

Analysis of Queueing Systems with Sample Paths and Simulation

Nicky D. van Foreest

March 31, 2019

CONTENTS

Introduction	v
1 CONSTRUCTION AND SIMULATION OF QUEUEING SYSTEMS	1
1.1 Preliminaries	1
1.2 Poisson Distribution	3
1.3 Queueing Processes in Discrete-Time	5
1.4 Exponential Distribution	11
1.5 Single-server Queueing Process in Continuous Time	12
2 ANALYTICAL MODELS	17
2.1 Kendall's Notation	17
2.2 Queueing Processes as Regulated Random Walks	18
2.3 Rate Stability and Utilization	21
2.4 Renewal Reward Theorem and load	23
2.5 (Limits of) Empirical Performance Measures	24
2.6 Level Crossing and Balance Equations	25
2.7 $M/M/1$ queue	31
2.8 $M(n)/M(n)/1$ Queue	33
2.9 Poisson Arrivals See Time Averages	35
2.10 Little's Law	36
2.11 $M^X/M/1$ Queue: Expected Waiting Time	38
2.12 $M/G/1$ Queue: Expected Waiting Time	40
2.13 $M^X/M/1$ Queue Length Distribution	43
2.14 $M/G/1$ Queue Length Distribution	46
3 APPROXIMATE MODELS	49
3.1 $G/G/c$ Queue: Approximations	49
3.2 Setups and Batch Processing	53
3.3 Non-preemptive Interruptions, Server Adjustments	55
3.4 Preemptive Interruptions, Server Failures	56
4 QUEUEING NETWORKS	59
4.1 Open Single-Class Product-Form Networks	59
4.2 Tandem queues	60
4.3 Gordon-Newell Networks	61
4.4 MVA Algorithm	62
Bibliography	65
Notation	67
Formula Sheet	69
Index	69

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). 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 some 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 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 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. By implementing these constructions in Python code we can then simulate and analyze such systems. Besides that, simulation provides ample motivation of why and how we deal with queueing systems, simulation is useful to analyzing realistic systems, as mathematical models have severe shortcomings in such cases. Consider, for example, the service process at a check-in desk of KLM. Business customers and economy customers are served by two separate queueing systems. The business customers are served by one server, server A say, while the economy class customers by three servers, say. What would happen to the sojourn time of the business customers if server A would be allowed to serve economy class customers when the business queue is empty? For the analysis of such complicated control policies, simulation is the most natural approach.

In Chapter 2 and Chapter 3 we derive exact and approximate models, respectively, for single-station queueing systems. The benefit of such models is that they offer structural insights into the behavior of the system and scaling laws, such as that the average waiting time scales (more or less) linearly in the variance of the service times of individual customers. The main idea 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. In fact, the performance measures defined for sample paths are precisely those we compute with simulation.

In Chapter 4 we construct algorithms to analyze open and closed queueing networks. Many of the sample path results developed for the single-station case can be applied to these networks. Thus, with sample-path methods we relate the theory, algorithms and simulation of queueing systems. For this part we refer to the book of Prof. Zijm; the present set of notes augments the discussion there.

Our aim is not to provide rigorous proofs for all results derived in the book. For this we refer to following books.

1. Bolch et al. [2006]

2. Capiński and Zastawniak [2003]
3. El-Taha and Stidham Jr. [1998]
4. Tijms [1994] and/or Tijms [2003]

EXERCISES






I urge you to try to make as many exercises as possible. The main text contains hardly any examples or derivations: the exercises *illustrate* the material and force you to *think* about the technical parts. The exercises require many of the tools you learned previously in courses on calculus, probability, and linear algebra. Here you can see them applied. Moreover, many of these tools will be useful for other, future, courses. Thus, the investments made here will pay off for the rest of your (student) career.

As a guideline to making the exercises I recommend the following approach. First read the notes. Then attempt to make an exercise for a few minutes by yourself. If by that time you have not obtained a good idea on how to approach the problem, check the hints and then the solution. Once you have understood the solution, try to repeat the arguments *with the solution manual closed*.

You'll notice that some of these problems are quite difficult, often not because the problem itself is difficult, but because you need to combine a substantial amount of knowledge all at the same time. All this takes time and effort. Realize that the exercises are not intended to be easy (otherwise we could have been satisfied with computing $1 + 1$). The problems should be doable, but hard.

The book is, admittedly, pretty big. The main reason is that the hints and solutions are very explicit, and spell out nearly every intermediate steps. 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. I also included exercises to show to obtain some result in several different ways. Thirdly, the numerical calculations show each intermediate numerical result along with the Python code. Like this, if you get stuck somewhere in the computations you can precisely check where you go wrong.

The following symbols are used to classify the type of exercise:

- : Computation.
- : Test some (simple) technical aspect.
- : This is a claim and it is up to you to decide whether it is correct or not.
- : Illustration.
- : hard, you can skip this if you run short of time.

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. Prof. dr. W.H.M. Zijm allowed me to use the first few chapters of his book. 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.

CONSTRUCTION AND SIMULATION OF QUEUEING SYSTEMS

In this chapter we start with a discussion of the Poisson process. We then construct queueing processes in discrete time and apply the Poisson process to model the number of arrivals in periods of fixed length. In Section 1.4 we relate the exponential distribution to the Poisson distribution. The exponential distribution often serves as a good model for inter-arrival times of individual jobs. As such this is a key component of the construction of queueing processes in continuous time. As it turns out, both ways to construct queueing processes are easily implemented as computer programs, thereby allowing us to use simulation to analyze queueing systems. In passing we develop a number of performance measures to provide insight into the (transient and average) behavior of queueing processes.

We assume that you *know all* results of Section 1.1.

1.1 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 course.

Theory and Exercises

We use the notation:

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

where the last equation defines an *indicator variable*.

The function $f(h) = o(h)$ means that f is such $f(h)/h \rightarrow 0$ as $h \rightarrow 0$. If we write $f(h) = o(h)$ it is implicit that $|h| \ll 1$. We call this *Small o notation*.

1.1.1. 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

$$(a + b)^n = \sum_{i=0}^n \binom{n}{i} a^{n-i} b^i, \tag{1.1.1a}$$

$$e^x = \lim_{n \rightarrow \infty} (1 + x/n)^n, \tag{1.1.1b}$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots = \sum_{k=0}^{\infty} \frac{x^k}{k!}, \tag{1.1.1c}$$

$$\sum_{n=0}^N = \frac{1 - \alpha^{N+1}}{1 - \alpha}. \quad (1.1.1d)$$

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

$$X = \sum_{n=0}^{\infty} X \mathbb{1}_{X=n} = \sum_{n=0}^{\infty} n \mathbb{1}_{X=n}, \quad (1.1.2a)$$

$$E[X] = \sum_{n=0}^{\infty} n f(n), \quad (1.1.2b)$$

$$E[g(X)] = \sum_{n=0}^{\infty} g(n) f(n) \quad (1.1.2c)$$

$$E[\mathbb{1}_{X \leq x}] = P(X \leq x), \quad (1.1.2d)$$

$$V[X] = E[X^2] - (E[X])^2, \quad (1.1.2e)$$

1.1.2. Define *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, this idea makes the computation of certain expressions quite a bit easier.

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

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 shorthand for $f(x) dx$ ¹. Recall that

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

1.1.4. 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 $E[X^2] = 2 \int_0^{\infty} y G(y) dy$, where $G(x) = 1 - F(x)$, provide the second moment exists.

1.1.5. Now use integration by parts to show that for a continuous non-negative random variable X with distribution function F and survivor function $G = 1 - F$, $\int_0^{\infty} y G(y) dy = E[X^2]/2$.

You should know that for the *moment-generating function* $M_X(s)$ of a random variable X and s a real number sufficiently small that the expectation(s) below exists:

$$M_X(s) = E[e^{sX}], \quad (1.1.3a)$$

$$M_X(s) \text{ uniquely characterizes the distribution of } X, \quad (1.1.3b)$$

$$E[X] = M'_X(0) = \left. \frac{dM_X(s)}{ds} \right|_{s=0}, \quad (1.1.3c)$$

$$E[X^2] = M''_X(0), \quad (1.1.3d)$$

$$M_{X+Y}(s) = M_X(s) \cdot M_Y(s), \quad \text{if } X \text{ and } Y \text{ are independent,} \quad (1.1.3e)$$

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

¹ For the interested reader, $\int x dF(x)$ is a Lebesgue-Stieltjes integral with respect to the distribution function F .

1.1.6. We have one gift to give to one 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 the number wins the gift. Show that the probability of winning the gift is the same for each child.

You should know that:

$$P(A|B) = \frac{P(AB)}{P(B)}, \quad \text{if } P(B) > 0, \quad (1.1.4a)$$

$$P(A) = \sum_{i=1}^n P(AB_i) = \sum_{i=1}^n P(A|B_i)P(B_i), \quad \text{if } A = \bigcup_{i=1}^n B_i \text{ and } P(B_i > 0) \text{ for all } i. \quad (1.1.4b)$$

1.2 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 entered during the time interval $[0, t]$ and, with this, $N(s, t] = N(t) - N(s)$. 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

$$E[N(s, t)] = \lambda(t - s). \quad (1.2.1)$$

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 second 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 such that $p = P(B_i = 1)$ and $P(B_i = 0) = 1 - P(B_i = 1)$, and we assume that $\{B_i\}$ are independent. The total number of arrivals $N_n(t)$ that occur in n intervals is then *binomially distributed*, i.e.,

$$P(N_n(t) = k) = \binom{n}{k} p^k (1 - p)^{n-k}. \quad (1.2.2)$$

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.,

$$P(N(t) = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}, \quad (1.2.3)$$

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.

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

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

1.2.2. 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

$$C^2 = \frac{V[X]}{(E[X])^2}. \quad (1.2.4)$$

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

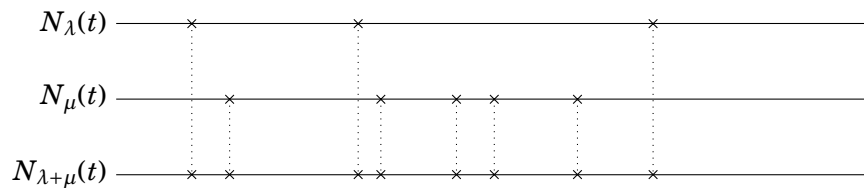
1.2.3. Show that the SCV of $N(t) \sim P(\lambda t)$ is equal to $1/(\lambda t)$. What does this mean for t large?

1.2.4. Show that

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

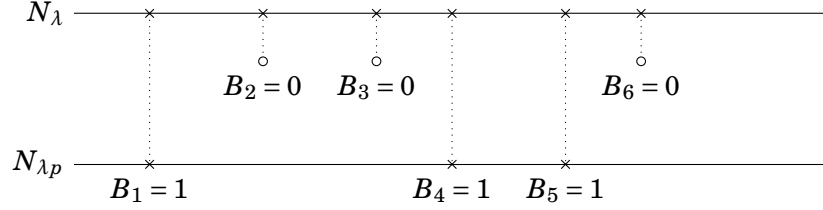
when $N(t) \sim P(\lambda t)$ and h is small.

Merging Poisson processes occurs often in practice. We have two Poisson processes, for instance, the arrival processes N_λ of men and N_μ 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 whose rate is the sum of the rates of the two original Poisson processes.



1.2.5. 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 $(\lambda + \mu)t$.

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.



1.2.6. Show with moment-generation 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.

1.3 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. We close the section with a large number of exercises to develop recursions for a large number of different queueing systems, and illustrate how powerful this approach is.

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 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 after each other, 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:

$$\begin{aligned}
 a_k &= \text{number of jobs that arrive in period } k, \\
 c_k &= \text{the capacity, i.e., the maximal number of jobs that can be served, during period } k, \\
 d_k &= \text{number of jobs that depart in period } k, \\
 L_k &= \text{number of jobs in the system at the end of period } k.
 \end{aligned}
 \tag{1.3.1}$$

In the sequel we also call a_k the size of the batch arriving in period k . 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 .

Let L_{k-1} be 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

$$d_k = \min\{L_{k-1}, c_k\}, \tag{1.3.2a}$$

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

$$L_k = L_{k-1} - d_k + a_k. \tag{1.3.2b}$$

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 section ‘queue length’ and ‘number of jobs in the system’ coincide.

1.3.1. Show that the scheme

$$\begin{aligned} L_k &= [L_{k-1} + a_k - c_k]^+, \\ d_k &= L_{k-1} + a_k - L_k. \end{aligned} \tag{1.3.3}$$

is equivalent to a modification of (1.3.2) 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 (1.3.2). 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

As a first step we model the arrival process of patients as a Poisson process, cf., Section 1.2. 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's 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 (1.3.2) 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 Figure 1. It is apparent that Suggestions 1 and 2 above do not significantly affect the behavior of the queue length process.

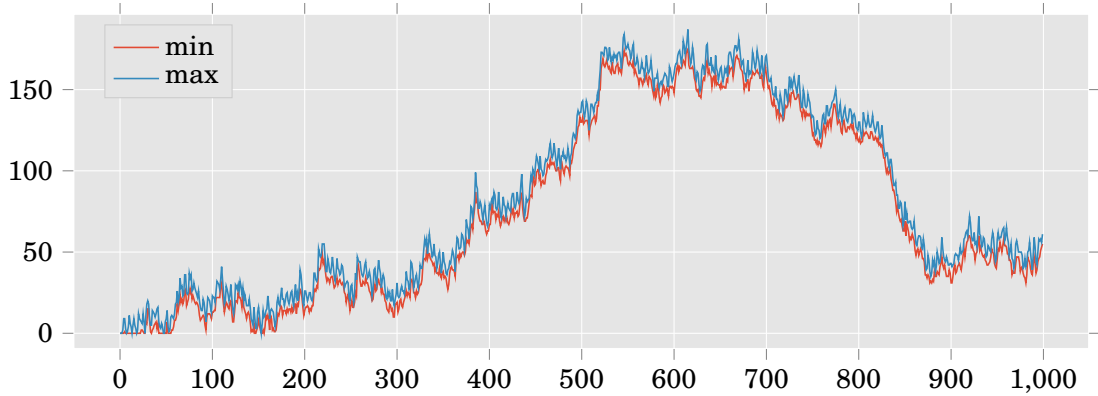


Figure 1: Effect of capacity and holiday plans. We plot for each time point the maximum and the minimum queue length for each of the policies. Apparently, the effect of each of these policies is, for all practical purposes, negligible.

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 to implement this is by considering last week's queue L_{n-1} : if $L_{n-1} < 12$, i.e., the service capacity of one week, then do e intakes less. When $L_{n-1} > 24$, i.e., larger than two weeks of intakes, do e intakes more. Here, $e = 1$ or 2 , or perhaps a larger number; it corresponds to the amount of control we want to 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.

Figure 2 shows a striking difference indeed. The queue does not explode any more, 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 be-

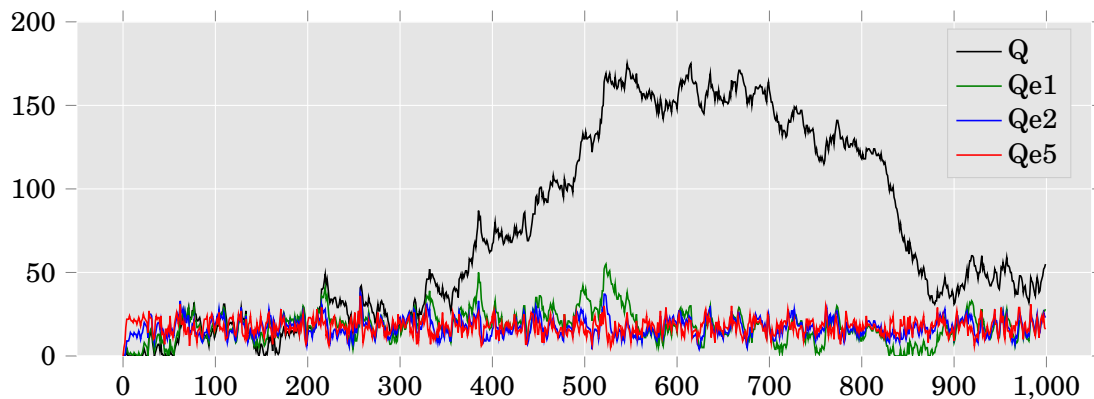


Figure 2: Controlling the number of intakes. Clearly, adapting the service rate ‘does wonders’ to control the queue length.

havior. However, controlling the service rate as a function of the queue length improves the situation quite dramatically.

Conclusion

Observe that, even with these (deceitfully) simple recursions, we can obtain considerable insight into this, otherwise, very complicated controlled queueing process². As a matter of fact, with such simple recursions we can analyze many practical queueing situations. Together with students the author applied it numerous times, for instance:

- Should a certain hospital invest in a new MRI scanner to reduce waiting times?
- When to switch on and off a tin bath at an electronics component factory?
- How to route post parcels in a post sorting center?

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 performance measures of relevance, 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 add service capacity. With simulation it is easy to study the effect of, hence evaluate, such decisions.

Exercises

The reader should understand from the above case that, once we have the recursions, we can analyze the system and make plots to evaluate suggestions for improvement. Thus, getting the

² If the reader doubts the value of simulation, s/he should try to develop a mathematical method to analyze multi-server queueing systems with vacations, of which this case is an example.

recursions is crucial to construct, i.e., model, queueing processes. For this reason, most of the exercises below focus on obtaining recursions for many different queueing systems.³

1.3.2 (Queue with Blocking). 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. 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?

1.3.3 (Estimating the lead time distribution). 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 $W_{+,k}$.

1.3.4 (Cost models). 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.

1.3.5 (Priority queueing). 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.

1.3.6 (Queue with protected service capacity and lost capacity). 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 .

1.3.7 (Tandem networks). 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 (1.3.2) to simulate this situation.

1.3.8 (Merging departure streams). 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.

1.3.9 (Inventory control). The recursions used in the exercises above can also be applied to analyze inventory control policies. Consider a production system that can produce maximally M_k

³It 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.

items per week during normal working hours, and maximally N_k items during extra (weekend and evening) hours. Let, for period k ,

- D_k = Demand in week k ,
- S_k = Sales, i.e., number of items sold, in week k ,
- r_k = Revenue per item sold in week k ,
- X_k = Number of items produced in week k during normal hours,
- Y_k = Number of items produced in week k during extra hours,
- c_k = Production cost per item during normal hours,
- d_k = Production cost per item during extra hours,
- h_k = Holding cost per item, due at the end of week k ,
- I_k = On hand inventory level at the end of week k .

Management needs a production plan that specifies for the next T weeks the number of items to be produced per week. Formulate this problem as an LP problem, taking into account the inventory dynamics. Assume that demand must be met from on-hand inventory.

1.4 EXPONENTIAL DISTRIBUTION

In Section 1.2 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 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 $1/\mu$ 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.* (*independent and identically distributed*) and exponentially distributed with mean $1/\lambda$, 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.

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

1.4.2. 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. $\sim \text{Exp}(\lambda)$.

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

1.4.4. If the inter-arrival times $\{X_i\}$ are i.i.d. $\sim \text{Exp}(\lambda)$, prove that the number $N(t)$ of arrivals during 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 .

1.4.5. 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].

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

In Section 1.3 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 (1.3.2) 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. Specially, 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 define

$$A_k = \min\{t : A(t) \geq k\}, \quad (1.5.1)$$

and set $A_0 = 0$ ⁴. 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

$$X_k = A_k - A_{k-1}. \quad (1.5.2)$$

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

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

⁴ If we want to be mathematically precise, we must take \inf rather than \min . However, in this set of notes we do not want to distinguish between subtleties.

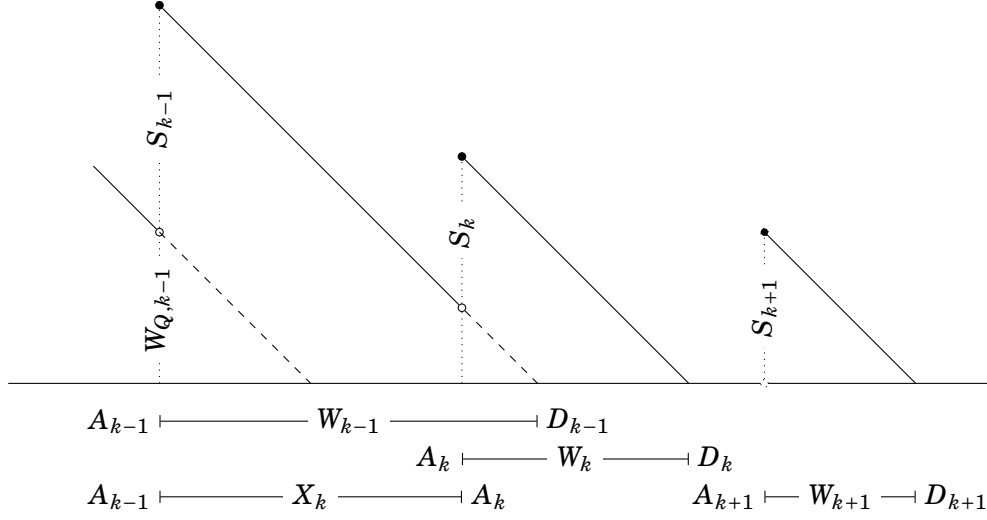


Figure 3: Construction of the single-server queue in continuous time. The sojourn time W_k of the k th job is the sum of the work in queue $W_{Q,k}$ at its arrival epoch A_k and its service time S_k ; its departure time is then $D_k = A_k + W_k$. The waiting time of job k is clearly equal to $W_{k-1} - X_k$. 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 .

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

$$A(t) = \max\{k : A_k \leq t\}. \quad (1.5.3)$$

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

We next construct the *waiting time in queue* $\{W_{Q,k}\}$ as seen by the arrivals. From Figure 3 it is evident that

$$W_{Q,k} = [W_{Q,k-1} + S_{k-1} - X_k]^+. \quad (1.5.4)$$

To use this, set $W_{Q,0} = 0$, from which we can compute $W_{Q,1}$, and then $W_{Q,2}$ and so on.

1.5.1. 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

$$W_k = D_k - A_k = \tilde{A}_k + S_k - A_k = W_{Q,k} + S_k, \quad (1.5.5)$$

where each of these equations has its own interpretation.

1.5.2. Explain the following recursions for the $G/G/1$ queue:

$$\begin{aligned} A_k &= A_{k-1} + X_k, \\ D_k &= \max\{A_k, D_{k-1}\} + S_k, \\ W_k &= D_k - A_k. \end{aligned} \quad (1.5.6)$$

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.

1.5.3. 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. Figure 4,

$$L(t) = A(t) - D(t) + L(0), \quad (1.5.7)$$

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 simulate a queueing process in continuous time, A bit of experimentation with computer programs such as R or Python will reveal that this is easy.

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

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

$$L(A_k -). \quad (1.5.8)$$

(Since in queueing systems we are concerned with processes with jumps, we need to be quite particular about left and right limits at jump epochs.)

1.5.5. Explain the following (algorithmic efficient) procedure to compute the number of jobs in the system as seen by arrivals:

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

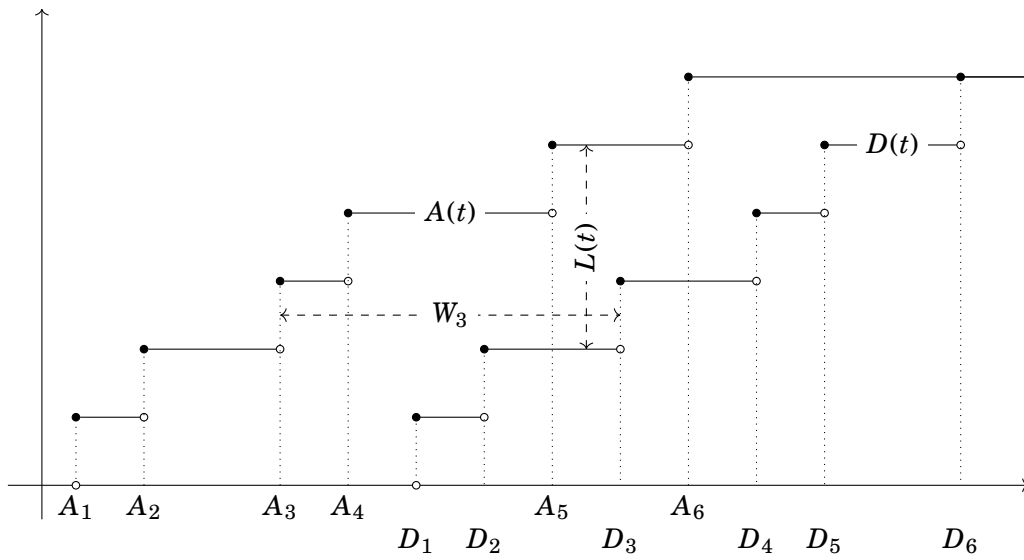


Figure 4: Relation between the arrival process $\{A(t)\}$, the departure process $\{D(t)\}$, the number in the system $\{L(t)\}$ and the waiting times $\{W_k\}$. In particular, $L(t)$ is the difference graphs of $A(t)$ and $D(t)$.

ANALYTICAL MODELS

In this chapter we focus on developing analytic models for various queueing systems in steady-state. In the analysis we use sample-path and level-crossing arguments to count how often certain events occur as a function of time. Then we define probabilities in terms of limits of fractions of these counting processes. Like this the performance measures can be explicitly computed for the statistical analysis of (simulations of) queueing systems.

We start with developing a useful set of shorthands to distinguish between different queueing models. Then we include a section to motivate why we focus on a steady-state analysis of queueing systems.

As a reminder, we keep the discussion in these notes mostly at an intuitive level, and refer to El-Taha and Stidham Jr. [1998] for proofs and further background.

2.1 KENDALL'S NOTATION

As became apparent in Sections 1.3 and 1.5, 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-serve (FCFS).

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. What is the meaning of $M(n)/M(n)/1$?

2.1.2. What is the meaning of $M^X/M/1$?

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

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

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.

2.1.5. 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 the construction of queueing processes as set out in Section 1.3 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., (1.3.3). Observe that in relation (1.3.3) 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. In this section we show how to build the queueing process $\{L_k\}$ from the random walk $\{Z_k\}$ using a device called a *reflection map*, which gives an elegant construction of a queueing process. Moreover, we can use the probabilistic tools that have been developed for the random walk to analyze queueing systems. One example is the distribution of the time until an especially large queue length is reached; these times can be formulated as *hitting times* of the random walk. Another example is the average time it takes to clear a large queue.

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 Figure 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 Figure 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 $a_k = c_k$.

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. For this, we first relate the random walk of the type (2.2.3) to a random walk in continuous time.

2.2.2. Let $N_{\lambda+\mu}$ be a Poisson process with rate $\lambda + \mu$. If $\{a_k\}$ is an i.i.d. sequence of Bernoulli random variables such that $P(a_k = 1) = \lambda/(\lambda + \mu) = 1 - P(a_k = 0)$, show that the random variable

$$N_\lambda(t) = \sum_{k=1}^{\infty} a_k \mathbb{1}_{k \leq N_{\lambda+\mu}(t)},$$

has a Poisson distribution with rate λt .

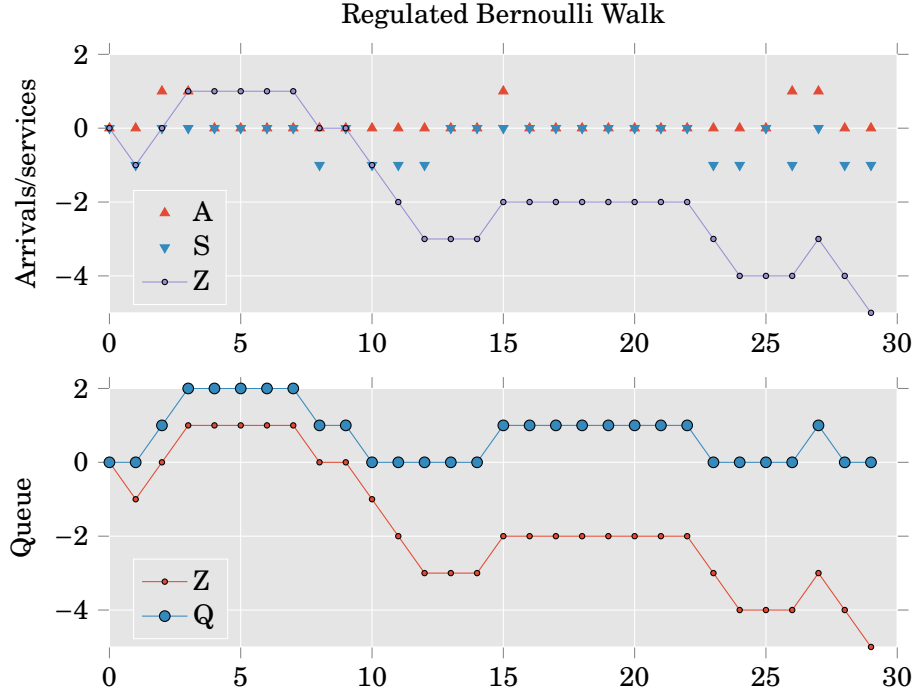


Figure 5: The upper panel shows a graph of the random walk Z . An upward pointing triangle corresponds to an arrival, a downward triangle to a potential service. The lower panel shows the queueing process $\{L_k\}$ as a random walk with reflection.

Similarly, let

$$N_\mu(t) = N_{\lambda+\mu}(t) - N_\lambda(t) = \sum_{k=1}^{\infty} (1 - a_k) \mathbb{1}_{N_{\lambda+\mu}(t) \leq k};$$

but this is $N_{\lambda+\mu}(t)$ thinned by the Bernoulli random variables $\{1 - a_k\}$. Let $N_\lambda = \{N_\lambda(t)\}$ and $N_\mu = \{N_\mu(t)\}$ be the associated Poisson processes.

With the processes N_λ and N_μ constructed above from the sequence $\{a_k\}$ and the Poisson process $N_{\lambda+\mu}$ we can define the process $Z = \{Z(t)\}$ such that

$$Z(t) = Z(0) + N_\lambda(t) - N_\mu(t).$$

Thus, we let N_λ correspond to job arrivals and N_μ to departures. Observe that the times $\{T_k\}$ at which Z makes jumps are such that $T_k - T_{k-1}$ have exponential distribution with mean $1/(\lambda + \mu)$. At the jump times, $Z(T_k) = Z_k$, where Z_k satisfies (2.2.3) with $P(a_k = 1) = \lambda/(\lambda + \mu)$. We call Z the *free M/M/1 queue* as, contrary to the real M/M/1 queue, Z can take negative values.

2.2.3. Show that

$$P_m(Z(t) = n) = e^{-(\lambda+\mu)t} \left(\frac{\lambda}{\mu}\right)^{(n-m)/2} \sum_{k=0}^{\infty} \frac{(t\sqrt{\lambda\mu})^{2k+m-n}}{k!(k+m-n)!},$$

where $P_m(\cdot)$ means that the random walk starts at m , i.e., $Z(0) = m$. As an aside, the summation includes negative factorials when $k + m - n < 0$. The tacit assumption is to take $n! \in \{\pm\infty\}$ for $n \in \mathbb{Z}_-$. Another way to get around this problem is to take $k = \max\{0, m - n\}$.

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



Figure 6: Another example of a reflected random walk.

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. And the $M/M/1$ is but the simplest queueing system; other queueing systems will be more complicated yet. We therefore give up the analysis of such transient queueing systems and we henceforth contend ourselves with the analysis of queueing systems in the limit as $t \rightarrow \infty$. The limiting random variable L is known as the *steady-state limit* of the sequence of random variables $\{L_k\}$, and the distribution of L is known as the *limiting distribution* or *stationary distribution* of $\{L_k\}$. Taking these limits warrants two questions: what type of limit is actually meant here, and what is the rate of convergence to this limiting situation? In these notes we sidestep all these fundamental issues, as the details require measure theory and more advanced probability theory than we can deal with in this course.

We illustrate the rate of convergence to the limiting situation by means of an example. Specifically, we consider the sequence of waiting times $\{W_{Q,k}\}$ to a limiting random variable W_L , where $W_{Q,k}$ is constructed according to the recursion Eq. (1.5.4). 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 Eq. (1.5.4) to compute the *exact* distribution of $W_{Q,k}$ for $k = 1, 2, \dots, 20$, cf., the left panel in Figure 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. Here the *empirical distribution* for the simulation is defined as

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

where $W_{Q,k}$ is obtained by simulation. 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 *stationarity*. The study of queueing systems in stationary state will occupy us for the rest of the book.

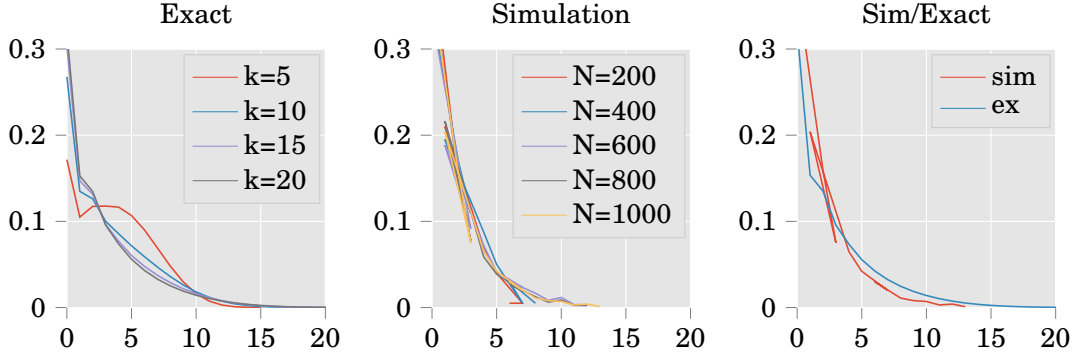


Figure 7: The density of $W_{Q,k}$ for $k = 5, 10, 15, 20$ computed by an exact method as compared the density obtained by simulation of different run lengths $N = 200, 400, \dots, 1000$. The right panel compares the exact density of $W_{Q,20}$ to the density obtained by simulation for $N = 1000$.

2.2.4. 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.

2.2.5. Validate the results of Figure 7 with simulation.

2.3 RATE STABILITY AND UTILIZATION

In the analysis of any queueing process the first step should be to check the relations between the arrival, service and departure rates. The concept of rate is crucial because it captures our intuition that when, on the long run, jobs arrive faster than they can leave, the system must ‘explode’. As we will see, when the arrival rate is smaller than the service rate the system is stable. Thus, the first performance measure we need to estimate for a queueing system is the ratio between the arrival and service rate. In this section we develop some concepts and notation to formalize these ideas. We will use these concepts throughout the remainder of the book.

We first formalize the *arrival rate* and *departure rate* in terms of the *counting processes* $\{A(t)\}$ and $\{D(t)\}$. The *arrival rate* is the long-run average number of jobs that arrive per unit time, 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 *independent and identically distributed (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 1.5; the strong law of large numbers then guarantees that the above limit exists.

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 Eq. (2.3.1) that the average inter-arrival time between two consecutive jobs is

$$E[X] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n X_k = \lim_{n \rightarrow \infty} \frac{A_n}{A(A_n)} = \lim_{t \rightarrow \infty} \frac{t}{A(t)} = \frac{1}{\lambda}, \quad (2.3.2)$$

where we take $t = A_n$ in the limit for $t \rightarrow \infty$. In words, the above states that the arrival rate λ is the inverse of the expected inter-arrival time.

The development of the departure times $\{D_k\}$ is entirely analogous to that of the arrival times; define the *departure rate* as

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

Assume now that there is a single server. Let S_k be the required service time of the k th job to be served, and define

$$U_n = \sum_{k=1}^n S_k$$

as the total service time required by the first n jobs. With this, let

$$U(t) = \max\{n : U_n \leq t\}$$

and define the *service rate* or *processing rate* as

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

In the same way as we derived that $E[X] = 1/\lambda$, we obtain for the expected (or average) service time required by an individual job

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

Now observe that, if the system is empty at time 0, it must be that at any time 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,

$$\delta = \lim_t \frac{D(t)}{t} \leq \lim_t \frac{A(t)}{t} = \lambda. \quad (2.3.4)$$

We call a system *rate stable* if

$$\lambda = \delta,$$

in other words, the system is stable if, on the long run, jobs leave the system just as fast as they arrive. Of course, if $\lambda > \delta$, the system length process $L(t) \rightarrow \infty$ as $t \rightarrow \infty$.

It is also evident that jobs cannot depart faster than they can be served, hence, $D(t) \leq U(t)$ for all t . Combining this with the fact that $\delta \leq \lambda$, we get

$$\delta \leq \min\{\lambda, \mu\}.$$

When $\mu \geq \lambda$ the above inequality reduces to $\delta = \lambda$ for rate-stable systems². As it turns out, when $\mu = \lambda$ and the variance of the service time $V[S] > 0$ or $V[X] > 0$ the queue length process can behave in a very peculiar way. For this reason we henceforth (and implicitly) require that $\mu > \lambda$.

² It would be interesting to prove this.

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

2.3.2. Consider a paint factory which 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 (rate) 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, so that $\mu < \lambda$. Thus, it is necessary to identify minimal batch sizes to ensure that $\mu > \lambda$. Motivate that the following linear program can be used to determine the minimal batch sizes:

minimize T

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

2.4 RENEWAL REWARD THEOREM AND LOAD

We start with stating and proving (graphically) the *renewal reward theorem*. In the sequel we will see many applications of this theorem. In this section we use it to relate the fraction of time the server is busy in a $G/G/1$ queue to the job arrival rate and the expected job service time.

The renewal reward theorem is very useful, and states intuitively that when customers arrive at rate λ and each customer pays an average amount X , then the system earns money at rate $Y = \lambda X$. Figure 8 provides graphical motivation about why this theorem is true; El-Taha and Stidham Jr. [1998] gives a (simple) proof.

Theorem 2.4.1 (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 $Y = \lambda X$.

Define the *load* or *utilization* as the limiting fraction of time the server is busy, i.e.,

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

2.4.1. Use the renewal reward theorem to prove that $\rho = \lambda E[S]$.

2.4.2. 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 complete the argument.

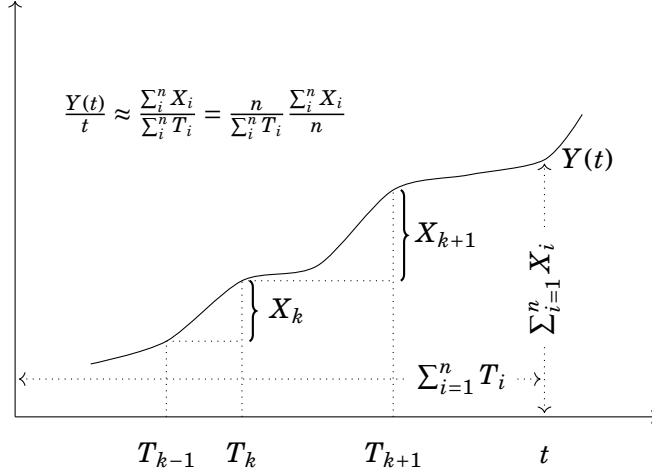


Figure 8: A graphical ‘proof’ of $Y = \lambda X$. Here $Y(t)/t \rightarrow Y$, $n/\sum_i^n T_i \rightarrow \lambda$ and $n^{-1}\sum_i^n X_i \rightarrow X$. (Observe that in the figure X_k does not represent an inter-arrival time; instead it corresponds to the increment of (the graph of) $Y(t)$ between two consecutive epochs T_{k-1} and T_k at which $Y(t)$ is observed.)

From the identities $\lambda^{-1} = E[X]$ and $\mu^{-1} = E[S]$, we get a further set of relations:

$$\rho = \lambda E[S] = \frac{\lambda}{\mu} = \frac{E[S]}{E[X]}.$$

Thus, the load has also the interpretation as the rate at which jobs arrive times the average amount of work per job. Finally, recall that for a system to be rate-stable, it is necessary that $\mu > \lambda$, implying in turn that $\rho < 1$. 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., $E[S] < E[X]$. In fact, when $\mu < \lambda$, it is easy to check with simulation that the queue length grows roughly linearly with slope $\lambda - \mu$.

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

2.5 (LIMITS OF) EMPIRICAL PERFORMANCE MEASURES

If the arrival and service processes are such that the queueing system is rate-stable, we can sensibly define other performance measures such as the average waiting time. In this section we define the second most important performance measures; recall that the most important is the utilization ρ . At the end we provide an overview of the relations between these performance measures in Figure 17.

With the construction of queueing processes in Section 1.5 we can compute the waiting time as observed by the first n , say, jobs. We therefore define the *expected waiting time* as

$$E[W] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n W_k, \quad (2.5.1)$$

and the expected time in queue as

$$E[W_Q] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n W_{Q,k}. \quad (2.5.2)$$

Note that these performance measures are limits of *empirical* measures. Note also that these statistics are as *observed by arriving jobs*: the first job has a waiting time W_1 at its arrival epoch, the second a waiting time W_2 , and so on. For this reason we colloquially say that $E[W]$ is the average waiting time as ‘seen by arrivals’. The *distribution of the waiting times at arrival times* can be found by counting:

$$P(W \leq x) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{W_k \leq x}. \quad (2.5.3)$$

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

$$E[L] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n L(A_k -), \quad (2.5.4)$$

where $L(A_k -)$ is the number of jobs in the system at the arrival epoch of the k th job. The *distribution of $\{L(t)\}$ as seen by customers upon arrival*, is

$$P(L \leq m) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{L(A_k -) \leq m}. \quad (2.5.5)$$

We call $P(L > m)$ the *excess probability*.

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 sampling (observation) moments. Assuming the limit exists we use (1.5.7) to define the *time-average number of jobs* as

$$E[L] = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t L(s) ds. \quad (2.5.6)$$

Observe that, notwithstanding that the symbols are the same, this expectation need not be the same as (2.5.4). 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

$$P(L \leq m) = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t \mathbb{1}_{L(s) \leq m} ds. \quad (2.5.7)$$

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

2.5.1. Design a queueing system to show that average number of jobs in the system as seen by the server can be very different from what the customers see.

2.5.2. Consider a discrete-time model of a queueing system, as we developed in Section 1.3. The average number of customers that *see upon arrival* more than m customers in the system cannot be defined as (2.5.5). Provide a better definition.

2.6 LEVEL CROSSING AND BALANCE EQUATIONS

Consider a system at which customers arrive and depart in single entities, such as customers in a shop or jobs at some machine. If the system starts empty, then we know that the number $L(t)$ in the system at time t is equal to $A(t) - D(t)$. To illustrate:

$$\longrightarrow A(t) \longrightarrow \boxed{L(t) = A(t) - D(t)} \longrightarrow D(t) \longrightarrow$$

What goes in the box (i.e., $A(t)$) minus what has already gone out (i.e., $D(t)$) must still be in the box, hence $L(t) = A(t) - D(t)$.

Let us denote an arrival as an ‘up-crossing’ and a departure as a ‘down-crossing’. Then, clearly $L(t)$ is the number of up-crossings up to time t minus the number of down-crossings up to time t . If $L(t)$ remains finite, or, more generally, $L(t)/t \rightarrow 0$ as $t \rightarrow \infty$, then it must be that

$$\lambda = \lim_{t \rightarrow \infty} \frac{A(t)}{t} = \lim_{t \rightarrow \infty} \frac{D(t) + L(t)}{t} = \lim_{t \rightarrow \infty} \frac{D(t)}{t} + \lim_{t \rightarrow \infty} \frac{L(t)}{t} = \delta.$$

Hence, when $L(t)/t \rightarrow 0$, the *up crossing rate* $\lim_{t \rightarrow \infty} A(t)/t = \lambda$ is equal to the *down-crossing rate* $\lim_{t \rightarrow \infty} D(t)/t = \delta$. We will generalize these notions of up- and downcrossing in this section to derive the *stationary*, or *long-run time average* or *steady-state*, distribution $p(n)$ that the system contains n jobs.

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 *crosses level* n at time t when its state changes from n to $n + 1$, either ‘from below’ due to an arrival, or ‘from above’ due to a departure, cf. Figure 9.

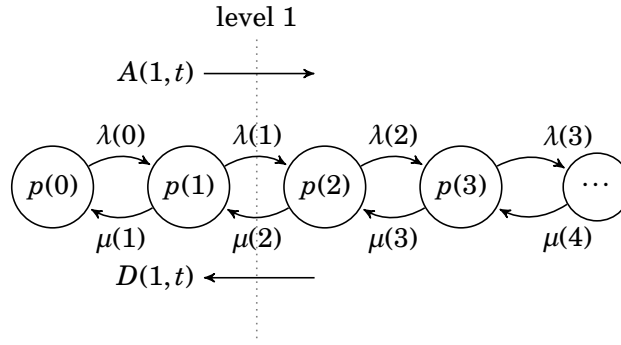


Figure 9: $A(1, t)$ counts the number of jobs up to time t that saw 1 job in the system upon arrival, and right after such arrivals the system contains 2 jobs. Thus, each time $A(1, t)$ increases by one, level 1 (the dotted line separating states 1 and 2) is crossed from below. Similarly, $D(1, t)$ counts 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) the number of times it is crossed from above. (We introduce $\lambda(n)$, $\mu(n)$ and $p(n)$ below.)

To establish the section’s main result Eq. (2.6.5) 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. Once we have the proper definitions, the above result will follow straightaway. Figure 18 at the end of the section summarizes all concepts we develop here.

LEVEL CROSSING Define

$$A(n, t) = \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq t} \mathbb{1}_{L(A_k-) = n} \quad (2.6.1a)$$

as the number of arrivals up to time t that saw n customers in the system at their arrival.

Next, let

$$Y(n, t) = \int_0^t \mathbb{1}_{L(s) = n} ds \quad (2.6.1b)$$



Figure 10: 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 L .)

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

$$p(n, t) = \frac{1}{t} \int_0^t \mathbb{1}_{L(s)=n} ds = \frac{Y(n, t)}{t}, \quad (2.6.1c)$$

be the fraction of time that $L(s) = n$ in $[0, t]$. Figure 10 illustrates the relation between $Y(n, t)$ and $A(n, t)$.

2.6.1. 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?)

2.6.2 (Continuation of 2.6.1). Find an expression for $Y(n, t)$.

Define also the limits:

$$\lambda(n) = \lim_{t \rightarrow \infty} \frac{A(n, t)}{Y(n, t)}, \quad p(n) = \lim_{t \rightarrow \infty} p(n, t), \quad (2.6.2)$$

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.

2.6.3 (Continuation of 2.6.2). 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 9 shows how $A(n, t)$ and $\lambda(n)$ relate to $D(n + 1, t)$ and $\mu(n)$.

2.6.4 (Continuation of 2.6.3). 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 $|A(n, t) - D(n, t)| \leq 1$ (think about this). From this observation it follows immediately that

$$\lim_{t \rightarrow \infty} \frac{A(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{D(n, t)}{t}. \quad (2.6.3)$$

With this equation we can obtain two nice and fundamental identities. The first we develop now; the second follows in Section 2.9.

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

$$\lim_{t \rightarrow \infty} \frac{A(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{A(n, t)}{Y(n, t)} \frac{Y(n, t)}{t} = \lambda(n)p(n), \quad (2.6.4a)$$

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

$$\lim_{t \rightarrow \infty} \frac{D(n, t)}{t} = \lim_{t \rightarrow \infty} \frac{D(n, t)}{Y(n + 1, t)} \frac{Y(n + 1, t)}{t} = \mu(n + 1)p(n + 1). \quad (2.6.4b)$$

Combining this with (2.6.3) we arrive at the *level crossing equations*

$$\lambda(n)p(n) = \mu(n + 1)p(n + 1). \quad (2.6.5)$$

2.6.5 (Continuation of 2.6.4). Compute $\lambda(n)p(n)$ for $n \geq 0$, and check $\lambda(n)p(n) = \mu(n + 1)p(n + 1)$.

Result (2.6.5) turns out to be exceedingly useful, as will become evident from Section 2.7 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

$$p(n + 1) = \frac{\lambda(n)}{\mu(n + 1)} p(n). \quad (2.6.6)$$

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

$$p(n + 1) = \frac{\lambda(n)\lambda(n - 1) \cdots \lambda(0)}{\mu(n + 1)\mu(n) \cdots \mu(1)} p(0). \quad (2.6.7)$$

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*

$$G = 1 + \sum_{n=0}^{\infty} \frac{\lambda(n)\lambda(n - 1) \cdots \lambda(0)}{\mu(n + 1)\mu(n) \cdots \mu(1)}. \quad (2.6.8)$$

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 (2.6.5), 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 average number of items becomes

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

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

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

2.6.6. Derive $E[L] = \sum_{n=0}^{\infty} np(n)$ from (2.5.6).

Finally, the following two exercises show that level crossing arguments extend well beyond the queueing systems modeled by Figure 9.

2.6.7. 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.

2.6.8. Use the graph of exercise 2.6.7 and a level crossing argument to express the steady-state probabilities $p(n), n = 0, \dots, 3$ in terms of λ and μ .

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. Figure 11. 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.

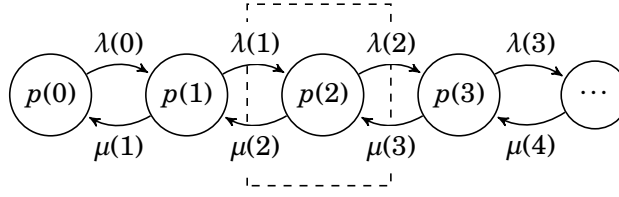


Figure 11: For the balance equations we count how often a box around a state is crossed from inside and outside. On the long run the entering and leaving rates should be equal. For the example here, the rate out is $p(2)\lambda(2) + p(2)\mu(2)$ while the rate in is $p(1)\lambda(1) + p(3)\mu(3)$.

INTERPRETATION The definitions in (2.6.1) 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 Eq. (2.6.1a) 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.

2.7 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 (2.6.6) 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 $N_\lambda(t)$ and the departures (provided the number $L(t)$ in the system is positive) are generated according to the Poisson process $N_\mu(t)$. 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, (2.6.6) 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 (2.6.8) and (1.1.1d) that

$$p(0) = 1 - \rho, \quad p(n) = (1 - \rho)\rho^n. \quad (2.7.1)$$

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,

$$E[L] = \frac{\rho}{1 - \rho}, \quad V[L] = \frac{\rho}{(1 - \rho)^2}, \quad P(L > n) = \rho^{n+1}. \quad (2.7.2)$$

2.7.1. Derive (2.7.2) by differentiating the left-hand and right-hand side of the standard formula for a geometric series: $\sum_{n=0}^{\infty} \rho^n = (1 - \rho)^{-1}$ for $|\rho| < 1$.

Let us interpret expression (2.7.2). 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.)

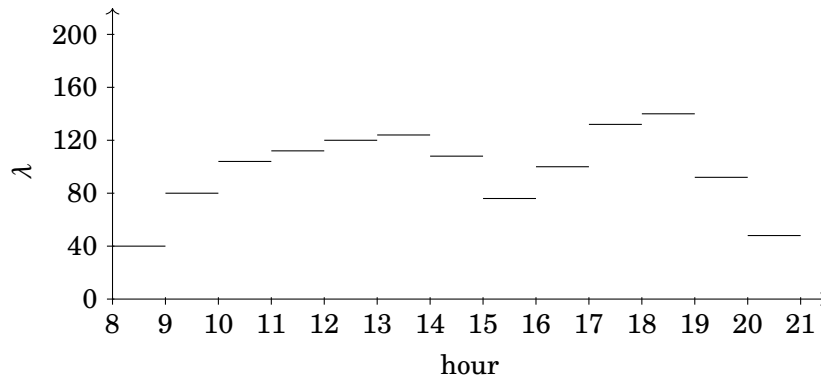


Figure 12: A demand profile of the arrival rate λ modeled as constant over each hour.

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 small probability to go the supermarket on 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. Figure 12. 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 Exercise 2.8.3 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 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 :

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

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 Figure 12 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 not 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 + indicate a working hour and – a break of an hour. Then generate shift types for each of these plans with starting times 8am, 9am, 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 .

2.8 $M(n)/M(n)/1$ 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 (2.6.6). For these queueing systems we just present the results. In the exercises we ask you to derive the formulas—the main challenge is not to make computational errors.

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 the jobs in the system. Specifically, consider the departure time D_k of the k th job, so that $A_{A(D_k)+1}$ is the arrival time of the next job. If $L(D_k) = n$, then we require that for every k

$$P(A_{A(D_k)+1} - D_k \leq x) = 1 - e^{-f(n)\lambda x},$$

for some function $f : \mathbb{N} \rightarrow [0, \infty)$. Next, since $D_{D(A_k)+1}$ is the time until the next departure after A_k , we assume for all k ,

$$P(D_{D(A_k)+1} - A_k \leq x) = 1 - e^{-g(n)\mu x},$$

if $L(A_k) = n$ (not $L(A_k -) = n$) for some function $g : \mathbb{N} \rightarrow [0, \infty)$.

2.8.1. 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.

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

2.8.3. 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 $M/M/q$ 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 $\rho \uparrow 1$.

2.8.4. 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.

2.8.5. 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.

2.8.6. 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.

2.8.7. 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 2.9 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 in Section 2.9.

2.8.8. 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 ?

2.9 POISSON ARRIVALS SEE TIME AVERAGES

Suppose the following limit exists:

$$\pi(n) = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{k=1}^m \mathbb{1}_{L(A_k-) = n}, \quad (2.9.1)$$

then $\pi(n)$ is the long-run fraction of jobs that observe n customers in the system at the moment an arbitrary job arrives. It is natural to ask whether $\pi(n)$ and $p(n)$, as defined by (2.6.2), 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 Ex. 2.9.4 for a proof)

$$\begin{aligned} \pi(n) &= \lim_{t \rightarrow \infty} \frac{1}{A(t)} \sum_{k=1}^{A(t)} \mathbb{1}_{L(A_k-) = n} = \lim_{t \rightarrow \infty} \frac{1}{A(t)} \sum_{k=1}^{\infty} \mathbb{1}_{A_k \leq t, L(A_k-) = n} \\ &= \lim_{t \rightarrow \infty} \frac{A(n, t)}{A(t)}, \end{aligned} \quad (2.9.2)$$

where we use (2.6.1a) in the last row. But, with (2.3.1),

$$\frac{A(n, t)}{t} = \frac{A(t)}{t} \frac{A(n, t)}{A(t)} \rightarrow \lambda \pi(n), \quad \text{as } t \rightarrow \infty, \quad (2.9.3)$$

while by (2.6.4),

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

Thus

$$\lambda \pi(n) = \lambda(n) p(n). \quad (2.9.4)$$

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 average is equal to the time-average. 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 typically called PASTA: Poisson Arrivals See Time Averages. Thus, for the $M/M/1$ queue in particular,

$$\pi(n) = p(n) = (1 - \rho) \rho^n.$$

2.9.1. Show for the case of Ex. 2.6.2 that $\pi(0) = 1$ and $\pi(n) = 0$, for $n > 0$.

2.9.2. Check that (2.9.4) holds for the system of Ex. 2.9.1.

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 (2.9.2),

$$\delta(n) = \lim_{t \rightarrow \infty} \frac{D(n, t)}{D(t)} \quad (2.9.5)$$

as the long-run fraction of jobs that leave n jobs *behind*. From Eq. (2.6.3)

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

Taking limits at the left and right, and using (2.3.3), we obtain for (queueing) systems in which customers arrive and leave as single units that

$$\lambda \pi(n) = \delta \delta(n). \quad (2.9.6)$$

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

$$\lambda = \delta \iff \pi(n) = \delta(n). \quad (2.9.7)$$

2.9.3. When $\lambda \neq \delta$, is $\pi(n) \geq \delta(n)$?

2.9.4. There is a subtle problem in the transition from (2.9.1) to (2.9.2) and the derivation of (2.9.3): $\pi(n)$ 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 (2.9.2) is valid.

With the PASTA property we can determine the distribution of the inter-departure times of the $M/M/1$ queue. Observing that in a network of queues the departures from one queueing station form the arrivals at another station, we can use this result to analyze networks of queues

2.9.5. Try to prove *Burke's law* which states that the departure process of the $M/M/1$ queue is a Poisson process with rate λ .

2.10 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*:

$$E[L] = \lambda E[W]. \quad (2.10.1)$$

Ex. 2.10.1 provides a proof of this under some simple conditions. In the forthcoming sections we will apply Little's law often. Part of the usefulness of Little's law is that it applies to all input-output systems, whether it is a queueing system or an inventory system or some much more general system.

We start with defining a few intuitively useful concepts. From (1.5.7), 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,

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

With Fig. 4 we can relate 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 $s \leq T$,

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

2.10.1. 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.

2.10.2. 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 simple 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.,

$$E[W] = E[W_Q] + E[S]. \quad (2.10.2)$$

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

$$E[W_Q] = E[L] E[S]. \quad (2.10.3)$$

2.10.3. Use Little's law to show for the $M/M/1$ queue that

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

2.10.4. Why is (2.10.3) *not* true in general for the $M/G/1$ queue?

The following problems show how combining PASTA with Little's law allows the analysis of some non-trivial practical queueing situations.

2.10.5 (Hall 5.10). A repair/maintenance facility would like to determine how many employees should be working in its tool crib. The service time is exponential, with mean 4 minutes, and customers arrive by a Poisson process with rate 28 per hour. The customers are actually maintenance workers at the facility, and are compensated at the same rate as the tool crib employees. What is $E[W]$ for $c = 1, 2, 3$, or 4 servers? How many employees should work in the tool crib?

2.10.6 (Hall 5.22). 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?

2.11 $M^X/M/1$ QUEUE: EXPECTED WAITING TIME

It is not always the case that jobs arrive in single units, they can also arrive in batches. For instance, when a car and or 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 $M^X/M/1$. In this section we derive the expressions for the load and the expected waiting time and queue length of the $M^X/M/1$ queue.

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, i.e., the number of items that job k brings into the system. We assume that $\{B_k\}$ is a sequence of independent discrete random variables each distributed as the generic random variable B . Let the pmf be $P(B = k) = f(k)$, and $G(k)$ the survivor function. The service time of each item $E[S] = 1/\mu$. Thus, the average time to serve the entire batch is $E[B]E[S]$. Define the load as

$$\rho = \lambda E[B]/\mu.$$

Henceforth we require that $\rho < 1$.

2.11.1. Use the renewal reward theorem to explain that work arrives at rate $\lambda E[B]$.

The aim of the remainder of the section is to derive two cornerstones of queueing theory. The first is the expected time an item spends in queue:

$$E[W_Q] = \frac{1 + C_s^2}{2} \frac{\rho}{1 - \rho} E[B] E[S] + \frac{1}{2} \frac{\rho}{1 - \rho} E[S], \quad (2.11.1)$$

where C_s^2 is the SCV of the batch size distribution. We obtain the second by applying (2.10.3), which evidently also applies here, to the above to find that the expected number of items in the system is

$$E[L] = \frac{E[W_Q]}{E[S]} = \frac{1 + C_s^2}{2} \frac{\rho}{1 - \rho} E[B] + \frac{1}{2} \frac{\rho}{1 - \rho}. \quad (2.11.2)$$

Thus, to compute the average number of items in the system, we only need to know the first and second moment (or the variance) of the batch size B .

2.11.2. Show that when the batch size is 1, the expression $E[L(M^X/M/1)]$, i.e., the system length for the $M^X/M/1$ queue, reduces to $E[L(M/M/1)]$, i.e., the system length for the $M/M/1$ queue. Realize the importance of such checks.

2.11.3. If the batch size is p geometrically distributed, what is $E[L]$?

2.11.4. 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.

Let us now focus on deriving expressions for the expected time batches wait in queue $E[W_Q]$. Assume that a batch upon arrival joins the end of the queue (if present), and once the queue is cleared, the entire batch moves from the queue to the server. Thus, all items in one batch spends 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. Let $E[L_S]$ be the average number of items of a batch at the server.

2.11.5. Show that

$$E[W_Q] = \frac{E[L_S]}{1-\rho} E[S].$$

Clearly, it remains to find an expression for $E[L_S]$; for this we can again use the renewal reward theorem. Let $L_S(s)$ be the number of items (of the batch in service) at the server, so that

$$Y_i(t) = \int_0^t \mathbb{1}_{L_S(s)=i} ds$$

is the total time up to t there are i items at the server.

2.11.6. Let \tilde{A}_k be the moment the k th batch moves to the server and D_k its departure time. Use Fig. 13 to show that

$$\int_{\tilde{A}_k}^{D_k} \mathbb{1}_{L_S(s)=i} ds = S_{k,i} \mathbb{1}_{B_k \geq i},$$

where $S_{k,i}$ is the service time of the i th item of this batch.

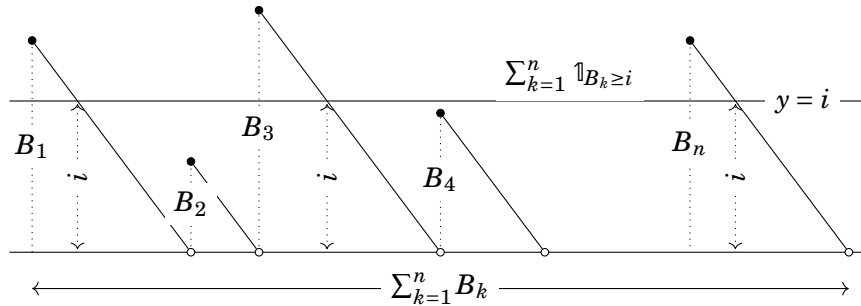


Figure 13: The remaining job size as a function of time. The total number of service periods, which is equal to the number of items arrived, is $\sum_{k=1}^n B_k$. A batch crosses the line $y = i$ iff it contains at least i items. Thus, during the service of a batch with i or more items, there is precisely one period during which the i -th item of a batch is waiting in queue. Consequently, $\sum_{k=1}^n \mathbb{1}_{B_k \geq i}$ is the number of periods in which there are precisely i items waiting at the server. The fraction of periods there are i items is therefore $\sum_{k=1}^n \mathbb{1}_{B_k \geq i} / \sum_{k=1}^n B_k$.

2.11.7. Use the previous exercise to show that

$$Y_i(D_n) = \sum_{k=1}^n \mathbb{1}_{B_k \geq i} S_{k,i}.$$

2.11.8. Now use the renewal reward theorem to show that the (time-average) fraction of time there are i items at the server is equal to

$$P(L_S = i) = \lambda E[S] G(i-1) = \rho \frac{G(i-1)}{E[B]}.$$

With the above exercises we conclude that

$$E[L_S] = \sum_{i=0}^{\infty} i P(L_S = i) = \frac{\rho}{E[B]} \sum_{i=1}^{\infty} i G(i-1).$$

2.11.9. Brush up the above expression for $E[L_S]$ to arrive at (2.11.1).

2.11.10. 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?

2.12 M/G/1 QUEUE: EXPECTED WAITING TIME

Let's focus on one queue in a supermarket, served by one cashier, and assume that customers do not jockey, i.e., change queue. What can we say about the average waiting time in queue if service times are not exponential, like in the $M/M/1$ queue, but have a more general distribution? One of the celebrated results of queueing theory is the Pollaczek-Khinchine formula by which we can compute the average waiting time formula for the $M/G/1$ queue. In this section we derive this result by means of sample path arguments.

To find an expression for $E[W_Q]$ we need the concept of expected *remaining service time* $E[S_r]$ which is defined as the expected time it takes to complete the job in service at the time a job arrives. In Section 2.12 we give a precise meaning to this idea.

2.12.1. Show for the $M/G/1$ queue that the expected time in queue is

$$E[W_Q] = E[S_r] + E[L_Q] E[S]. \quad (2.12.1)$$

2.12.2. ($M/G/1$) Use the PASTA property to show that

$$E[S_r] = \rho E[S_r | S_r > 0]. \quad (2.12.2)$$

2.12.3. It is an easy mistake to think that $E[S_r] = E[S]$ when service times are exponential. Why is this wrong?

2.12.4. What would you guess for $E[S_r | S_r > 0]$ for the $M/D/1$ queue?

2.12.5. Use Exercise 2.12.1 and Little's law to derive for the $M/G/1$ queue that

$$E[W_Q] = \frac{E[S_r]}{1 - \rho}. \quad (2.12.3)$$

Recall that Eq. (2.12.3) states that

$$E[W_Q] = \frac{E[S_r]}{1 - \rho}.$$

It remains to compute the average remaining service time $E[S_r]$ for generally distributed service times. For this, we use the renewal reward theorem, again.

For the $M/M/1$ queue the situation becomes significantly simpler as then the service times are exponential, hence memoryless, which implies that $E[S_r | S_r > 0] = E[S]$.

Consider the k th job of some sample path of the $M/G/1$ queueing process. This job requires S_k units of service, let its service time start at time \tilde{A}_k so that it departs the server at time $D_k = \tilde{A}_k + S_k$.

2.12.6 (Δ). Use Figure 14 to explain that job k the remaining service time at time s is given by

$$R_k(s) = (D_k - s) \mathbb{1}_{\tilde{A}_k \leq s < D_k}.$$

2.12.7 (Δ). 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 .

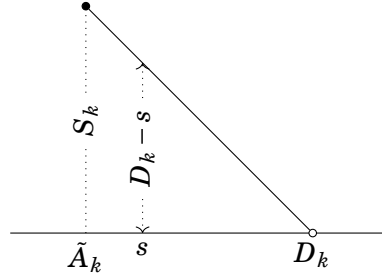


Figure 14: Remaining service time.

2.12.8 (\boxtimes). Use the renewal reward theorem to prove that

$$E[S_r] = \frac{\lambda}{2} E[S^2]. \quad (2.12.4)$$

2.12.9 (\boxtimes). Use $C_s^2 = V[S]/(E[S])^2$ to show that $\lambda E[S^2] = (1 + C_s^2)\rho E[S]$

With the above exercises we have obtained the fundamentally important *Pollaczek-Khinchine formula* for the average waiting time in queue:

$$E[W_Q] = \frac{E[S_r]}{1 - \rho} = \frac{1}{2} \frac{\lambda E[S^2]}{1 - \rho} = \frac{1 + C_s^2}{2} \frac{\rho}{1 - \rho} E[S]. \quad (2.12.5)$$

The problems below will illustrate how useful this result is.

2.12.10 (Δ). Show from Eq. (2.12.4) that

$$E[S_r | S_r > 0] = \frac{E[S^2]}{2E[S]}. \quad (2.12.6)$$

2.12.11 (Δ).

$$E[S_r | S_r > 0] = \frac{E[S^2]}{2E[S]} \implies E[S_r | S_r > 0] = \frac{\alpha}{2}$$

when $V[S] = 0$.

2.12.12 (\boxtimes). Show that $E[S_r | S_r > 0] = \alpha/3$ when $S \sim U[0, \alpha]$.

2.12.13. (Δ) Show that $E[S_r | S_r > 0] = \mu^{-1}$ when $S \sim \text{Exp}(\mu)$.

2.12.14 (Δ). Show that when services are exponential, the expected waiting time $E[W_Q(M/G/1)]$ reduces to $E[W_Q(M/M/1)]$.

2.12.15 (▣). Compute $E[W_Q]$ and $E[L]$ for the $M/D/1$ queue. Assume that the service time is always T . Compare $E[L(M/D/1)]$ to $E[L(M/M/1)]$ where the mean service time is the same in both cases.

2.12.16 (▣). Compute $E[L]$ for the $M/G/1$ queue with $S \sim U[0, \alpha]$.

2.12.17 (▣). 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 $E[L_Q]$, $E[L]$, $E[W_Q]$, $E[W]$.

2.12.18 (▣). Consider a workstation with just one machine. Jobs arrive, roughly, in accordance with 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?

2.12.19 (▣). (Hall 5.16) The manager of a small firm would like to determine which of two people to hire. One employee is fast, on average, but is 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 $E[L]$?

2.12.20 (▣). Show that for the $M/G/1$ queue, the expected idle time is $E[I] = 1/\lambda$.

2.12.21 (▣). What is the utilization of the $M/G/1/1$ queue?

2.12.22 (▣). For the $M/G/1/1$ queue what is the fraction of jobs rejected (hence, what is the fraction of accepted jobs)?

2.12.23 (▣). Why is the fraction of lost jobs at the $M/G/1/1$ queue not necessarily the same as for a $G/G/1/1$ queue with the same load?

2.13 $M^X/M/1$ QUEUE LENGTH DISTRIBUTION

In Sections 2.11 and 2.12 we established the Pollaczek-Khinchine formula for the waiting times of the $M^X/M/1$ queue and $M/G/1$ queue, respectively. To compute more difficult performance measures, for instance the loss probability $P(L > n)$, we need expressions for the stationary distribution $\pi(n) = P(L = n)$ of the number of jobs in the system. Here we present a numerical, recursive, scheme to compute these probabilities.

To find $\pi(n)$, $n = 0, 1, \dots$, we turn again to level-crossing arguments. However, the reasoning that lead to the level-crossing equation (2.6.3) need 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 it's 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, at the arrival of job k , all levels between states 3 and 23 are crossed.

The left panel in Figure 15 demonstrates the up- and down-crossings in more general terms. Level n can be up-crossed from below from many states, in fact from any level m , $0 \leq m < n$. However, it can only be down-crossed from state $n + 1$.

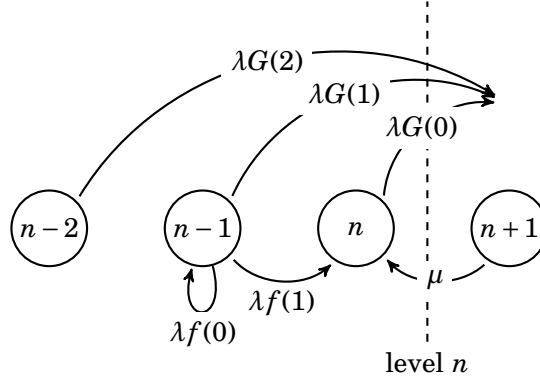


Figure 15: Level crossing of level n . Observe that when the system is in state $n-2$, the arrival of any batch larger than 2 ensures that level n is crossed from below. The rate at which such events happen is $\lambda G(2)$. Similarly, in state $n-1$ the arrival of any batch larger than one item ensures that level n is crossed, and this occurs with rate $\lambda G(1)$, and so on.

As always with level-crossing arguments, we turn to counting how often level n is up-crossed and down-crossed as a function of time. The down-crossing rate is easy: there is just one arrow from right to left in Figure 15 to down-cross level n , namely from $n+1$ to n . Hence, the down-crossing rate is exactly the same as for the $M/M/1$ queue, i.e., (2.6.4b).

Counting up-crossings requires quite some more work. Observe that $\mathbb{1}_{L(A_k-) \leq n} = 1$ only when the k th job sees n or less items in the system, and $\mathbb{1}_{L(A_k) > n} = 1$ only after the k th arrival the system contains more than n items. Thus, $\mathbb{1}_{L(A_k-) \leq n} \mathbb{1}_{L(A_k) > n} = 1$ iff the k th arrival generates an up-crossing of level n .

From Figure 15 we see that an up-crossing can be decomposed into:

$$\begin{aligned} & \mathbb{1}_{L(A_k-) \leq n} \mathbb{1}_{L(A_k) > n} \\ &= \mathbb{1}_{L(A_k-) = n} \mathbb{1}_{B_k > 0} + \mathbb{1}_{L(A_k-) = n-1} \mathbb{1}_{B_k > 1} + \cdots + \mathbb{1}_{L(A_k-) = 0} \mathbb{1}_{B_k > n} \\ &= \sum_{m=0}^n \mathbb{1}_{L(A_k-) = m} \mathbb{1}_{B_k > n-m}. \end{aligned}$$

In other words, sample paths that up-cross level n require that any job that sees m ($m \leq n$) in the system upon arrival must bring a batch larger than $n-m$ items.

In view of the above, 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 $n-m$.

2.13.1 (A). Show that $A(n, n, t) = A(n, t)$, where $A(n, t)$ is defined by Eq. (2.6.1a).

As in Section 2.6, we are primarily interested in long-run averages. For this purpose, observe that we can write

$$\frac{A(m, n, t)}{t} = \frac{A(t)}{t} \frac{A(m, t)}{A(t)} \frac{A(m, n, t)}{A(m, t)}. \quad (2.13.1)$$

By the assumptions of Section 2.9, $A(t)/t \rightarrow \lambda$ and $A(m, t)/A(t) \rightarrow \pi(m)$. Now, provided the limit exists, we define

$$\lim_{t \rightarrow \infty} \frac{A(m, n, t)}{A(m, t)} = \lim_{t \rightarrow \infty} \frac{\sum_{k=1}^{A(t)} \mathbb{1}_{L(A_k-) = m, B_k > n-m}}{\sum_{k=1}^{A(t)} \mathbb{1}_{L(A_k-) = m}} = P(B > n - m | L(A-) = m), \quad (2.13.2)$$

where the random variable $L(A-)$ denotes the number in the system seen by an arbitrary arrival.

2.13.2 (\triangle). Show that $P(B > n - m | L(A-) = m) = P(B > n - m)$.

By the above exercise,

$$\lim_{t \rightarrow \infty} \frac{A(m, n, t)}{A(m, t)} = P(B > n - m) = G(n - m).$$

By combining the above and making the usual assumptions about the existence of all limits involved we find

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

2.13.3 (\triangle). Provide an interpretation of the above result in terms of a thinned Poisson arrival process.

The last step is to relate the up- and down-crossing rates. Clearly, $\sum_{m=0}^n A(m, n, t)$ is the total number of times level n is up-crossed up to time t . By level crossing,

$$\sum_{m=0}^n A(m, n, t) \approx D(n, t).$$

Thus, taking the limit $t \rightarrow \infty$ in this equation, we conclude that

$$\lambda \sum_{m=0}^n G(n - m) \pi(m) = \mu \pi(n + 1), \quad (2.13.3)$$

where we use PASTA in (2.6.4b) to see that $\mu(n + 1)p(n + 1) = \mu \pi(n + 1)$.

2.13.4 (\triangle). Show that Eq. (2.13.3) reduces to $\mu \pi(n + 1) = \lambda \pi(n)$ for the $M/M/1$ case.

2.13.5 (\boxtimes). With $\alpha = \lambda/\mu$, show that *unnormalized* state probabilities are given by


$$\begin{aligned} \pi(0) &= 1 & \pi(1) &= \alpha \\ \pi(2) &= \alpha^2 + \alpha G(1), & \pi(3) &= \alpha[\alpha^2 + 2\alpha G(1) + G(2)]. \end{aligned}$$

We leave the rest to the computer to continue with this.


It is left to find the normalization constant. As this recursion does not lead to a closed form expression for $\pi(n)$, such as Eq. (2.7.1), we need to use a criterion to stop this iterative procedure. Finding general conditions when to stop is not directly easy, but a pragmatic approach is simple: stop at some (large) number N such that $\pi(k) \ll \pi(0)$ and steadily decreases for $k > N$.³ Then take $G = \sum_{i=0}^N \pi(i)$ as the normalization constant, so that $\pi(0) = 1/G$, $\pi(1) = \alpha/G$, and so on. While this is a practical approach, getting formal bounds on a proper size of N requires more work than we can do here.

Once again with $\pi(n)$ we can compute all performance measures we need, and study the influence on the batch size distribution and λ and μ on the system's performance.

³ An interesting question, why should it decrease monotonically after some, large, N ?

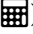
2.13.6 () Why is Eq. (2.9.7), i.e., $\pi(n) = \delta(n)$, not true for the $M^X/M/1$ batch queue?

Let us use recursion Eq. (2.13.3) for $\pi(n)$ to derive an expression for the expected number of units of work $E[L]$ in the system.


2.13.7 () Show that


$$\mu E[L] = \mu \sum_{n=0}^{\infty} n \pi(n) = \lambda \frac{E[B^2]}{2} + \lambda E[B] E[L] + \lambda \frac{E[B]}{2}. \quad (2.13.4)$$

With this we can check a result of Section 2.11.

2.13.8 () Use Eq. (2.13.4) and the definition $\rho = \lambda E[B]/\mu$ to show that


$$(1 - \rho)E[L] = \frac{\lambda}{\mu} \frac{E[B^2]}{2} + \frac{\rho}{2}.$$

2.13.9 () Implement the recursion (2.13.3) in a computer program for the case $f(1) = f(2) = f(3) = 1/3$ (recall $P(B = k) = f_k$). Take $\lambda = 1$ and $\mu = 3$.

2.13.10 () (Finite queueing systems) 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, it is necessary to specify an acceptance, or equivalently a rejection, policy. 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 otherwise.

Derive a set of recursions, analogous to Eq. (2.13.3), to compute $\pi(n)$ for these three different acceptance rules.

2.13.11 () An interesting extension is to consider a queueing process with batch services, i.e., the $M/M^Y/1$ queue. Constructing a recursion for the steady-state probabilities $\pi(n)$ for this case is not hard, in fact, mostly analogous to Eq. (2.13.3). However, solving the recursion appears to be quite a bit harder; we will not discuss this further here.

2.14 M/G/1 QUEUE LENGTH DISTRIBUTION

In Section 2.13 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’, which are job departure epochs as we will see below. 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 relation (2.9.7) 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 an i.i.d. sequence of random variables, the sequence $\{Y_k\}$ is 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.

2.14.1 (🖼️). Explain that if the service time is constant and equal to s , then

$$P(Y_k = j | S = s) = e^{-\lambda s} \frac{(\lambda s)^j}{j!}. \quad (2.14.1)$$

2.14.2 (📐). Explain that

$$P(Y_k = j) = \int_0^\infty e^{-\lambda x} \frac{(\lambda x)^j}{j!} dF(x), \quad (2.14.2)$$

where F is the distribution of the service times.

2.14.3 (📐). If S is deterministic and equal to s , show that Eq. (2.14.2) reduces to Eq. (2.14.1).

2.14.4 (📊). If $S \sim \text{Exp}(\mu)$, show that

$$f(j) = P(Y_k = j) = \frac{\mu}{\lambda + \mu} \left(\frac{\lambda}{\lambda + \mu} \right)^j. \quad (2.14.3)$$

2.14.5 (📐). If $S \sim \text{Exp}(\mu)$, show that

$$G(j) = \sum_{k=j+1}^{\infty} f(k) = \left(\frac{\lambda}{\lambda + \mu} \right)^{j+1}. \quad (2.14.4)$$

2.14.6 (🖼️). Design a suitable numerical method to evaluate Eq. (2.14.2) for more general distribution functions F .

Let us concentrate on a down-crossing of level n , see Figure 16; 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

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

2.14.7. Use a similar derivation as in (2.13.2) to show that

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

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

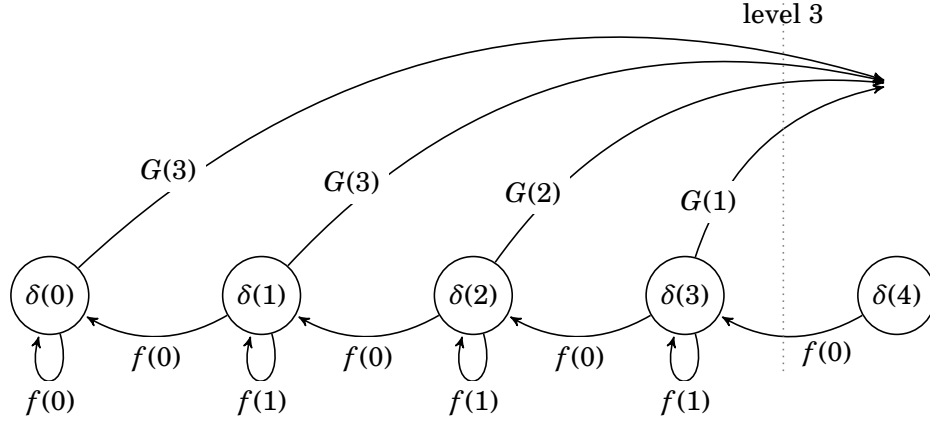


Figure 16: Level 3 is crossed from below with rate $\delta\delta(0)G(3) + \delta\delta(1)G(3) + \cdots \delta\delta(3)G(1)$ and crossed from above with rate $\delta\delta(4)f(0)$.

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

2.14.8 (Δ). 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),

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

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

$$\delta\delta(0)G(n) + \delta \sum_{m=1}^n \delta(m)G(n-m+1). \quad (2.14.5)$$

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 $\pi(n) = \delta(n)$, which follows from (2.9.7) and the fact that the $M/G/1$ queue length process has one-step transitions, we arrive at

$$f(0)\pi(n+1) = \pi(0)G(n) + \sum_{m=1}^n \pi(m)G(n+1-m). \quad (2.14.6)$$

Clearly, we have again obtained a recursion by which we can compute, iteratively, the state probabilities, and follow the approach sketched below Exercise 2.13.5.

2.14.9. Provide the details behind the derivation of Eq. (2.14.5).

2.14.10 (▣). 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 (2.14.6).

APPROXIMATE MODELS

In this chapter we first consider the very useful formula of Sakasegawa to approximate the average waiting time in queue for the $G/G/c$ queue. We then illustrate how to use this formula to estimate waiting time in three examples in which the server is interrupted. In the first case the server has to produce jobs from different families, and there is a switch-over time required to change the production family. Such setups reduce the time the server has available to serve jobs. To reduce the load the server produces in batches of fixed sizes. In the second case, the server requires sometime a small adjustment, for instance, to prevent its quality to degrade below a certain level. Such adjustments are not necessary to perform 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 (the number of jobs served between two interruptions) are constant. In the third example, quality problems or break downs can occur during a job's service. For each case we develop a model to analyze the influence of the interruption duration and frequency on average waiting times.

3.1 $G/G/c$ QUEUE: APPROXIMATIONS

Theory and Exercises

In manufacturing settings it is quite often the case that the arrival process at a station is not Poisson. 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. Hence, in this case, the SCV of the arrivals at the second station $C_{a,2}^2$ is most probably smaller than 1. As a consequence the Pollaczek-Khinchine formula for the $M/G/1$ queue can no longer be reliably used to compute the average waiting times. 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. Thus, the $M/G/1$ waiting time formula should not be naively applied to approximate the average waiting time of the $G/G/1$ queue.

There is no formula as yet by which the average waiting times for the $G/G/1$ queue can be computed; only approximations are available. One such simple and robust approximation is based on the following observation. Recall the waiting time in queue for the $M/M/1$ queue:

$$E[W_Q(M/M/1)] = \frac{\rho}{1-\rho} E[S] = \frac{1_a + 1_s}{2} \frac{\rho}{1-\rho} E[S],$$

where we label the number 1 with a and s . When generalizing this result to the $M/G/1$ queue we get

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

Thus, 1_s in the expression for the $M/M/1$ queue is replaced by C_s^2 in the expression for the $M/G/1$ queue. As a second generalization, Kingman proposed to replace 1_a in this formula by the SCV of the inter-arrival times

$$C_a^2 = \frac{V[X]}{(E[X])^2},$$

resulting in

$$E[W_Q(G/G/1)] \approx \frac{C_a^2 + C_s^2}{2} \frac{\rho}{1-\rho} E[S].$$

This formula is reasonably accurate; for related expressions we refer to Bolch et al. [2006] and Hall [1991]. With Little's law we can compute $E[L_Q]$ from the above. Moreover, $E[W] = E[W_Q] + E[S]$, and so on, cf., Section ??.

It is crucial to memorize the *scaling* relations that can be obtained from the $G/G/1$ waiting time formula. 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 very 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. $E[W_Q] \sim E[S]$. Thus, working in smaller job sizes reduces the waiting time as well. The average queue length does not decrease by working with smaller batches, but jobs are more 'uniformly spread' over the queue. This effect lies behind the idea of 'lot-splitting', i.e., rather than process large jobs, split jobs into multiple small jobs (assuming that setup times are negligible), so that the waiting time per job can be reduced.

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.

For the $G/G/c$ queue we can use Sakasegawa's approximation, Sakasegawa [1977],

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

to estimate the time in queue, where

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

is the load of the station, not of the individual machines. We refer to Hopp and Spearman [2008] for a discussion of this formula and its many applications.

Even though the above results are only approximate, they prove to be exceedingly useful when designing queueing systems and analyzing the effect of certain changes, in particular changes in capacity, variability and service times.

3.1.1 (A). Show that the approximation (3.1.1) reduces to the result known for the $M/M/1$ and $M/G/1$ queues.

3.1.2 (A). Is Eq. (2.12.1) also valid for the $G/G/1$ queue? Why (not)?

3.1.3 (\triangle). Consider a queue with n servers, with generally distributed inter-arrival times, generally distributed service times, and the system can contain at most K customers, i.e., the $G/G/n/K$ 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

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

3.1.4 (\triangle). Consider a single-server queue at which every minute a customer arrives, precisely at the first second. Each customer requires precisely 50 seconds of service. What are ρ , $E[L]$, C_a^2 , and C_s^2 ?

3.1.5 (\boxtimes). Consider the same single-server system as in the previous exercise, but now the customer service time is stochastic: with probability 1/2 a customer requires 1 minute and 20 seconds of service, and with probability 1/2 the customer requires only 20 seconds of service. What are ρ , C_a^2 , and C_s^2 ?

It is crucial to remember from the above exercises that knowledge of the utilization is not sufficient to characterize the average queue length.

3.1.6 (\triangle). For the $G/G/1$ queue, prove that the fraction of jobs that see n jobs in the system is the same as the fraction of departures that leave n jobs behind. What condition have you used to prove this?

3.1.7 (\boxtimes). (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.

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 your previous answer is smaller.

Clearly, Kingman's equation requires an estimate of the SCV C_a^2 of the inter-arrival times and the SCV C_s^2 of the service times. 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, mostly just the date or the hour, perhaps, that a customer arrived. Hence, it is often not possible to estimate C_a^2 and C_s^2 from information that is available. However, when for instance the number of arrivals per day have been logged for some time so that we know $\{a_n, n = 1, \dots, N\}$ 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 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

Theory and Exercises

With the $G/G/1$ waiting time formula (3.1.1) we can compute, approximately, the waiting time in queue for many non-trivial queueing situations. In this section we focus on the effect of change-overs, or setups. Consider, for instance, a machine that paints red and blue bikes. When the machine requires a color change, a clean-up time is necessary. As we will see it is necessary in such situations to produce in batches. Other examples are ovens that need warm up or cool down times when different item types require different temperatures. In service settings, when servers have to move from a part of a building to another, the time spend moving cannot be spent on serving customers.

Specifically, we analyze the following queueing situation. There are two job families, e.g., red and blue. Jobs arrive at rate λ_r and λ_b , respectively, so that the arrival rate of jobs is $\lambda = \lambda_b + \lambda_r$. For ease we assume that the job's service time, S_0 , has the same distribution for both colors. The change-over time is given by a random variable R , which is independent of the normal job service times.

Jobs of each color are assembled into batches of size B . Once a batch is complete, the batch enters a queue (of batches). Once a batch reaches the head of the queue, the machine performs a setup, and then starts processing each job individually until the batch is complete. If there is another batch in queue, a new setup time is required. Otherwise the machine just switches off. 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 finish.

We analyze in steps the total average time a job spends in the system.

First we consider the time it takes to form a batch.

3.2.1 (A). Show that the total time to form a red batch is $(B-1)/\lambda_r$. Hence, the average time a red job spends waiting until the batch is complete is

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

Now that we know how long jobs spend to form batches, we turn to finding an estimate for the average time a batch has to spend in queue, for which we use Eq. (3.1.1). Recall that for

this formula, we need the arrival rate, the average service time and the SCVs. These elements we will compute now.

It is evident that the rate at which batches arrive is

$$\lambda_B = \frac{\lambda}{B},$$

since both job colors have the same batch size.

Observe next that the machine not only serves jobs, part of the time it is occupied with setups. This leads to the idea to *incorporate* the effects of the setup times in the service times. For this we distinguish between a job's *net service time* S_0 and its *effective processing time* S which also include setup times.

3.2.2 (A). Show that

$$E[S] = E[S_0] + \frac{E[R]}{B}.$$

Now that the batch arrival rate and the service time per batch are known, the load can be written as

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

where the first equality has the interpretation of the batch arrival rate times the work per batch, while the second is the job arrival rate times the effective work per job.

3.2.3 (A). Show that the requirement $\rho < 1$ leads to the following constraint on the minimal batch size B

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

The next element is to find the SCVs. To obtain $C_{a,B}^2$, i.e., the SCV of the inter-arrival times of the batches, recall that jobs are first assembled into batches, and then these batches are sent to the queue.

3.2.4 (A). Show that

$$C_{a,B}^2 = \frac{C_a^2}{B},$$

with C_a^2 the SCV of inter-arrival times of individual jobs.

The last element is to find the SCV $C_{s,B}^2$ of the service times of the batches.

3.2.5 (A). Show that

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

Finally, when the batch is taken into service, there can be various rules to determine when the job's service finished. If the job has to wait until all jobs in the batch are served, the time a job spends at the server is $B E[S_0] + E[R]$.

3.2.6 (A). In our model we assume that jobs can leave right after being served. Show for this case that the expected time until a job leaves the server is

$$E[R] + \frac{B-1}{2} E[S_0] + E[S_0].$$

Clearly, we now have all elements to compute the average time in the system. Let's illustrate this.

3.2.7 (📐). Jobs arrive at $\lambda = 3$ per hour at a machine with $C_a^2 = 1$; service times are exponential with average 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. Suppose the batch size $B = 30$ jobs. What is the minimal batch size? What is the average time a red job spends in the system?

In summary, to find the average queueing time in the system, we need to find the arrival rate, the effective service times and the SCVs, so that we can fill in the $G/G/1$ waiting time formula. The main idea is to incorporate the setup times into the job service times. Observe also that the times to form and process batches are linear functions of the batch size B , while the load is, for small batch sizes, very sensitive to the batch size. Thus, batch sizes should not be too small. Overall, batch sizes need to be tuned to minimize the total average time jobs spend in the system. When the batch sizes are small, the load ρ is near to one (in other words, the server spends a relatively large fraction of its time on setups), so that the queueing times are long, but the times to form a batch are small. If, however, the batch sizes are large, the queueing times will be relatively short, but the times to form and unpack batches will be large.

3.3 NON-PREEMPTIVE INTERRUPTIONS, SERVER ADJUSTMENTS

Theory and Exercises

In Section 3.2 we studied the effect of setup times between batches of jobs. In this model we assumed that the number of jobs between two setups is fixed to the batch size B . In other words, there are no setups between any two jobs, setups are *planned* between B jobs. However, other types of interruptions can occur, such as a machine requiring a small adjustment after just a few jobs. In fact, such random interruptions can happen between any two jobs. As such outages do not interrupt the processing of a job in service, we call this *non-preemptive outages*. In this section we develop a simple model to understand the impact of such outages on the mean time jobs spend in the system.

Let us assume that adjustments occur geometrically distributed between any two jobs with a mean of B jobs between any two repairs. 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.) The adjustments form an i.i.d. sequence of random variables distributed as the common random variable R and independent of S_0 . The main idea is to incorporate the effects of these outages in the mean and SCV of job service times, so that we can use the $G/G/1$ waiting time formula.

Define the *effective processing time* S as the time the server is occupied with processing a job including a potential adjustment, and write S_0 for the net service time of a job, i.e., the service time required to serve just the job.

3.3.1 (📐). Show that the average effective processing time satisfies $E[S] = E[S_0] + E[R]/B$. Conclude that the effective server load including down-times is $\rho = \lambda E[S]$.

The next step is to find an expression for $E[S^2]$ from which $V[S]$ will follow easily.

3.3.2 (📐). Show that

$$E[S^2] = E[S_0^2] + 2 \frac{E[S_0] E[R]}{B} + \frac{E[R^2]}{B}.$$

3.3.3 (▣). Use the above to find that

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

With the above we can compute $C_s^2 = V[S]/(E[S]^2)$ of the effective job processing times. We have now all elements to fill in the $G/G/1$ waiting time formula!

3.3.4 (▣). A machine requires an adjustment with average 5 hours and standard deviation of 2 hours. 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 hour. Interruptions occur on average between 30 jobs. Compute the average waiting time in queue.

Observe that with these formulas we can obtain quantitative insights into the effects of reducing adjustment times, or the variability of these adjustments times. For instance, we might decide to do less adjustments, so that p decreases, but the average outage time (or its variance) may increase as a function of this decision. We now have tools to analyze the consequences of such decisions without needing to actually do the experiments in real life to see the effects.

3.4 PREEMPTIVE INTERRUPTIONS, SERVER FAILURES

Theory and Exercises

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.

Let us assume that a job's normal service time, without interruptions, is given by S_0 . The durations of the interruptions are given by the i.i.d. random variables $\{R_i\}$ and have common mean $E[R]$ and variance $V[R]$. If N interruptions occur, the effective service time will then be

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

Observe that to use the $G/G/1$ waiting time formula it suffices to find expressions for $E[S]$ and $V[S]$. Thus, this will be our task for the rest of the section. We remark in passing that the results and the derivation are of general interest.

We first aim to find an expression for $E[S]$. Write $S_N = \sum_{i=1}^N R_i$ for the total duration of the interruptions, so that the total job duration becomes $S = S_0 + S_N$.

3.4.1 (▴). Suppose that $N = n$, show that $E[S_n] = n E[R]$.

Let $p_n = P(N = n)$; then it is reasonable that $E[S_N] = \sum_{n=0}^{\infty} E[S_n] p_n$. (Compare the definition of $E[f(X)] = \sum_n f(n) p_n$.)

3.4.2 (▣). Use the above to show that $E[S_N] = E[R] E[N]$. (This result is known as *Wald's equation*.)

Thus, with the above,

$$E[S] = E[S_0 + S_N] = E[S_0] + E[R] E[N].$$

To make further progress, we need some additional assumptions. A common assumption is that the time between two interruptions is $\text{Exp}(\lambda_f)$, hence is memoryless. Consequently, the number of interruptions N that occur during the net service time S_0 is Poisson distributed with mean $E[N] = \lambda_f E[S_0]$.

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.

3.4.3 (A). Show that for our model of interruptions,

$$A = \frac{1}{1 + \lambda_f E[R]}$$

3.4.4 (A). Show that

$$E[S] = \frac{E[S_0]}{A} = E[S_0](1 + \lambda_f E[R]).$$

An intuitive way to obtain this result is by noting that A is the fraction of time the server is working. As the total service time of a job is $E[S]$, the net work done is $A E[S]$. But this must be the time needed to do the real job, hence $A E[S] = E[S_0]$.

It is important to realize that

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

hence the load increases due to failures.

We can use similar ideas to derive an expression for the variance of S . The next exercise helps to understand why this derivation is a bit more involved.

3.4.5 (A). Why is $V[S] \neq V[S_0] + V[\sum_{i=0}^N R_i]$?

So let us first consider $E[S^2]$; recall that $V[S] = E[S^2] - (E[S])^2$, and we already know that $E[S] = E[S_0]/A$.

3.4.6 (A). Show that

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

To simplify this, we assume at first that S_0 is known, so that the number of failures that occur during a service time S_0 is Poisson distributed, i.e., $N \sim P(\lambda_f S_0)$.

3.4.7 (A). Show that $E[S_0 \sum_{i=1}^N R_i | S_0] = \lambda_f S_0^2 E[R]$.

3.4.8 (A). Show that $E[\sum_{i=1}^N R_i^2 | S_0] = \lambda_f S_0 E[R^2]$.

3.4.9 (A). Show that $E[\sum_{i=1}^N \sum_{j \neq i} R_i R_j | S_0] = \lambda_f^2 S_0^2 (E[R])^2$.

3.4.10 (▣). Combine the above to see that $E[S^2 | S_0] = \frac{S_0^2}{A^2} + \lambda_f E[R^2] S_0$. From this,

$$E[S^2] = \frac{E[S_0^2]}{A^2} + \lambda_f E[R^2] E[S_0].$$

3.4.11 (▲). Next, show that

$$V[S] = \frac{V[S_0]}{A^2} + \lambda_f E[R^2] E[S_0].$$

3.4.12 (▣). Finally, show that

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

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

If we assume that repair times are exponentially distributed with mean $E[R]$, we can simplify this yet further.

3.4.13 (▣). With the above assumption on the distribution of R , show that

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

Again, we have all elements ready to use the $G/G/1$ waiting time formula. Let's illustrate this.

3.4.14 (▣). 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?

3.4.15 (▣). Suppose we could buy another machine that never fails. What is the average sojourn time?

QUEUEING NETWORKS

We refer to the relevant sections of Zijm's book for background. Here we just include the solutions and repair a few typos.

4.1 OPEN SINGLE-CLASS PRODUCT-FORM NETWORKS

Theory and Exercises

The remark above Zijm.Eq.2.11 is not entirely correct. Remove the sentence: 'These visit ratios satisfy ... up to a multiplicative constant'.

I don't like the derivation of Zijm.Eq.2.20. The appearance of the visit ratios λ_i/γ seems to come out of thin air. The argument should be like this. Consider the entire queueing network as one 'box' in which jobs enter at rate $\gamma = \sum_{i=1}^M \gamma_i$. Assuming that there is sufficient capacity at each station, i.e., $\lambda_i < c_i \mu_i$ at each station i , the output rate of the 'box' must also be γ . Thus, by applying Little's law to the 'box', we have that

$$E[L] = \gamma E[W].$$


It is also evident that the average total number of jobs must be equal to the sum of the average number of jobs at each station:


$$E[L] = \sum_{i=1}^M E[L_i].$$

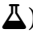
Applying Little's law to each station separately we get that $E[L_i] = \lambda_i E[W_i]$. Filling this into the above,

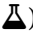
$$E[W] = \frac{E[L]}{\gamma} = \sum_{i=1}^M \frac{E[L_i]}{\gamma} = \sum_{i=1}^M \frac{\lambda_i E[W_i]}{\gamma},$$


where we recognize the visit ratios.


4.1.1 ( Linear algebra refresher). Can you find an example to show for two matrices A and B that $AB \neq BA$, hence $xA \neq Ax$.


4.1.2 ( Linear algebra refresher 2). Suppose the matrix A has an eigenvalue 0. What is the geometric meaning of this fact?

4.1.3 (). Zijm.Ex.2.2.1

4.1.4 (). Zijm.Ex.2.2.2

4.1.5 (). Zijm.Ex.2.2.3

4.1.6 (). Zijm.Ex.2.2.4

4.1.7 (). Zijm.Ex.2.2.5. The problem is not entirely correctly formulated. It should be, if for at least one i , $\sum_{j=1}^M P_{ij} < 1 \dots$

4.1.8 (▣). Zijm.Ex.2.2.6

4.1.9 (▣). Show that Zijm.Eq.2.13 and 2.14 can be written as

$$f_i(n_i) = \frac{1}{G(i)} \frac{1}{\prod_{k=1}^{n_i} \min\{k, c_i\}} \left(\frac{\lambda_i}{\mu_i} \right)^{n_i}.$$

4.1.10 (▣). 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. Compute λ . What happens if $\alpha \rightarrow 1$ or $\beta_1 \rightarrow 0$?

4.1.11 (▣). Zijm.Ex.2.2.8

4.2 TANDEM QUEUES

Theory and Exercises

Consider two $M/M/1$ stations in tandem. Suppose we can remove the variability in the service processing times at one, but not both, of the servers. Which one is the better one to spend it on, in terms of reducing waiting times? After we obtained some insights into this question, we will provide a model to approximate the waiting time in a tandem of $G/G/1$ queues.

4.2.1. Assuming that jobs arrive at the first station at rate λ , and are served at rate μ_i at station i , show that the average queueing time for the tandem of two $M/M/1$ queues is given by

$$E[W_Q] = \frac{\rho_1}{1-\rho_1} \frac{1}{\mu_1} + \frac{\rho_2}{1-\rho_2} \frac{1}{\mu_2}, \quad (4.2.1)$$

where $\rho_i = \lambda/\mu_i$ and $E[S_i] = 1/\mu_i$, for $i = 1, 2$.

4.2.2. Suppose we can remove all variability of service process at the second station. Show that in this case the total time in queue is equal to

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

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

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

is a reasonable approximation of the queueing time. Compare this to the queueing time of the reference situation.

4.2.4. What do you conclude from the above exercises?

For a tandem network of $G/G/1$ queues, observe that the SCV of the departure process $C_{d,i}^2$ of the i th station is the SCV of the arrival process $C_{a,i+1}^2$ at station $i+1$. Thus, if we have $C_{d,i}^2$ we can compute the average waiting time at station $i+1$ by means of the $G/G/1$ waiting time approximation.

To obtain an estimate for $C_{d,i}^2$ we reason as follows. Suppose that the load ρ_i at station i is very high. Then the server will seldom be idle, so that the departure process must be reasonably well approximated by the service process. If, however, the load is small, the server

will be idle most of the time, and inter-departure times must be approximately distributed as the inter-arrival times. Based on this, we interpolate between these two extremes to get the approximation

$$C_{d,i}^2 \approx (1 - \rho_i^2)C_{a,i}^2 + \rho_i^2 C_{s,i}^2. \quad (4.2.2)$$

4.2.5. What is C_d^2 for the $D/D/1$ queue according to (4.2.2)?

4.2.6. What is C_d^2 for the $M/M/1$ queue according to (4.2.2)?

4.2.7. Use (4.2.2) to show for the $G/D/1$ that $C_d^2 < C_a^2$.

4.2.8. 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? What is the average number of jobs in the network?

For a $G/G/c$ queue, we can use the following approximation

$$C_{d,i}^2 = 1 + (1 - \rho_i^2)(C_{a,i}^2 - 1) + \frac{\rho_i^2}{\sqrt{c_i}}(C_{s,i}^2 - 1). \quad (4.2.3)$$

4.2.9. Show that (4.2.3) reduces to (4.2.2) for the $G/G/1$ queue.

For the interested reader we refer to Zijm, Section 2.4.2, for a discussion of an extension for $G/G/c$ queues in tandem, and to networks. In particular, in networks we need to be concerned with output streams merging into a single input stream at one station, and the splitting of the output stream of a station to several other stations. The algorithm discussed in Zijm, Section 2.4.2, is mainly useful for numerical analysis. We will not discuss it here.

4.3 GORDON-NEWELL NETWORKS

Theory and Exercises

4.3.1. Provide an interpretation of a single-server queueing server with a finite calling population in terms of a closed network.

The formula with the visit ratios should be like this:

$$V_k = \sum_{j=0}^M V_j P_{jk},$$

i.e., the sum should start at index 0. This is to include the load/unload station.

Also, assume that the load/unload station has just one server.

You should realize that the algorithms discussed in this section are meant to be carried out by computers. Thus the results will be numerical, not in terms of formulas.

Mind the order of V and P in the computation of the visit ratios: do not mix up $VP = V$ with $PV = V$, as in general, $VP \neq PV$. We use $VP = V$.

4.3.2. Compute the visit ratios for a network with three stations such that all jobs from station 0 move to station 1, from station 1 all move to station 2, and from station 2 half of the jobs move to station 0 and the other half to station 1.

4.3.3. Zijm.Ex.3.1.1

4.3.4. Zijm.Ex.3.1.2

4.3.5. Zijm.Ex.3.1.3

4.3.6. Relate Zijm.Eq.3.3 to the form of the steady-state distribution of the number of jobs in an $M/M/c$ queue.

4.3.7. Zijm.Ex.3.1.4

4.3.8. Zijm.Ex.3.1.5

4.4 MVA ALGORITHM

Theory and Exercises


4.4.1. Consider two stations in tandem, stations 0 and 1. The service times are $E[S_0] = 2 = 1/\mu_0$ and $E[S_1] = 3 = 1/\mu_1$ hours. The routing matrix is

$$P = \begin{pmatrix} 0 & 1 \\ 1/2 & 1/2 \end{pmatrix}.$$

Apply the MVA algorithm to this case.

4.4.2. Zijm.Ex.3.1.13. Assume that all stations have just one server.

4.4.3. Zijm.Ex.3.1.14

4.4.4 (). Implement the MVA algorithm in your preferred computer language and make Figure 3.2.

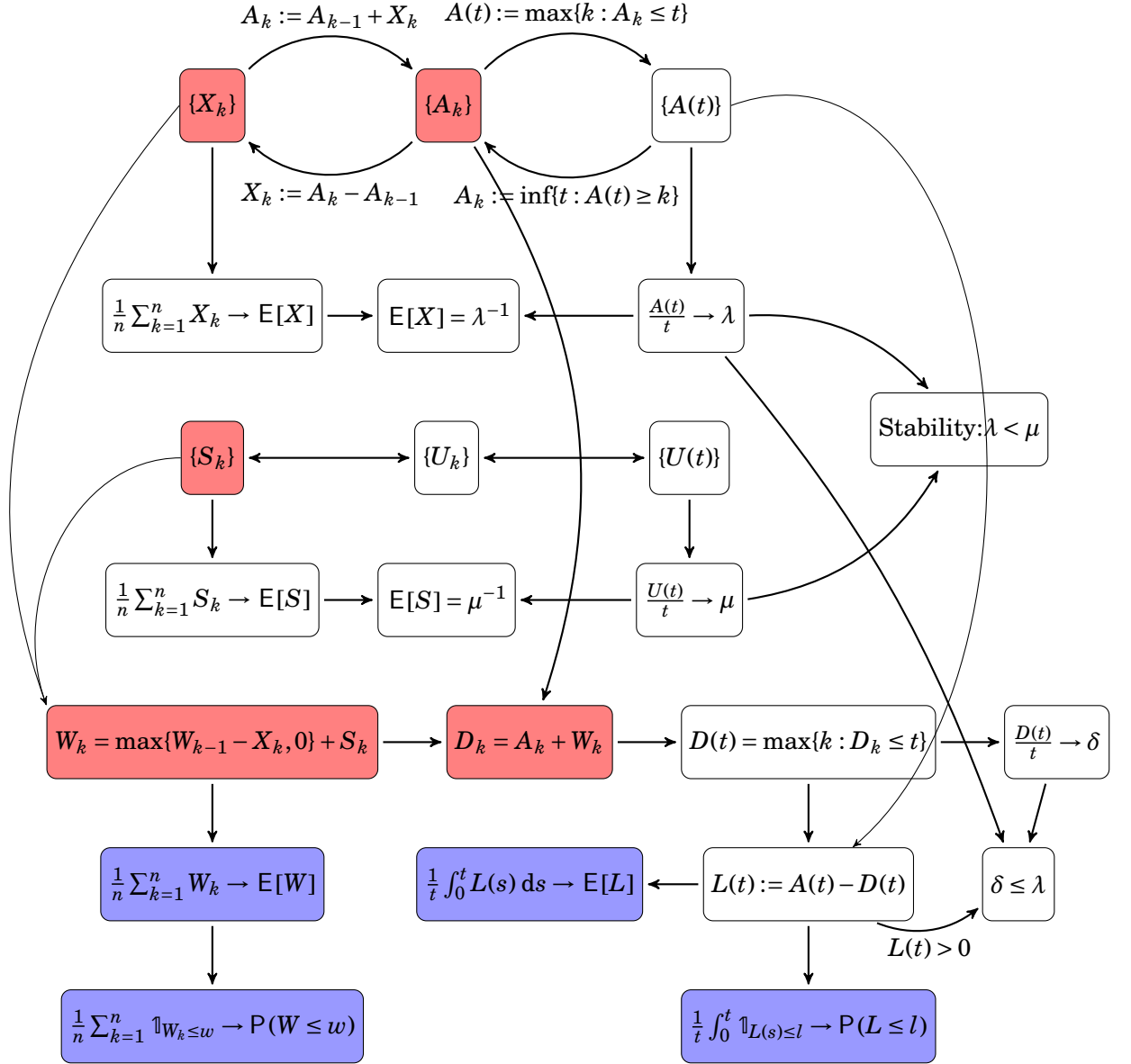


Figure 17: Here we sketch the relations between the construction of the $G/G/1$ queue from the primary data, i.e., the inter-arrival times $\{X_k; k \geq 0\}$ and the service times $\{S_k; k \geq 0\}$, and different performance measures.

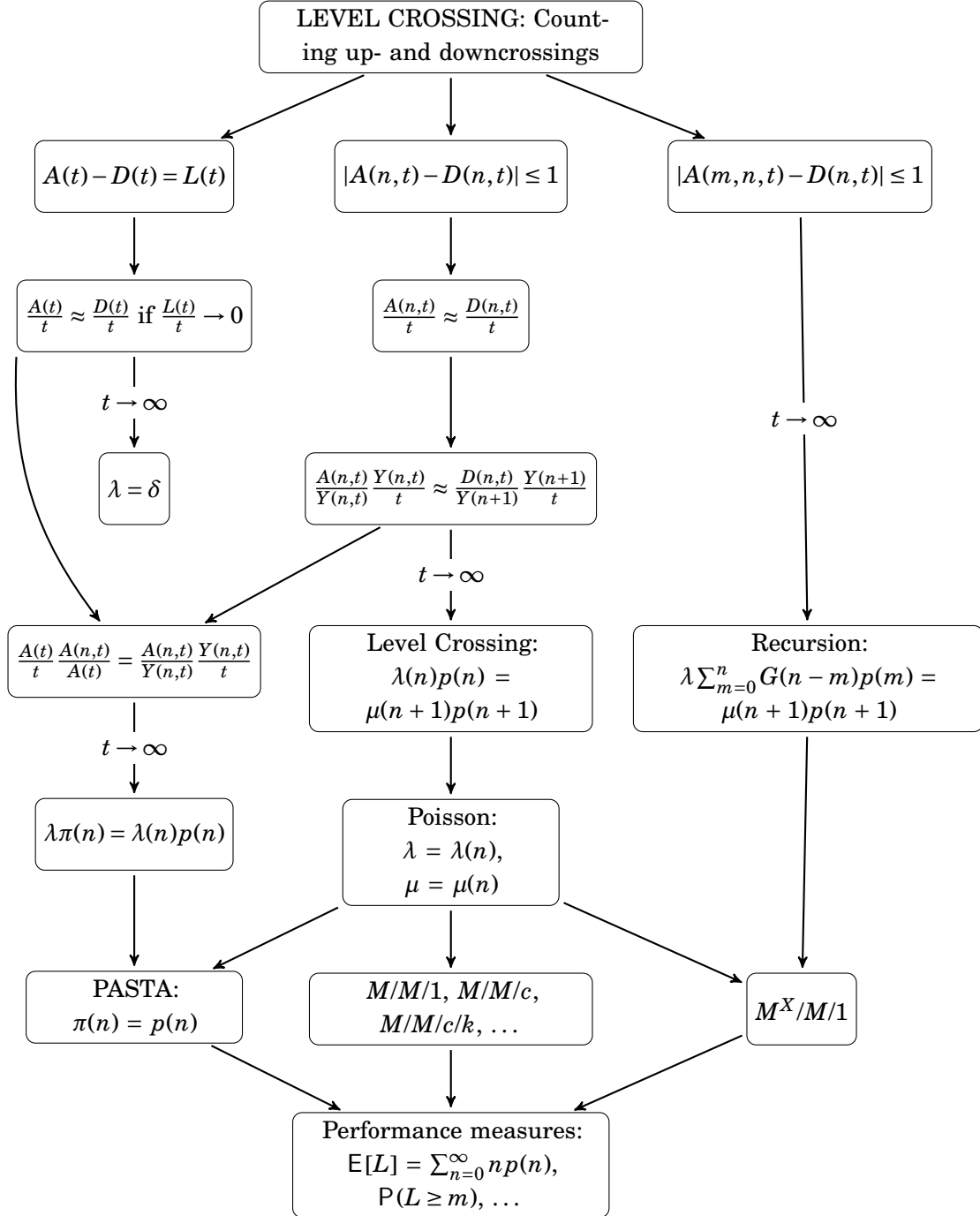


Figure 18: With level crossing arguments we can derive a number of useful relations. This figure presents an overview of these relations that we derive in this and the next sections.

BIBLIOGRAPHY

- F. Baccelli and W.A. Massey. A sample path analysis of the $M/M/1$ queue. *Journal of Applied Probability*, 26(2):418–422, 1988.
- 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.
- M. El-Taha and S. Stidham Jr. *Sample-Path Analysis of Queueing Systems*. Kluwer Academic Publishers, 1998.
- R.W. Hall. *Queueing Methods for Services and Manufacturing*. Prentice Hall, 1991.
- W.J. Hopp and M.L. Spearman. *Factory Physics*. Waveland Press, Inc., 3rd edition, 2008.
- H. Sakasegawa. An approximation formula $l_q = \alpha\beta^\rho/(1 - \rho)$. *Ananals of the Institute for Statistical Mathematics*, 29:67–75, 1977.
- 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
 $E[Q]$ = 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 = Generic service time of a job
 t = Time
 W_k = Time in the system 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

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

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

$$\text{Batching: } C_{sB}^2 = \frac{B V[S_0] + V[T]}{(B E[S_0] + E[T])^2}$$

$$\text{Nonpreemptive: } V[S] = V[S_0] + \frac{V[T]}{B} + \frac{B-1}{B^2} (E[T])^2$$

$$\text{Preemptive: } A = \frac{m_f}{m_r + m_f}, C_s^2 = C_0^2 + 2A(1-A) \frac{m_r}{E[S_0]}$$

$$C_{di}^2 = 1 + (1 - \rho_i^2)(C_{ai}^2 - 1) + \frac{\rho_i^2}{\sqrt{c_i}}(C_{si}^2 - 1)$$

$$f_i(n_i) = \begin{cases} G(i)^{-1} (c_i \rho_i)^{n_i} (n_i!)^{-1}, & \text{if } n_i < c_i, \\ G(i)^{-1} c_i^{c_i} \rho_i^{n_i} (c_i!)^{-1}, & \text{if } n_i \geq c_i \end{cases}$$

$$\text{with } G(i) = \sum_{n=0}^{c_i-1} \frac{(c_i \rho_i)^n}{n!} + \frac{(c_i \rho_i)^{c_i}}{c_i!} \frac{1}{1 - \rho_i}$$

$$E[L_i] = \frac{(c_i \rho_i)^{c_i}}{c_i! G(i)} \frac{\rho_i}{(1 - \rho_i)^2} + c_i \rho_i$$

$$f_i(n_i) = \frac{1}{\prod_{k=1}^{n_i} \min\{k, c_i\}} \left(\frac{V_i}{\mu_i} \right)^{n_i}, i = 0 \dots M$$

$$V_i = (VP)_i = \sum_{j=0}^M V_j P_{ji}$$

INDEX

- arrival process, 12
- arrival rate, 3, 21
- arrival times, 12
- average number of jobs, 25

- balance equations, 29
- balking, 34
- binomially distributed, 3
- Burke's law, 36

- conditional probability, 2

- departure rate, 22
- departure time of the system, 13

- effective processing time, 54, 55
- excess probability, 25
- expected waiting time, 24
- exponentially distributed, 11

- i.i.d., 11
- independent and identically distributed, 11
- indicator variable, 1
- inter-arrival times, 12

- Kendall's abbreviation, 17

- level crossing equations, 28
- limiting distribution, 20
- load, 23

- memoryless, 12
- Merging, 4
- moment-generating function, 2

- net service time, 54

- non-preemptive outages, 55
- normalization constant, 28
- number of jobs in the system, 14

- PASTA, 35
- Poisson arrivals see time averages, 35
- Poisson distributed, 4
- Poisson process, 4
- Pollaczek-Khinchine formula, 41
- probability mass function, 2
- processing rate, 22

- rate stable, 22
- remaining service time, 40
- renewal reward theorem, 23

- SCV, 4
- service rate, 22
- Small o notation, 1
- sojourn time, 14
- square coefficient of variation, 4
- stationary and independent increments, 3
- steady-state limit, 20
- survivor function, 2

- time-average number of jobs, 25

- up crossing rate, 26
- utilization, 23

- virtual waiting time process, 14

- waiting time in queue, 13
- Wald's equation, 56