memoryless, hence do not restart at arrival times, we cannot choose any moment we like to apply level-crossing. Thus, for the $M/G/1$ queue we need to focus on moments in time in which the system ‘restarts’. As we will see below, the appropriate moments are job departure epochs. All in all, the argumentation to find the recursion for $\{p(n)\}$ is quite subtle, as it uses an interplay of the PASTA property and (4.3.7) between $\pi(n)$, $p(n)$ and $\delta(n)$.

An important role below is played by the number of arrivals Y_k during the service time of the k th job. Since the service times of the jobs form a sequence of i.i.d. random variables, the elements of the sequence $\{Y_k\}$ are also i.i.d. Let Y be the common random variable with probability mass $f(j) = P(Y = j)$; write $G(j) = P(Y_k > j)$ for the survivor function.

5.6.1. [5.6.2] Explain that

$$P(Y_k = j) = \int_0^\infty e^{-\lambda x} \frac{(\lambda x)^j}{j!} dF(x), \quad (5.6.1)$$

where F is the distribution of the service times.

5.6.2. [5.6.5] If $S \sim \text{Exp}(\mu)$, show that

$$G(j) = \sum_{k=j+1}^{\infty} f(k) = \left(\frac{\lambda}{\lambda + \mu} \right)^{j+1}. \quad (5.6.2)$$

5.6.3. [5.6.6] Design a suitable numerical method to evaluate (5.6.1) for more general distribution functions F .

5.6.4. [5.6.7] Show that $E[Y] = \lambda E[S] = \rho$.

Let us concentrate on a down-crossing of level n , see Fig. 17; recall that level n lies between states n and $n + 1$. For job k to generate a down-crossing of level n , two events must take place: job ‘ $k - 1$ ’ must leave $n + 1$ jobs behind after its service completion, and job k must leave n jobs behind. Thus,

$$\text{Down-crossing of level } n \iff \mathbb{1}_{L(D_{k-1})=n+1} \mathbb{1}_{L(D_k)=n} = 1.$$

Let us write this in another way. Observe that if $L(D_{k-1}) = n + 1$ and no other jobs arrive during the service time S_k of job k , i.e., when $Y_k = 0$, it must also be that job k leaves n jobs behind. If, however, $Y_k > 0$, then $L(D_k) \geq n + 1$. Thus, we see that

$$\text{Down-crossing of level } n \iff \mathbb{1}_{L(D_{k-1})=n+1} \mathbb{1}_{Y_k=0} = 1.$$

Consequently, the number of down-crossings of level n up to time t is

$$D(n + 1, 0, t) = \sum_{k=1}^{D(t)} \mathbb{1}_{L(D_{k-1})=n+1} \mathbb{1}_{Y_k=0}.$$

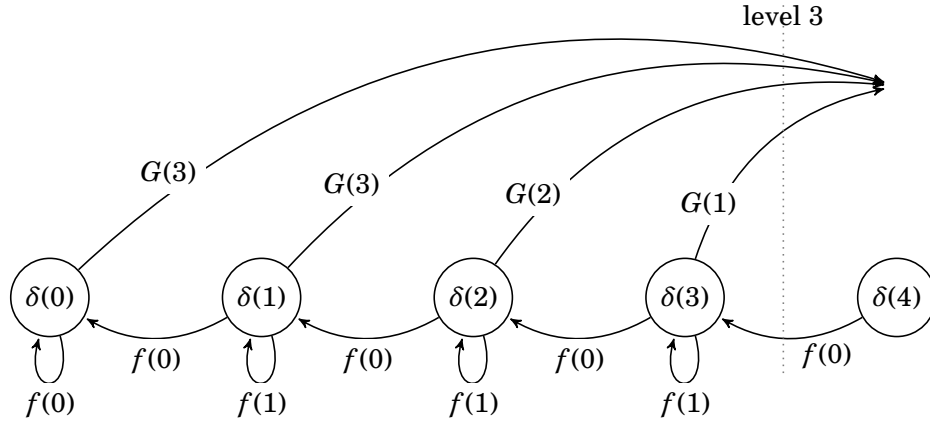


Figure 17: Level 3 is crossed from below with rate $\delta\delta(0)G(3) + \delta\delta(1)G(3) + \dots + \delta\delta(3)G(1)$ and crossed from above with rate $\delta\delta(4)f(0)$.

5.6.5. [5.6.8] Show that

$$\lim_{t \rightarrow \infty} \frac{D(n+1, 0, t)}{t} = \delta\delta(n+1)f(0),$$

where $f(0) = P(Y = 0)$.

Before we deal with the up-crossing, it is important to do the next exercise.

5.6.6. [5.6.9] Suppose that $L(D_{k-1}) > 0$. Why is $D_k = D_{k-1} + S_k$? However, if $L(D_{k-1}) = 0$, the time between D_{k-1} and D_k is not equal to S_k . Why not? Can you find an expression for the distribution of $D_k - D_{k-1}$ in case $L(D_{k-1}) = 0$?

For the up-crossings, assume first that $L(D_{k-1}) = n > 0$. Then an up-crossing of level $n > 0$ must have occurred when $L(D_k) > n$, i.e.,

$$\mathbb{1}_{L(D_{k-1})=n} \mathbb{1}_{L(D_k)>n} = 1 \implies \text{Up-crossing of level } n.$$

Again, we can convert this into a statement about the number of arrivals Y_k that occurred during the service time S_k of job k . If $Y_k = 0$, then job k must leave $n - 1$ jobs behind, so no up-crossing can happen. Next, if $Y_k = 1$, then job k leaves n jobs behind, so still no up-crossing occurs. In fact, level n can only be up-crossed from level n if more than one job arrives during the service of job k , i.e.,

$$\mathbb{1}_{L(D_{k-1})=n} \mathbb{1}_{Y_k>1} = 1 \implies \text{Up-crossing of level } n.$$

More generally, level n is up-crossed from level m , $0 < m \leq n$ whenever

$$\mathbb{1}_{L(D_{k-1})=m} \mathbb{1}_{Y_k>n-m+1} = 1 \implies \text{Up-crossing of level } n.$$

However, if $m = 0$ (think about this),

$$\mathbb{1}_{L(D_{k-1})=0} \mathbb{1}_{Y_k>n} = 1 \implies \text{Up-crossing of level } n.$$

Again we define proper counting functions, divide by t , and take suitable limits to find for the up-crossing rate

$$\delta\delta(0)G(n) + \delta \sum_{m=1}^n \delta(m)G(n-m+1). \quad (5.6.3)$$

Equating the down-crossing and up-crossing rates and dividing by δ gives

$$f(0)\delta(n+1) = \delta(0)G(n) + \sum_{m=1}^n \delta(m)G(n+1-m).$$

Noting that $\pi(n) = \delta(n)$, as this holds for any rate-stable $G/G/1$ queue, cf., (4.3.7), hence in particular for the $M/G/1$ queue length process, we arrive at

$$f(0)\pi(n+1) = \pi(0)G(n) + \sum_{m=1}^n \pi(m)G(n+1-m). \quad (5.6.4)$$

Clearly, we have again obtained a recursion with which we can compute the state probabilities.

5.6.7. [5.6.10] *Provide the details behind the derivation of (5.6.3).*

5.6.8. [5.6.11] *Clearly, the $M/M/1$ queue is a special case of the $M/G/1$ queue. Check that the queue length distribution of the $M/M/1$ queue satisfies (5.6.4).*

QUEUEING CONTROL AND OPEN NETWORKS

In the queueing systems we analyzed up to now, the server is always present to serve jobs in the system. However, this condition is not always satisfied. As an example, consider a queueing system in which there is a cost associated with switching on and off the server. For instance, in some cases the server has to be set-up for operation; in other cases, the operator of a machine has to move from one place in the factory to another. To reduce the cost, the so-called N -policy of 1.2.5 can be used. Recall that this policy works as follows. As soon as the system becomes empty (and the server idle), we switch off the server. Then we wait until N or more jobs have arrived, and then we switch on the server. The server processes jobs until the system is empty again, switches off, and remains idle until a sufficient number of new jobs have arrived, and so on. Thus, we use an N -policy to *control* the queueing system, in particular the server, and the task is to find a switching threshold N that minimizes the long-run average cost.

Observe that under such policies the server also has longer busy and idle times. In fact, this sometimes seems to be the policy at dentists or hospitals: do something else until the waiting room is quite full, and then start serving patients. Like this, in the example of a GP, the server (GP) does not have to wait for short times for patients that might be late, but instead can collect idle times into one long stretch, and do something useful instead.

In this chapter we first study the $M/M/1$ queue and $M/G/1$ under a N -policy. As a second topic we study open networks of stations of $M/M/c$ stations and jobs arrive as a Poisson process. As we will see, the analysis of the N -policy requires to solve an equation of the type $v = c + Pv$, where v and c are vectors and P a (stochastic) matrix, while for the network we need to solve an equation of the type $\lambda = \gamma + \lambda P$, where λ and γ are vectors and P is again a (stochastic) matrix. As these equations are (nearly) the same, we therefore concentrate in Section 6.4 on the solution. The analysis in this chapter allows us to illustrate many tools and results of the previous chapters such as the renewal-reward theorem. In a sense, everything comes together here.

We finally point out that the techniques developed in this chapter extend (way) beyond just queueing theory; they are worth memorizing. The concepts we introduce here can for instance be generalized to (optimal) stopping problems, which find many applications beyond queueing, such in finance, inventory theory, decision theory, and so on. As another set of extensions, it is possible to make the matrix P and the vector c depend on an action one can take in certain states. This idea underlies Markov decision theory, which in turn provides the theoretical basis of a number of machine learning tools such as Q learning, reinforcement learning, and so on. Thus, while this chapter closes our journey on the study of queueing system, it is a first step toward a much longer journey on the diverse applications of the probability theory.

6.1 N -POLICIES FOR THE $M/M/1$ QUEUE

Let us consider the $M/M/1$ queue in which the server switches off as soon as it becomes idle and it costs K to switch on. Supposing that each job charges h Euros per unit time while in the system, it makes sense to build up a queue of jobs while the server is idle, and after some

time switch on the server to process jobs until the system is empty again. In particular, we analyze the influence of the N -policy on the long-run average cost. For this we make a cost model in several steps and at the end we discuss how to minimize the cost as a function of the threshold N . In passing, we obtain a third way to compute the time-average number $E[L]$ of jobs in the system; the first resulted from the analysis of the $M/M/1$ queue in Section 5.1, the second from Little's law, cf. 4.4.3

Suppose the server is on, and there are q jobs in the system. Let us write $T(q)$ for the expected time to clear the system. Now one of two events happens first. Either a new job enters the system, or the job in service leaves. It follows from 1.3.9 that $\alpha = \lambda/(\lambda + \mu)$ is the probability the first event occurs and $\beta = \mu/(\lambda + \mu)$ is the probability the second occurs. Moreover, from 1.3.8 we see that the expected time to either an arrival or a departure, whichever is first, is $1/(\mu + \lambda)$. Therefore, $T(q)$ must satisfy the following recursion:

$$T(q) = \alpha T(q+1) + \beta T(q-1) + \frac{1}{\lambda + \mu}. \quad (6.1.1)$$

The reader should note that this type of recursion is a difference equation. Moreover, the ideas behind its derivation are very similar to the ideas used in dynamic programming.

6.1.1. [6.1.1] *Provide an intuitive explanation for (6.1.1) (where is the memoryless property used?), and show that $T(q) = q/(\mu - \lambda)$ solves (6.1.1).*

With the same line of reasoning we can compute the expected cost $V(q)$ to clear the system. Noting that the queueing cost is hq per unit time when there are q jobs in the system, $V(q)$ must satisfy the relation

$$V(q) = \alpha V(q+1) + \beta V(q-1) + h \frac{q}{\lambda + \mu}. \quad (6.1.2)$$

To solve this, observe first that in (6.1.1), which we use to compute $T(q)$, the last term is a constant, and that $T(q)$ is a linear function in q . In the recursion for $V(q)$ we see that the last term is linear in q . So let us guess that $V(q)$ is quadratic in q , i.e., $V(q) = aq^2 + bq + c$. Thus, we substitute this into the above expression and then try to solve for a, b and c .¹ As $V(0) = 0$, it follows already that $c = 0$. Thus, it remains to find a and b .

6.1.2. [6.1.2] *Use the ideas of 6.1.1 to show that*

$$V(q) = aq^2 + bq = \frac{h}{2} \frac{1}{\mu - \lambda} q^2 + \frac{h}{2} \frac{\lambda + \mu}{(\mu - \lambda)^2} q.$$

Interestingly, the above turns out to be immediately useful. Suppose we switch on the server when $N = 1$, that is, directly at the arrival of the first job after the system became empty. Write $C(1)$ for the expected duration of a cycle that starts when the system becomes idle and stops when the system becomes idle again (and after at least one job has arrived). In other words, the cycle consists of an idle and a busy period.

6.1.3. [6.1.3] *Explain that $C(1) = 1/\lambda + T(1)$ for the $M/M/1$ queue.*

¹ The reader might wonder about the uniqueness of the solution. Noting that $V(q+1)$ follows directly from $V(q)$ and $V(q-1)$, there can be just one solution when the boundary conditions are given.

6.1.4. [6.1.4] Use the renewal-reward theorem to explain the relation

$$\frac{V(1)}{C(1)} = \frac{V(1)}{1/\lambda + T(1)} = h E[L],$$

where $E[L]$ is the expected number of jobs in an M/M/1 queue and given by (5.1.2). Then use the above expressions for $V(1)$ and $C(1)$ to verify this.

With this result, we have found yet another way to compute the expected number of jobs in the system.

It remains to generalize to a general threshold N ; just above we already covered the case with $N = 1$. As we already have expressions for the cost and time for the time the server is on, we only have to consider the cost and time while the server is off. Note that right after the server switches off, we need N independent inter-arrival times to reach level N . By 1.3.5, the expected time to switch on must be equal to N/λ .

For the cost while building up the queue, we use again a recursive procedure. Write $W(q)$ for the accumulated queueing cost from the moment the server becomes idle up to the arrival time of the q th job (the job that sees $q - 1$ jobs in the system). Then, by the next exercise,

$$W(q) = W(q - 1) + h \frac{q - 1}{\lambda} = h \frac{q(q - 1)}{2\lambda}. \quad (6.1.3)$$

6.1.5. [6.1.5] Explain the above recursion, and derive the right-hand side.

It remains to assemble all results. Let us assume that the switching cost is K . Then, by the renewal-reward theorem, the time-average cost of the N -policy is equal to

$$\frac{W(N) + K + V(N)}{C(N)},$$

where $C(N) = N/\lambda + T(N)$.

Finding the optimal N is easy (from a practical point). Observe that $V(N)$ and $W(N)$ are quadratic in N , while $C(N)$ is linear in N . Hence, the average cost is a convex function of N , and locating the minimizer of a one-dimensional convex function is simple.

6.2 N-POLICIES FOR THE M/G/1 QUEUE

Interestingly, we can extend the analysis of Section 5.6 and Section 6.1 to compute the average cost of the M/G/1 queue under an N -policy. Thus, rather than exponentially distributed service times, we now consider general service times. For the rest, the model and the problem is the same as in Section 6.1.

Similar to Section 5.6, we consider the M/G/1 queueing process at departure moments. It is easy to find an expression for the time to clear the queue right after the departure of a job that leaves q jobs behind. Analogous to (6.1.1) and using the definition of $f(k) = P(Y = k)$ as the probability that $Y = k$ jobs arrive during an arbitrary service time, we see that $T(q)$ must satisfy the relation

$$T(q) = \sum_{k=0}^{\infty} f(k)T(q + k - 1) + E[S].$$

To see this, observe that when $k = 0$ the next level is $q - 1$ so that it takes $T(q - 1)$ time to hit level 0. When $k = 1$, the level returns to q after the departure of the job currently in service; consequently, we have not made any progress, hence we still need $T(q)$ units of time, and so on. Similar to the reasoning in Section 6.1, we guess that $T(q) = aq + b$. Since $T(0) = 0$, we already have $b = 0$.

6.2.1. [6.2.1] Substitute $T(q) = aq$ in the above recursion for $T(q)$, and solve for a to obtain

$$T(q) = \frac{E[S]}{1 - E[Y]} q = \frac{E[S]}{1 - \lambda E[S]} q.$$

Let $V(q)$ be the cost to clear the system right after a departure that leaves q jobs behind. Thus, $V(0) = 0$. Again, analogous Section 5.6 and the derivation above for $T(q)$, we see that $V(q)$ must satisfy the relation

$$V(q) = \sum_{k=0}^{\infty} f(k)V(q+k-1) + H(q), \quad (6.2.1)$$

where $H(q)$ is the queueing, or holding, cost, which we will determine next.

The queueing cost $H(q)$ consists of two components. The first is the cost to keep q jobs in the system while there is a job in service. Clearly, the expected cost of this is $hqE[S]$. The second component is the cost of new jobs that arrive during the service. While this is slightly harder to determine, we can combine the ideas underlying the derivation in (6.1.3) and 1.1.3 and 1.3.9. Specifically, assume it is given that the service time is s . Then the expected total queueing cost $U(s)$ of the new arrivals must satisfy for some $0 < \delta \ll 1$

$$U(s) = U(s - \delta) + (1 - \lambda\delta) \cdot 0 + \lambda\delta hs + o(\delta).$$

This follows because during the interval $[0, \delta]$ a job arrives with probability $\lambda\delta$ and stays s time units in the system, while with probability $1 - \lambda\delta$ no job arrives, and the cost of this is 0. By the next problem,

$$U(s) = \lambda h \frac{s^2}{2}.$$

6.2.2. [6.2.2] Assuming that $U(\cdot)$ is differentiable, we can use that $(U(s) - U(s - \delta))/\delta \approx U'(s)$. Use this in the above to arrive at the differential equation $U'(s) = \lambda hs$ that U must satisfy with condition $U(0) = 0$ (why is this so?). Solve this DE to conclude that $U(s)$ is as given above.

Since $U(s)$ is the expected total cost given that the service time is s , it follows from 5.6.4 that $E[U(S)] = \lambda h E[S^2]/2$ is the expected queueing cost of the new jobs that arrive during the service. By combining the first and second component of $H(q)$, we obtain

$$H(q) = hqE[S] + \frac{1}{2}\lambda h E[S^2].$$

To solve (6.2.1), note that as in (6.1.2), the cost $H(q)$ has a term linear in q and a constant term. Thus, once again, we guess that $V(q)$ is a quadratic function in q . Then we substitute this form into (6.2.1), assemble terms with the same power in q , and try to solve for the coefficients in $aq^2 + bq + c$. If we succeed, the solution we found must be correct, as a difference equation of the type (6.2.1) can have just one solution once the boundary conditions are given. Again, $V(q) = 0 \implies c = 0$. In fact, observe that we have all elements to compute the costs; except for (quite a bit of) algebra, we are done!

6.2.3. [6.2.6] Show that

$$a = \frac{h E[S]}{2(1 - \rho)}, \quad b = \frac{h E[S]}{2(1 - \rho)^2} (1 + \rho C_s^2).$$

(In a sense, this is trivial, as it is just algebra, but it is hard to get the details right.)

6.2.4. [6.2.8] As in Section 6.1, show that we can obtain the Pollaczek-Khinchine equation (5.4.1) with the results we have obtained up to now.

For the $M/G/1$ queue it is evident that the expected queueing cost while the server is idle is also given by (6.1.3). With this, and noting that there is a cost K to switch on the server, it follows from 6.2.5 that the long-run time-average average costs are equal to

$$\frac{V(N) + K + W(N)}{C(N)} = \frac{h}{2} \frac{\rho^2}{1 - \rho} (1 + C_s^2) + h\rho + h \frac{N-1}{2} + K \frac{\lambda(1-\rho)}{N}. \quad (6.2.2)$$

Note that, if $N = 1$ and $K = 0$, this reduces to $hE[L]$, as it should by our work above. Finally, minimizing over N gives that

$$N^* \approx \sqrt{\frac{2\lambda(1-\rho)K}{h}}.$$

Remark 6.2.1. The expression for N^* is a famous result in inventory theory: it is the optimal order size for a machine that produces items with holding cost h per item per unit time, a cost K to switch on the machine, demand arrives at rate λ , and the machine produces at rate $\mu = 1/E[S]$. N^* is known as the Economic Production Quantity (EPQ). Taking $\mu \rightarrow \infty$, it reduces to the Economic Order Quantity (EOQ).

6.2.5. [6.2.9] Derive (6.2.2).

6.3 OPEN SINGLE-CLASS PRODUCT-FORM NETWORKS

Up to now our analysis focused on single-station queueing systems. In many practical situations, however, jobs in a factory or patients in a hospital have to undergo several process steps before they are ‘finished’. One of the simplest models to analyze such situations is to assume that jobs arrive as a Poisson processes, and the service times are exponentially distributed. We will see that it is possible to obtain closed-form expressions for the stationary distribution of jobs at each station. To establish this we will first concentrate on two stations in tandem, and then extend to general networks. We remark that in Section 3.5 we considered tandem networks of $G/G/c$ queues, but there we could only obtain insight in the expected times, not the full distribution of the number of jobs at each station.

Theory and Exercises

Tandem Queues

In 6.3.1 (and the intermediate exercises leading to that result) we (ask you to) prove that the inter-departure times of an $M/M/1$ queue are also exponentially distributed with rate λ . This is useful because when the first station is an $M/M/1$ queue, this implies that the arrival process at the second station is also an $M/M/1$ queue. Stated differently, from the perspective of the second station, it is as if there is no first station. And, if there is a third station with exponentially distributed service times, this also behaves as an $M/M/1$ queue, and so on. With this insight, it is easy to see that the average total waiting time in a tandem network of M stations equals

$$E[W] = \sum_{i=1}^M \frac{E[S_i]}{1 - \rho_i},$$

where $\rho_i = \lambda E[S_i]$, and $E[S_i]/(1 - \rho_i)$ is the expected waiting time at station i ; note that the arrival rate at each station is λ , due to the topology of the network, i.e., a tandem network.

6.3.1. [6.3.8] Assuming that inter-departure times are independent, prove Burke's law which states that the departure process of the $M/M/1$ queue is a Poisson process with rate λ .

Open Networks of $M/M/1$ queues

It is not difficult to extend the above result for tandem networks to general networks of $M/M/1$ queues. For this, we first need to model such networks more formally. In particular, we assume that the probability that a job moves to station j after completing its service at station i is independent of anything else, and is given by the number $P_{ij} \in [0, 1]$. (This is called Markov routing.) We assemble all these probabilities in a *routing matrix* P such that P_{ij} is the element of P on the i th row and j th column. We require that $\sum_{j=1}^M P_{ij} \in [0, 1]$ for each row i .

6.3.2. [6.3.9] Why do we require this? What is the interpretation of $P_{i0} = 1 - \sum_{j=1}^M P_{ij}$? What does it mean when $P_{i0} > 0$?

Consider station i , say, and assume that jobs arrive as a Poisson process with rate λ_i . Since service times are exponentially distributed, it follows from the previous section that the departure process is also Poisson with rate λ_i . Then, after departure, jobs are sent with probability P_{ij} to station j , independent of anything else. But then we can use 1.1.9 to conclude that the jobs sent to station j form a Poisson processes with rate $\lambda_i P_{ij}$. Now take the perspective of some station j . Suppose this station receives such thinned Poisson 'streams' from all other stations. Then observe that, by 1.1.8, this merged process is also a Poisson process with the combined rate of the individual 'streams'. Assuming that new jobs arrive at station j as a Poisson process with rate γ_j , we can merge this process with the departure processes of the other stations to obtain the total arrival process at station j , and this must again be Poisson and has rate

$$\gamma_j + \sum_{i=1}^M \lambda_i P_{ij}.$$

Finally, since jobs arrive at station j as a Poisson process, and service times are exponential, the departure process of this station is also Poisson, and so on for all the stations in the network. Thus, it is intuitively clear that we can model this network as a set of $M/M/1$ queues; below we will give a formal proof of this fact for two stations. Note that we allow for external jobs arriving at any station. Therefore this network is *open*. This differs from so-called *emph* closed networks; in such networks jobs do not enter or leave.

It is evident that, when the network is stable (so that queues do not keep on increasing over time), all jobs that enter the network must eventually leave. This insight leads us to the *traffic (rate) equations*, which state that for all stations i ,

$$\lambda_i = \gamma_i + \sum_{j=1}^M \lambda_j P_{ji}, \quad i = 1, \dots, M, \quad (6.3.1)$$

where the left-hand side represents the rate at which jobs depart and the right-hand side the rate at which jobs arrive.

Let us for the moment assume that we can solve the traffic equations, in other words, we can find a set of numbers $\lambda = (\lambda_1, \dots, \lambda_M)$ such that (6.3.1) is satisfied for given $\gamma = (\gamma_1, \dots, \gamma_M)$ and routing matrix P , cf. Section 6.4. Then, we can define the load at station i as $\rho_i = \lambda_i E[S_i]$.

Clearly, we assume that $\rho_i < 1$ for all stations i . Moreover, station i being an $M/M/1$ queue, it follows that the sojourn time is $E[W_i] = E[S_i]/(1 - \rho_i)$. Write $|\gamma| = \sum_{i=1}^M \gamma_i$ as the total external arrival rate. Then, using that the average total number of jobs is equal to the sum of the average number of jobs at each station, it must follow that

$$E[L] = \sum_{i=1}^M E[L_i].$$

Then with an application of Little's law to the network as a whole and to each station individually, we get

$$|\gamma| E[W] = E[L] = \sum_{i=1}^M E[L_i] = \sum_{i=1}^M \lambda_i E[W_i].$$

As a last step, by dividing by $|\gamma|$, we can express the average sojourn time in the network in terms of the *visiting ratios* $\lambda_i/|\gamma|$ as

$$E[W] = \sum_{i=1}^M \frac{\lambda_i}{|\gamma|} E[W_i].$$

6.3.3. *Provide an interpretation for the above expression.*

6.3.4. [6.3.11] *We have a two-station single-server open network. Jobs enter the network at the first station with rate γ . A fraction α returns from station 1 to itself; the rest moves to station 2. At station 2 a fraction β_2 returns to station 2 again, a fraction β_1 goes to station 1.*

First, compute λ , then analyze what happens if $\alpha \rightarrow 1$ or $\beta_1 \rightarrow 0$.

Stationary distributions

Above we derived expressions for the average waiting time in a network of $M/M/1$ queues. In fact, it is possible to obtain the much stronger result that the stationary probability that the system contains $n = (n_1, n_2, \dots, n_M)$ at stations $1, \dots, M$ takes the form

$$P(N_1 = n_1, \dots, N_M = n_M) = p(n) = \prod_{i=1}^M p(n_i) = \prod_{i=1}^M (1 - \rho_i) \rho_i^{n_i},$$

where $p(n_i) = (1 - \rho_i) \rho_i^{n_i}$ is the stationary probability that station i contains n_i jobs, compare (5.1.1). In words, $p(n)$ is equal to the product of the probabilities $p(n_i)$ of all of the stations $i = 1, \dots, M$. But, this implies that $p(n_i)$ is *independent* of the state of the other stations. For notational ease, we will prove this for the case of two stations in tandem. The general result is known as the fact that *Jackson networks*, i.e., open networks of $M/M/c$ networks, admit a *product-form solution*. Note that these probabilities are useful to estimate excess probabilities such as $P(L_1 > n_1, L_2 > n_2)$.

Write $p(i, j) = P(N_1 = i, N_2 = j)$ for the state of the two-station network to denote that station 1 contains i jobs and station 2 contains j jobs. We have to show that the probabilities $p(i, j)$ satisfy the balance equations for all $i, j \geq 0$. Recall that the balance equations express that, in steady state, the total (probability) rate out of a state must be equal to the (probability) rate into this state, cf., Section 4.2.

6.3.5. *Provide the balance equations for states (i, j) with $i, j \geq 0$ and check that these are satisfied by $p(i, j)$ for $i, j \geq 0$.*

6.4 ON $\lambda = \gamma + \lambda P$

Here we study the existence of a solution for the equation $\lambda = \gamma + \lambda P$, and its inverse $V = c + PV$.

Theory and Exercises

From 6.3.2 we know that $P_{i0} = 1 - \sum_{j=1}^M P_{ij}$ is the probability that a job departing from node i also leaves the network, in other words, with probability P_{i0} a job is finished after its service at station i . Next, consider a station k with $P_{ki} > 0$. Then the probability that a job starting at k , moving to i and then leaving the network, must be equal to $P_{ki}P_{i0}$. As $P_{ki} > 0$ and $P_{i0} > 0$, the probability that a job leaves the network from node k in two steps is positive. As a matter of fact, $P_{k0}^2 = \sum_{j=0}^M P_{kj}P_{j0} \geq P_{ki}P_{i0} > 0$. More generally, we assume that the (finite) matrix P is transient, which means that it is possible to leave the network from any station in at most M steps. In other words, for any station j there is a sequence of intermediate stations j_1, j_2, \dots, j_{M-1} such that $P_{j0}^M \geq P_{jj_1}P_{j_1j_2} \cdots P_{j_{M-1}0} > 0$.

The next few exercises provide two different ways to prove that a finite transient matrix $P^n \rightarrow 0$ geometrically fast ((element wise) as $n \rightarrow \infty$).

6.4.1. [6.4.2] Use that $\sum_{j=1}^M P_{ij}^M < 1$ for all $i = 1, \dots, M$ to prove that $P^n \rightarrow 0$ geometrically fast in n .

Before we can provide the second type of proof, do the next exercise.

6.4.2. [6.4.4] What is the geometric interpretation of an eigenvector and eigenvalue of a matrix P , say? Specifically, what happens if an eigenvalue has modulus less than 1?

Let us make the simplifying assumption that P is a diagonalizable matrix with M different eigenvalues.² In this case, there exists an invertible matrix V with the (left) eigenvectors of P as its rows and a diagonal matrix Λ with the eigenvalues on its diagonal such that $VP = \Lambda V$. Hence, premultiplying with V^{-1} , $P = V^{-1}\Lambda V$. But then $P^2 = V^{-1}\Lambda V \cdot V^{-1}\Lambda V = V^{-1}\Lambda^2 V$, and in general $P^n = V^{-1}\Lambda^n V$. Clearly, if each eigenvalue λ_i is such that its modulus $|\lambda_i| < 1$, then $\Lambda^n \rightarrow 0$ geometrically fast, hence $P^n \rightarrow 0$ geometrically fast.

So, let us prove that all eigenvalues of a finite, transient routing matrix P have modulus less than 1. For this we use *Gerschgorin's disk theorem*. Define the Gerschgorin disk of the i th row of the matrix P as the disk in the complex plane:

$$B_i = \left\{ z \in \mathbb{C}; |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| \right\}.$$

In words, this is the set of complex numbers that lies within a distance $\sum_{j \neq i} |a_{ij}|$ of the point a_{ii} . Next, assume for notational simplicity that for each row i of P we have that $\sum_j a_{ij} < 1$ (otherwise apply the argument to P^M .) Then this implies for all i that

$$1 > \sum_{j=1}^M a_{ij} = a_{ii} + \sum_{j \neq i} a_{ij}.$$

Since all elements of P are non-negative, so that $|a_{ij}| = a_{ij}$, it follows that

$$-1 < a_{ii} - \sum_{j \neq i} a_{ij} \leq a_{ii} + \sum_{j \neq i} a_{ij} < 1.$$

² The argument below applies just as well matrices reduced to Jordan normal form, but only adds notational clutter.

With this and using that a_{ii} is a real number (so that it lies on the real number axis) it follows that the disk B_i lies strictly within the complex unit circle $\{z \in \mathbb{C}; |z| \leq 1\}$. As this applies to any row i , the union of the disks $\cup_i B_i$ lies strictly within the complex unit circle. Now Gerschgorin's theorem states that all eigenvalues of the matrix P must lie in $\cup_i B_i$. We conclude that all eigenvalues of P also lie strictly in the unit circle, hence all eigenvalues have modulus smaller than 1.

With the above results, we can show that the equation $\lambda = \gamma + \lambda P$ has a unique solution when P is a transient matrix. For this, define iteratively,

$$\lambda^0 = 0, \quad \lambda^n = \gamma + \lambda^{n-1}P, \quad \text{for } n \geq 1.$$

Then, by substituting $\lambda^j = \gamma + \lambda^{j-1}P$ a sufficient number of times,

$$\lambda^n = \gamma + \lambda^{n-1}P = \gamma + (\gamma + \lambda^{n-2}P)P = \gamma \sum_{i=0}^{n-1} P^i,$$

where we take $P^0 = 1$, i.e., equal to the identity matrix. By the result of the above reasoning, there exists an N and $\epsilon > 0$ such that $P_{ij}^n < (1 - \epsilon)^n$ for all $n > N$ and $1 \leq i, j \leq M$. Therefore, and using (0.2.1d), each element i, j of the sequence of matrices $\sum_{k=0}^n P^k$ increases monotonically as $n \rightarrow \infty$ to a finite limit. Consequently,

$$\lambda = \gamma \sum_{k=0}^{\infty} P^k \tag{6.4.1}$$

is well-defined, finite, and the solution of $\lambda = \gamma + \lambda P$.

We can apply the results here also to see why the recursions for T , V and W in Section 6.1 and Section 6.2 have solutions. For instance, observe that we can write (6.2.1) as

$$V = PV + H,$$

with

$$P = \begin{pmatrix} 0 & 0 & 0 & 0 & \dots \\ f(0) & f(1) & f(2) & 0 & \dots \\ 0 & f(0) & f(1) & f(2) & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots \end{pmatrix}, \quad H = \begin{pmatrix} 0 \\ H(1) \\ H(2) \\ \vdots \end{pmatrix}.$$

It is clear that $V = PV + H$ is the same equation (in transpose) as $\lambda = \gamma + \lambda P$. With the same type of reasoning, we can find the solution for V , for instance by defining iteratively $V^0 = 0$, and $V^n = PV^{n-1} + H$, for $n \geq 0$. Again, P is here a non-negative matrix, although it is infinite, which adds a few technical complications. These complications can be solved, and with this body of theory we can analyze all such systems.

We finish our discussion of queueing systems, but there many other interesting extensions to learn. Here are some nice references that are now accessible to you.

- You can find really nice discussion of networks of $M/M/\infty$, chemical reactions, population dynamics and Petri nets in [Baez and Biamonte \[2019\]](#), which is freely available on arXiv.
- Simple queueing networks (networks that satisfy so-called local balance) can be modeled as electrical networks. For this, see [Doyle and Laurie Snell \[1984\]](#), which you can download for free from the homepage of Doyle.

- In more general terms, queueing systems or networks are examples of Markov processes. A particularly nice book on these topics is [Norris \[1997\]](#). The material of this chapter can be couched in the theory of martingales and optimal stopping. Besides that this is nice theory, this is widely used in quantitative finance.

BIBLIOGRAPHY

- J.C. Baez and J. Biamonte. Quantum techniques for stochastic mechanics, 2019.
- G. Bolch, S. Greiner, H. de Meer, and K. S. Trivedi. *Queueing Networks and Markov Chains: Modeling and Performance Evaluation with Computer Science Applications*. John Wiley & Sons, 2006.
- M. Capiński and T. Zastawniak. *Probability through Problems*. Springer Verlag, 2nd edition, 2003.
- D.R. Cox, editor. *Renewal Theory*. John Wiley & Sons Inc, New York, 1962.
- P. Doyle and J Laurie Snell. *Random Walks and Electrical Networks*. Mathematical Association of America, 1984.
- M. El-Taha and S. Stidham Jr. *Sample-Path Analysis of Queueing Systems*. Kluwer Academic Publishers, 1998.
- J.R. Norris. *Markov Chains*. Cambridge University Press, 1997.
- H.C. Tijms. *Stochastic Models, An Algorithmic Approach*. J. Wiley & Sons, 1994.
- H.C. Tijms. *A First Course in Stochastic Models*. John Wiley & Sons, Chichester, 2003.
- A.A. Yushkevich and E.B. Dynkin. *Markov Processes: Theorems and Problems*. Plenum Press, 1969.

NOTATION

- a_k = Number of arrivals in the k th period
- $A(t)$ = Number of arrivals in $[0, t]$
- A_k = Arrival time of k th job
- \tilde{A}_k = Start of service of k th job
- c_n = Service/production capacity in the n th period
- d_n = Number of departures in the n th period
- c = Number of servers
- C_a^2 = Squared coefficient of variation of the inter-arrival times
- C_s^2 = Squared coefficient of variation of the service times
- $D(t)$ = Number of departures in $[0, t]$
- $D_Q(t)$ = Number of customers/jobs that departed from the queue in $[0, t]$
- D_k = Departure time of k th job
- F = Distribution of the service time of a job
- $L(t)$ = Number of customers/jobs in the system at time t
- $Q(t)$ = Number of customers/jobs in queue at time t
- $L_S(t)$ = Number of customers/jobs in service at time t
- $E[L]$ = Long run (time) average of the number of jobs in the system
- $E[L_Q]$ = Long run (time) average of the number of jobs in queue
- $E[L_S]$ = Long run (time) average of the number of jobs in service
- $N(t)$ = Number of arrivals in $[0, t]$
- $N(s, t)$ = Number of arrivals in $(s, t]$
- $p(n)$ = Long-run time average that the system contains n jobs
- Q_k = Queue length as seen by the k th job, or at the *end* of the k th period
- S_k = Service time required by the k th job
- $S(t)$ = Total service time available in $[0, t]$
- S = Service time of a generic job
- W_k = Sojourn time of k th job
- $W_{Q,k}$ = Time in the queue of k th job
- $E[W]$ = Sample average of the sojourn time
- $E[W_Q]$ = Sample average of the time in queue
- X_k = Inter-arrival time between job $k - 1$ and job k
- X = Generic inter-arrival time between two consecutive jobs
- δ = Departure rate
- λ = Arrival rate

μ = Service rate

$\pi(n)$ = Stationary probability that an arrival sees n jobs in the system

ρ = Load on the system

FORMULA SHEET

$$\rho = \lambda \frac{E[S]}{c}$$

$$E[W_Q] = \frac{C_a^2 + C_s^2}{2} \frac{\rho^{\sqrt{2(c+1)}-1}}{c(1-\rho)} E[S]$$

$$\text{Batching: } C_{sB}^2 = \frac{B V[S_0] + V[T]}{(B E[S_0] + E[T])^2}$$

$$\text{Nonpreemptive: } V[S] = V[S_0] + \frac{V[T]}{B} + \frac{B-1}{B^2} (E[T])^2$$

$$\text{Preemptive: } A = \frac{m_f}{m_r + m_f}, C_s^2 = C_0^2 + 2A(1-A) \frac{m_r}{E[S_0]}$$

$$C_{di}^2 = 1 + (1 - \rho_i^2)(C_{ai}^2 - 1) + \frac{\rho_i^2}{\sqrt{c_i}}(C_{si}^2 - 1)$$

$$f_i(n_i) = \frac{(c_i \rho_i)^{n_i}}{n_i! G(i)} \mathbb{1}_{n_i < c_i} + \frac{c_i^{c_i} \rho_i^{n_i}}{c_i! G(i)} \mathbb{1}_{n_i \geq c_i},$$

$$\text{with } G(i) = \sum_{n=0}^{c_i-1} \frac{(c_i \rho_i)^n}{n!} + \frac{(c_i \rho_i)^{c_i}}{c_i!} \frac{1}{1 - \rho_i}$$

$$E[L_i] = \frac{(c_i \rho_i)^{c_i}}{c_i! G(i)} \frac{\rho_i}{(1 - \rho_i)^2} + c_i \rho_i$$

INDEX

- arrival process, 11
- arrival rate, 2, 18
- arrival times, 11
- average number of jobs, 20

- balance equations, 42
- balking, 53
- binomially distributed, 2
- Burke's law, 68

- conditional probability, ix

- departure rate, 18
- departure time of the system, 12
- distribution function, ix

- Economic Order Quantity (EPQ), 67
- Economic Production Quantity (EPQ), 67
- effective processing time, 29
- expected sojourn time, 20
- expected waiting time in queue, 20
- exponentially distributed, 10

- Gerschgorin's disk theorem, 70

- i.i.d., 2
- indicator function, viii
- inter-arrival times, 11

- Jackson networks, 69

- Kendall's abbreviation, 15

- level-crossing equations, 40
- limiting, 21
- load, 20, 37

- memoryless, 10
- Merging, 3
- moment-generating function, ix

- net processing time, 29
- non-preemptive, 31
- normalization constant, 41
- number of jobs in the system, 13

- open, 68

- PASTA, 43
- PK formula, 55
- Poisson arrivals see time averages, 43
- Poisson distributed, 2
- Poisson process, 2
- Pollaczek-Khinchine formula, 55
- probability mass function, viii
- processing rate, 19
- product-form solution, 69

- rate-stable, 19
- remaining service time, 56
- renewal reward theorem, 37

- SCV, 3
- service rate, 19
- small o notation, viii
- sojourn time, 12
- square coefficient of variation, 3
- stationary, 21
- stationary and independent increments, 2
- steady-state, 21
- survivor function, ix

- time-average number of jobs, 21
- traffic (rate) equations, 68

- utilization, 37

- virtual waiting time process, 13

- waiting time in queue, 11