

Queues with a Dynamic Schedule

John Gilbertson

A thesis presented for the degree of
Master of Science (Mathematics and Statistics)

Supervised by Professor Peter Taylor
Department of Mathematics and Statistics
The University of Melbourne

October 2016

Declaration

This thesis is the sole work of the author whose name appears on the title page and it contains no material which the author has previously submitted for assessment at the University of Melbourne or elsewhere. To the best of my knowledge and belief, the thesis contains no material previously published or written by another person, in the form of unacknowledged quotations or mathematical workings or in any other form, except where due reference is made. I declare that I have read, and in undertaking this research I have complied with, the University's Code of Conduct for Research. I also declare that I understand what is meant by plagiarism and that this is unacceptable.

Signed

John Gilbertson

Abstract

Abstract goes here.

Contents

1	Introduction	6
2	Literature Review	8
3	Static Schedule	11
3.1	Objective Function	11
3.2	Expected Customer Waiting Times	12
3.3	Computing the Objective Function	13
3.4	Example Models	14
4	Dynamic Schedule	16
4.1	Aim	16
4.2	Assumptions	16
4.3	List of Variables	17
4.4	Objective Function	17
4.5	Example Models	23
5	Value of Dynamic Schedule	26
5.1	Expected Cost Comparison	26
5.2	Percentage Cost Saving	27
6	Simulation Studies	29
7	Conclusion	30

A	Dynamic Cost Derivation	31
A.1	Transition Probability	31
A.2	Expected Transition Cost	33

Chapter 1

Introduction

Queues with scheduled arrivals occur frequently in society. These are queues where instead of customers arriving randomly, their arrival times are scheduled in advance. A common example of such queues is a doctor's surgery where patients are given appointment times. Moreover, these queues also occur in shipping when ship docking times are scheduled.

Throughout this thesis, we make several assumptions about the underlying system. First, we assume a single server queue with exponential service times. We assume that customers arrive punctually at their scheduled arrival times. In addition, all customers have the same mean service time μ .

The objective is to find a schedule of customer arrival times that minimises a linear combination of the expected total waiting time of all the customers and the expected time where the server is available (i.e., the total time between the start of service and conclusion of service of the last customer).

A customer's waiting time is the time from the customer's arrival until the server begins serving that customer. Instead of queue size, we refer to the number of customers in the system at a given point in time. This number includes both any customers currently being served and any customers currently waiting for service.

Bailey was the first to study such queues. His ideas have been extended in various ways including allowing for some non-punctuality of customers. However, few authors have considered the problem of adjusting a given schedule.

Most authors assume that a schedule is fixed at the start of service and cannot be altered during service. This static schedule appears to be a restrictive assumption. Through this thesis, we look at the potential advantage of being able to alter a schedule during service.

We consider a special case whereby customer arrivals are scheduled iteratively (i.e., one by one). A customer's arrival time is only decided on arrival of the customer to be served immediately before them. The expected cost of this dynamic schedule must be at least as good as the static schedule

Due to the relatively small number of customers in these queues, they do not reach steady state.

Chapter 2

Literature Review

Health care providers are under a great deal of pressure to improve service quality and efficiency (Goldsmith, 1989). There is a large body of literature studying the potential of appointment systems to reduce patient waiting times, and waiting room congestion. Fomundam and Herrmann (2007), and Cayirli and Veral (2003) provide comprehensive surveys of research on appointment scheduling. There is a fundamental trade-off in appointment policies. If patients are scheduled to arrive close together, they experience long waiting times. However, if appointment times are spread further apart, the doctor's idle time increases.

Most of the papers on scheduled arrivals in health care can be classed into two categories. Those that design algorithms to find good schedules, and those that evaluate schedules using simulation. While simulation studies can easily model complicated patient flows, queuing models often provide more generic results than simulation (Green, 2006).

The foundation paper on modeling queues with scheduled arrivals is Bailey (1952). Bailey proposes that customers' waiting times can be reduced without a significant increase in doctor's idle time. The Bailey rule, which is commonly referenced in literature, is that patients should be scheduled to arrive at fixed intervals with two patients scheduled to arrive at the start of service. Bailey found that a great deal of time wasted by patients could be reduced without a significant increase in the doctor's idle time. Under the Bailey rule, patients with late appointments will wait longer than those with early appointments. This lack

of uniformity might be undesirable due to issues of fairness.

Pegden and Rosenshine (1990) extend on Bailey's paper. They present an algorithm to iteratively determine the optimal arrival times for n patients that need to be scheduled. The optimal arrival times are those that minimise a weighted sum of the expected patients' waiting time and the expected doctor's idle time. Pegden and Rosenshine prove that their objective function is convex for $n \leq 4$, thus their algorithm finds the optimal schedule. While they conjecture that the objective function is convex for $n \geq 5$, it hasn't been proven.

Stein and Côté (1994) apply Pegden and Rosenshine's model to obtain numerical results for situations with more than three patients. The optimal times between successive patients become near constant as n grows. This is the often observed dome-shape. Optimal appointment intervals exhibit a common pattern where they initially increase towards the middle of a session, and then decrease. Stein and Côté simplify the model by requiring the intervals between arriving patients to be held constant. This realistic restriction (used commonly in the literature) makes the model more easily applicable in practice without significant altering the results.

Stein and Côté apply queuing theory results to solve the model for the optimal arrival interval assuming the queue reaches its steady state distribution. This assumption greatly reduces the computation required. However, in practice, it is common to find services that never reach steady state. Babes and Sarma (1991) attempted to apply steady state queuing theory, but found their results tended to be very different from those observed in real operation.

These key papers by Bailey, Pegden and Rosenshine, and Stein and Côté provide the basis for a more realistic exploration of health care systems. DeLaurentis et al. (2006) point out that patient no-shows can lead to a waste of resources. Mendel (2006) incorporates the probability of a patient not showing up into the model presented by Pegden and Rosenshine. Unsurprisingly, no-shows lead to lower expected waiting times for patients who do show up.

The presence of walk-ins (regular and emergency) can disrupt a schedule. Gupta, Zoreda, and Kramer (1971) propose a system where non-routine requests are superimposed on top of routine scheduled requests. Fiems, Koole, and Nain (2007) investigate the effect of emergency requests on the waiting times of sched-

uled patients. Fiem models a system with deterministic service times and discrete time. Despite this research, Cayirli and Veral suggest that walk-ins are neglected in most studies. Further research could investigate their effect on optimal arrival times.

Mondschein and Weintraub (2003) observe that the majority of the literature assumes that demand is exogenous and independent of patients' waiting times. These papers assume the total number of patients n is fixed and independent of waiting times. The vast majority of servers are now private (including medical servers), so face competitive environments. Mondschein and Weintraub thus present a model where demand depends on the patients' expected waiting time.

Simulation is a useful tool to analyse the effectiveness of appointment policies. Kao and Tung (1981) use simulation to compliment their results obtained from queuing theory. Ho and Lau (1992) study the performance of eight different appointment rules under different scenarios. They find that no rule will perform well under all circumstances.

Case studies can test the real world applications of an appointment system. While they lack generalisation, they are necessary to compliment the theoretical research. Rockart and Hofmann (1969) show individual block systems lead to more punctual doctors and patients, and less no-shows. Walter (1973) indicates that the simple grouping of inpatients and outpatients results in a substantial improvement in doctor's idle time.

Unfortunately, Cayirli and Veral (2003) lament that despite much published work, the impact of appointment systems on outpatient clinics has been limited. Doctors are often unwilling to change their old habits. O'Keefe (1985) had their proposed appointment system of classifying patients rejected by staff. Huarng and Hou Lee (1996) were unable to implement their system due to staff resistance. Bennett and Worthington (1998) found their recommendations weren't implemented successfully. Future research must attempt to develop models that will be accepted and implemented in real health care services.

Chapter 3

Static Schedule

This chapter largely follows the results of Pegden and Rosenshine (1990). The aim is to derive a method for choosing an optimal static schedule. A static schedule is a sequence of n customer arrival times t_1, \dots, t_n chosen at the start of service and fixed for the duration of service.

For simplicity, these results assume the customer service times are independent and identically distributed (iid) and follow an exponential distribution with mean μ . There is a single server. All customers are punctual and arrive at their scheduled arrival time.

3.1 Objective Function

The optimal schedule minimises the expected cost, which is a linear combination of the expected total customers' waiting time and total expected server availability time.

Denote the expected waiting time of customer i as w_i . The expected total customers' waiting time is the sum of the individual customer's expected waiting times.

$$\mathbb{E}[\text{total customer's waiting time}] = \sum_{i=1}^n w_i \quad (3.1)$$

Instead of solving for the customer arrival times, it is easier to solve for the arrival time of the first customer t_1 and the customer interarrival times $\mathbf{x} =$

(x_1, \dots, x_{n-1}) where $x_i = t_{i+1} - t_i$. The expected total server availability time is the expected time from the start of service until the end of service. This is the sum of the last customer's scheduled arrival time, the last customer's expected waiting time and the last customer's expected service time.

$$\mathbb{E}[\text{total server availability time}] = \left(t_1 + \sum_{i=1}^{n-1} x_i \right) + w_n + \mu \quad (3.2)$$

Denote c_W and c_I as the per unit time cost of the expected total customer's waiting time and the per unit time cost of the expected total server availability time respectively. The objective function to be minimised is thus,

$$\phi(t_1, \mathbf{x}) = c_W \sum_{i=1}^n w_i + c_S \left[t_1 + \sum_{i=1}^{n-1} x_i + w_n + \mu \right] \quad (3.3)$$

The first customer should obviously be scheduled for the start of service so $t_1 = 0$. Moreover, can scale the objective function by dividing by $(c_W + c_S)$ and defining $\gamma = \frac{c_S}{c_W + c_S}$.

$$\phi(\mathbf{x}) = (1 - \gamma) \sum_{i=1}^n w_i + \gamma \left[\sum_{i=1}^{n-1} x_i + w_n + \mu \right] \quad (3.4)$$

The optimal static schedule is thus the interarrival times that minimise Equation 3.4

3.2 Expected Customer Waiting Times

We want to express w_i (i.e., the expected waiting time of customer i) as a function of the interarrival times \mathbf{x} . If there are j customers in the system just prior to the arrival of customer i , then $w_i = j\mu$ by the memoryless property of the exponential distribution. The number of customers in the system refers to both customers being served and customers waiting for service.

Denote the number of customers in the system just prior to the arrival of

customer i as N_i . Thus, the expected waiting time of customer i is given by

$$w_i = \sum_{j=0}^{i-1} (j\mu) \mathbb{P}(N_i = j) \quad (3.5)$$

The probability of a given number of customers in the system can be expressed recursively. This probability depends on the number of departures from the system between the arrival of customer $(i - 1)$ and customer i . The number of departures is the minimum of $(N_{i-1} + 1)$ and a Poisson random variable with mean $\frac{x_{i-1}}{\mu}$.

The full recursive expression for the probability of j customers in the system immediately before the arrival of customer i is

$$\mathbb{P}(N_i = j) = \begin{cases} 1 & \text{for } i = 1, j = 0 \\ \sum_{k=1}^{i-1} \mathbb{P}(N_{i-1} = k - 1) \left[1 - \sum_{l=0}^{k-1} \frac{x_{i-1}^l}{\mu^l l!} e^{-\frac{x_{i-1}}{\mu}} \right] & \text{for } i \geq 2, j = 0 \\ \sum_{k=0}^{i-j-1} \mathbb{P}(N_{i-1} = j + k - 1) \left[\frac{x_{i-1}^k}{\mu^k k!} e^{-\frac{x_{i-1}}{\mu}} \right] & \text{for } i \geq 2, j \geq 1 \end{cases} \quad (3.6)$$

3.3 Computing the Objective Function

Pegden and Rosenshine (1990) suggest a similar algorithm to Algorithm 1 for computing the value of the objective function given by Equation 3.4 for a given vector \mathbf{x} and parameters γ and μ .

Algorithm 1 Return $\phi(\mathbf{x})$ for a given vector \mathbf{x} , γ and μ

```

function OBJECTIVEFUNCTION( $\mathbf{x}, \gamma, \mu$ )
  for  $i = 1, 2, \dots, n$  do
    for  $j = 0, 1, \dots, (i - 1)$  do
      compute  $\mathbb{P}(N_i = j)$  by Equation 3.6
  for  $i = 1, 2, \dots, n$  do
    compute  $w_i$  by Equation 3.5
  return  $\phi(\mathbf{x})$  computed by Equation 3.4

```

3.4 Example Models

3.4.1 Model for Two Customers

We consider the simplest case of this model where there are two customers to be scheduled (i.e., $n = 2$). As the first customer is scheduled to arrive at the start of service, the only unknown variable is the optimal interarrival time between the first and second customer (i.e., x_1).

By Equation 3.5, the expected waiting times of the two customers are

$$w_1 = 0 \quad (3.7)$$

$$w_2 = \mu e^{\frac{-x_1}{\mu}} \quad (3.8)$$

By Equation 3.4, the objective function to be minimised is given by

$$\phi(x_1) = \mu \left[\gamma + \exp\left(\frac{-x_1}{\mu}\right) \right] + \gamma x_1 \quad (3.9)$$

This objective function is convex as

$$\forall x_1 \quad \phi''(x_1) = \frac{1}{\mu} \exp\left(\frac{-x_1}{\mu}\right) > 0 \quad (3.10)$$

Due to the convexity of the objective function, the optimal policy that minimises $\phi(x_1)$ can be found by solving:

$$\phi'(x_1) = 0 \implies -\exp\left(\frac{-x_1}{\mu}\right) + \gamma = 0 \quad (3.11)$$

Therefore, the optimal policy is:

$$x_1^* = \arg \min_{x_1} \phi(x_1) = -\mu \ln \gamma \quad (3.12)$$

As the server availability cost increases relative to the customer waiting cost (i.e., γ increases), the second customer is scheduled to arrive earlier (i.e., x_1^* decreases).

Unfortunately, as Pegden and Rosenshine (1990) found, no general algebraic

solution exists for more than two customers (i.e., $n \geq 3$). All other cases need to be solve numerically. Pegden and Rosenshine (1990) proved that the objective function $\phi(\mathbf{x})$ is convex for $n = 1, 2, 3, 4$. Moreover, they conjecture that it appears to be convex for general n , but are unable to prove it.

3.4.2 Model for 15 Customers

The more interesting case concerns scheduling a larger number of customers. We consider here the case of $n = 15$. Without loss of generality, can let $\mu = 1$. For $\mu \neq 1$, the solutions found are the optimal values of $\mu\mathbf{x} = (\mu x_1, \dots, \mu x_{n-1})$.

Minimising the objective function is a constrained optimization problem. A numeric solution can be found using `scipy.optimize.minimize` in Python.

Chapter 4

Dynamic Schedule

4.1 Aim

The aim is to choose a schedule of customer arrivals times that minimises the expected cost of the system. The expected cost is a linear combination of the total customers' waiting time and the server's idle time. Instead of choosing a fixed schedule at the start of service as is common in literature, the schedule will be chosen progressively. Immediately after a customer arrives and begins waiting for service, the scheduler chooses the arrival time of the next customer.

4.2 Assumptions

To simplify this problem, need to make several assumptions:

- Service times are independent and identically distributed (iid)
- Each service time has an exponential distribution with mean service time μ
- There is a single server
- The queue operates on a first in, first out (FIFO) basis
- Customers can be scheduled to arrive at any future (or present) time
- Customers are punctual and arrive at their scheduled time

4.3 List of Variables

μ	: mean service time of each customer
c_W	: cost of total customers' waiting time per unit time
c_S	: cost of total server's availability time per unit time
k	: current number of customers in the system
j	: number of customers in the system immediately after the next customer's arrival
n	: number of customers remaining to be scheduled
a	: time next customer is scheduled to arrive (relative to current time)
$C_n^*(k)$: the expected cost of having k customers in the system and n customers remaining to be scheduled
$C_n(a, k)$: the expected cost of having k customers in the system, n customers remaining to be scheduled and the next customer scheduled to arrive after a time units
$p_a(k, j)$: the probability of transitioning from k customers in the system to j customers in the system after a time units if the next customer is scheduled to arrive after a time units
$R_a(k, j)$: the expected cost of transitioning from k customers in the system to j customers in the system after a time units if the next customer is scheduled to arrive after a time units

4.4 Objective Function

The state (n, k) refers to n customers remaining to be scheduled and k customers currently in the system (i.e., either waiting or being service). The time the next customer is scheduled to arrive relative to the current time is a . The number of customers in the system immediately after that next customer's arrival is $j \in \{1, \dots, k+1\}$.

The expected cost of the current state is a linear combination of the total customers' expected waiting time and the expected total service availability time from now until the end of service. It is a function of the expected cost involved in transitioning to the next state, the expected cost of the next state and the probability of transitioning to the next state over all possible next states.

The expected cost of the state (n, k) where $n \geq 1$ is given by the following form of Bellman's equation:

$$C_n^*(k) = \min_{a \geq 0} C_n(a, k) = \min_{a \geq 0} \left[\sum_{j=1}^{k+1} p_a(k, j) (R_a(k, j) + C_{n-1}^*(j)) \right] \quad (4.1)$$

Equation 4.1 is a recursive equation involving C^* . The optimal solution is found by solving for each customer's arrival time a iteratively. The optimal policy a^* is the customer's arrival time that attains the minimum cost whereby

$$C_n^*(k) = C_n(a^*, k) = \min_{a \geq 0} C_n(a, k) \quad (4.2)$$

It is reasonably intuitive that the minimum cost cannot occur at $a = \infty$. As $a \rightarrow \infty$, the total server availability time converges to ∞ . As the per unit time cost of the server's availability time c_S is strictly positive, the overall expected cost must also converge to ∞ as $a \rightarrow \infty$. Thus, $\lim_{a \rightarrow \infty} C_n(a, k) = \infty$.

Moreover, as will be explained later in this chapter, for given values of n and k , $C_n(a, k)$ is continuous where $a \in [0, \infty)$.

Consider the set of possible policies \mathcal{A} given by

$$\mathcal{A} = \{0\} \cup \left\{ a > 0 : \frac{\partial}{\partial a} C_n(a, k) = 0 \right\} \quad (4.3)$$

Solving Equation 4.2 involves solving a nonlinear optimisation problem over a left-closed interval. As $C_n(a, k)$ is continuous of $a \in [0, \infty)$ and the minimum cannot occur at $a = \infty$, this solution is equivalent to the solution found by checking the left end point (where $a = 0$) and all points where $\frac{\partial}{\partial a} C_n(a, k) = 0$. Thus, the the optimal policy can be found by solving

$$C_n^*(k) = \min_{a \in \mathcal{A}} C_n(a, k) \quad (4.4)$$

As will be explained later, it is not possible to find a 'nice' closed form for $\frac{\partial}{\partial a} C_n(a, k)$ for general n and k . However, for given values of n and k , it is reasonably efficient to solve $\frac{\partial}{\partial a} C_n(a, k) = 0$. Thus, the expected cost can be found by computing $\frac{\partial}{\partial a} C_n(a, k)$ for each state (n, k) . Of course, this method becomes more

computationally inefficient, the larger the number of states.

4.4.1 Base Case

Solving Equation 4.1 requires a solution for the base case where $n = 0$. If $n = 0$, there are no customers remaining to be scheduled, thus the expected total server availability time is the expected time until the customer currently in the last position in the queue finishes service. This time is the summation of the expected service time of all the customers currently in the system.

Let w_i be the expected waiting time of the customer that is currently in position i in the queue, c_W be the cost of the customers' waiting time per unit time and c_S be the cost of the server's availability time per unit time. The waiting time of the customer currently in position i is the summation of the service times of the customers in positions $\{1, \dots, i-1\}$ (i.e., it doesn't include the service time of the customer in position i). The cost of the base case is thus given by

$$\begin{aligned} C_0^*(k) &= c_W \sum_{i=2}^k w_i + c_S \sum_{i=1}^k \mu \\ &= c_W \sum_{i=2}^k \mu(i-1) + c_S k \mu \\ &= \frac{c_W \mu k(k-1)}{2} + c_S k \mu \end{aligned}$$

Scale $C_0^*(k)$ by dividing by $(c_S + c_W)$ and substituting $\gamma = \frac{c_S}{c_S + c_W}$:

$$C_0^*(k) = (1 - \gamma) \cdot \frac{\mu k(k-1)}{2} + \gamma k \mu \quad (4.5)$$

4.4.2 Transition Probability

Let S_i be the service time of the customer that is currently in position i in the queue. The service times S_1, \dots, S_n are iid (independent and identically distributed) exponential random variables with mean μ .

For $r \geq 1$, the waiting time of the customer in position $(r + 1)$ in the queue is:

$$X = \sum_{i=1}^r S_i \sim \text{Erlang}(r, \mu) \quad (4.6)$$

which has the pdf:

$$f(x; r) = \frac{1}{\mu \cdot (r - 1)!} \left(\frac{x}{\mu} \right)^{r-1} \exp \left(-\frac{x}{\mu} \right) \quad (4.7)$$

Let W_t be a Poisson Point Process with $W_t \sim \text{Poisson} \left(\frac{t}{\mu} \right)$. For $r \geq 1$, the probability that the customer currently in position $(r + 1)$ in the queue waits longer than a time units before service is equal to the probability that W_a is smaller than r , such that

$$\mathbb{P}(X > a) = \mathbb{P}(W_a < r) \quad (4.8)$$

The transition probability $p_a(k, j)$ is the probability that the queue length changes from k customers initially to j customers on the arrival of the next customer after a time units. In other words, it is the probability that there are $k - (j - 1)$ departures from the queue over a time interval of length a . Computing this probability requires the cdf of the Erlang distribution, which is calculated (for $a > 0$) as follows:

$$\begin{aligned} F(a; r) &= \mathbb{P}(X \leq a) \\ &= 1 - \mathbb{P}(X > a) \\ &= 1 - \mathbb{P}(W_a < r) \\ &= 1 - \sum_{i=0}^{r-1} \mathbb{P}(W_a = i) \\ &= 1 - \sum_{i=0}^{r-1} \frac{1}{i!} \left(\frac{a}{\mu} \right)^i \exp \left(-\frac{a}{\mu} \right) \end{aligned}$$

Moreover, $F(0; r) = \mathbb{P}(X = 0) = 0$ by definition. Therefore, the cdf of the

Erlang distribution is defined as

$$F(a; r) = \begin{cases} 1 - \sum_{i=0}^{r-1} \frac{1}{i!} \left(\frac{a}{\mu}\right)^i \exp\left(\frac{-a}{\mu}\right) & \text{where } a > 0 \\ 0 & \text{where } a = 0 \end{cases} \quad (4.9)$$

For $r \geq 1$, $F(a; r)$ is continuous for $a \geq 0$ as $\lim_{a \rightarrow 0^+} F(a; r) = 0 = F(0; r)$. Can now compute the transition probability on a case by case basis. The full derivation is included in the appendix and only the final equation is presented here.

$$p_a(k, j) = \begin{cases} \mathbb{1}(j = 1) & \text{where } k = 0 \\ F(a; k) & \text{where } k \geq 1, j = 1 \\ F(a; k - j + 1) - F(a; k - j + 2) & \text{where } k \geq 1, 2 \leq j \leq k \\ 1 - F(a; 1) & \text{where } k \geq 1, j = (k + 1) \\ 0 & \text{otherwise} \end{cases} \quad (4.10)$$

4.4.3 Expected Transition Cost

The cost involved in transitioning from state (n, k) to state $(n - 1, j)$ in a time units is a linear combination of the expected total waiting time of the customers during the transition and the expected total server availability time during the transition as in Chapter 3. During the transition, the server is available for the entire time interval, so the total server availability time during the transition is always a .

The expected total waiting time of the customers depends on the conditional

expectation of the Erlang distribution, which is derived here (for $a > 0$).

$$\begin{aligned}
 G(a; k) &= \mathbb{E}[X|X \leq a] \\
 &= \int x \cdot \mathbb{P}(X \in dx|X \leq a) \\
 &= \int x \cdot \frac{\mathbb{P}(X \in dx, X \leq a)}{\mathbb{P}(X \leq a)} \\
 &= \frac{1}{F(a; r)} \int_0^a x \cdot f(x; r) dx \\
 &= \frac{1}{F(a; r)} \int_0^a x \cdot \frac{1}{\mu \cdot (r-1)!} \left(\frac{x}{\mu}\right)^{r-1} \exp\left(\frac{-x}{\mu}\right) \\
 &= \frac{\mu r}{F(a; r)} \int_0^a \frac{1}{\mu \cdot r!} \left(\frac{x}{\mu}\right)^r \exp\left(\frac{-x}{\mu}\right) \\
 &= \mu r \cdot \frac{F(a; r+1)}{F(a; r)}
 \end{aligned}$$

In addition, $G(0; r) = \mathbb{E}[X|X \leq 0] = 0$ by definition. Therefore, the conditional expectation of the Erlang distribution is defined as

$$G(a; r) = \begin{cases} \mu r \cdot \frac{F(a; r+1)}{F(a; r)} & \text{where } a > 0 \\ 0 & \text{where } a = 0 \end{cases} \quad (4.11)$$

This expression for the conditional expectation makes intuitive sense. The mean of the Erlang distribution is $\mathbb{E}[X] = \mu r$. In addition, for $a > 0$ and $r \geq 1$, $\frac{F(a; r+1)}{F(a; r)} < 1$. Thus, for $a > 0$ and $r \geq 1$, $\mathbb{E}[X|X \leq a] < \mathbb{E}[X]$ as expected. Moreover,

$$\lim_{a \rightarrow \infty} \frac{F(a; r+1)}{F(a; r)} = \frac{\lim_{a \rightarrow \infty} F(a; r+1)}{\lim_{a \rightarrow \infty} F(a; r)} = \frac{1}{1} = 1 \quad (4.12)$$

Thus, $\lim_{a \rightarrow \infty} \mathbb{E}[X|X \leq a] = \mathbb{E}[X]$ as expected.

In a similar way to the transition probability, the expected transition cost is derived on a case by case basis. The cost is scaled by defining $\gamma = \frac{c_S}{c_S + c_W}$. The full

derivation is included in the appendix and the final equation is presented here.

$$\begin{aligned}
 R_a(k, j) &= \begin{cases} \gamma a & \text{where } k \in \{0, 1\} \\ (1 - \gamma) \cdot \frac{G(a; k)(k-1)}{2} + \gamma a & \text{where } k \geq 2, j = 1 \\ (1 - \gamma) \left[a(j-2) + \frac{G(a; k-(j-1))[k-(j-2)]}{2} \right] + \gamma a & \text{where } k \geq 2, 2 \leq j \leq k \\ (1 - \gamma) \cdot a(j-2) + \gamma a & \text{where } k \geq 2, j = (k+1) \end{cases} \\
 &\quad (4.13)
 \end{aligned}$$

4.5 Example Models

4.5.1 State (3, 2)

Assume that the mean service time $\mu = 1$. In addition, assume that the per unit time costs of the customers' waiting time (c_W) and server's total availability time (c_S) are both equal such that $\gamma = 0.5$. Under these assumptions, the expected cost of the state (3, 2) for various values of a is given by the following expression, which is continuous for $a \geq 0$.

$$C_3(a, 2) = \begin{cases} 4.94 & \text{where } a = 0 \\ 2.76 + \frac{a}{2} + \exp(-a) \left(\frac{a^2}{4} + 1.65a + 2.18 \right) - \frac{a^2}{2[\exp(a)-1]} & \text{where } a > 0 \end{cases}$$

Figure 4.1 plots this cost for $a \in [0, 10]$ to visualise the relationship between the cost and a . The only value of a where $\frac{\partial}{\partial a} C_3(a, 2) = 0$ is 1.90. Thus, the set of possible policies \mathcal{A} is $\{0, 1.90\}$ whereby there are two possible policies labelled as orange dots on Figure 4.1. Moreover, note that as a increases beyond 2, the cost increases approximately linearly.

The optimal policy a^* that minimises $C_3(a, 2)$ is 1.90 where the cost is 8.64. If (on arrival of a customer) there are two customers in the system and three customers remaining to be scheduled, then the next customer should be scheduled to arrive in 1.90 time units. This is slightly below the expected service time of the two customers in the system to account for the availability cost of the server.

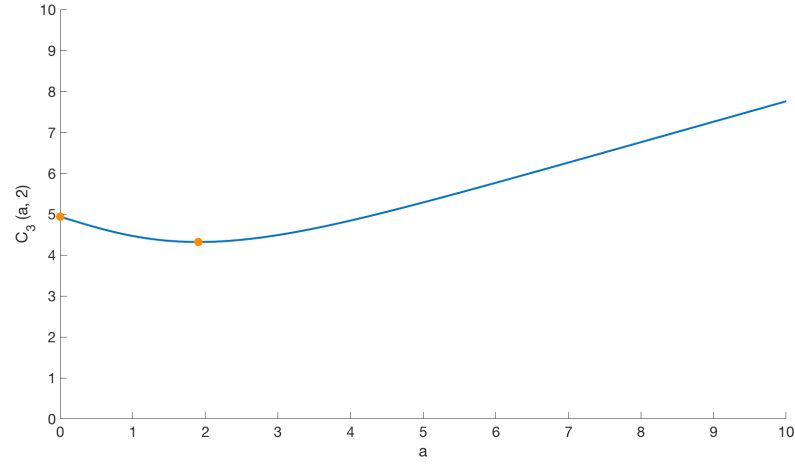


Figure 4.1: Cost of state with $n = 3$ and $k = 2$ for $a \in [0, 10]$ where $\mu = 1$ and $\gamma = 0.5$.

4.5.2 Model for Six Customers

Assume there are six customers that need to be scheduled for service who all have a mean service time $\mu = 1$. The initial state is $(6, 0)$, and the possible states during service are all states in the set

$$\left\{ (n, k) \in \{0, 1, \dots, 6\}^2 : n + k \leq 6 \right\} \quad (4.14)$$

		Current Number in System (k)						
		0	1	2	3	4	5	6
Customers to be Scheduled (n)	0	0	0.5	1.5	3	5	7.5	10.5
	1	0.5	1.35	2.49	4.02	5.98	8.37	
	2	1.35	2.26	3.41	4.94	6.89		
	3	2.26	3.18	4.32	5.86			
	4	3.18	4.09	5.24				
	5	4.09	5.01					
	6	5.01						

Table 4.1: Cost of possible states for 6 total customers

Table 4.1 displays the calculated expected costs for each possible state if there are six total customers (assuming $\gamma = 1$). The cost of the initial state is 5.01,

thus the expected cost of servicing six customers with a dynamic schedule is 5.01.

Note for $n \geq 1$, $C_n^*(0) = C_{n-1}^*(1)$. This is due to the fact that if there are no customers currently in the system (e.g., initially), then it is always optimal to schedule the next arrival immediately. The expected cost thus doesn't change as the next arrival occurs immediately.

The worst state on Table 4.1 is the state $(0, 6)$, which is all six customers in the system. This can occur if all six customers are scheduled to arrive immediately (to ensure minimal server availability time) or if the first customer has an extremely long service time.

		Current Number in System (k)					
		0	1	2	3	4	5
Customers to be Scheduled (n)	1	0	0.69	1.76	2.83	3.86	4.85
	2	0	0.83	1.90	2.95	3.96	
	3	0	0.83	1.90	2.95		
	4	0	0.83	1.90			
	5	0	0.83				
	6	0					

Table 4.2: Optimal policy for each possible states for 6 total customers

Table 4.2 displays the corresponding arrival times for the costs in Table 4.1. This table does not include any optimal times for $n = 0$, as there is no next customer to schedule in those states.

The first pattern to notice is (as discussed earlier) if there are no customers currently waiting (i.e., $k = 0$), then the optimal policy is to schedule the next arrival immediately. This makes intuitive sense as scheduling the next arrival immediately minimises the total server availability time without affecting the total customers' waiting times.

As k increases for fixed n , the optimal scheduled arrival time a^* appears to increase at a decreasing rate. The optimal a^* increases by approximately μ for each extra k .

In contrast, for $n \geq 2$, a^* appears to be constant at a value slightly less than μk (i.e., the expected time for the system to be empty). This is due to the server availability cost leading to a desire for the next customer to arrive just before the queue becomes empty.

Chapter 5

Value of Dynamic Schedule

5.1 Expected Cost Comparison

Assume that there are N customers to be scheduled. In the case of the static schedule, the expected cost of the optimal policy is given by $\phi(\mathbf{x}^*) = \phi(x_1^*, \dots, x_{N-1}^*)$. In the case of the dynamic schedule, the expected cost of the optimal policy is given by $C_N^*(0)$.

Clearly, as the optimal policy for the dynamic schedule can match the optimal policy for the static schedule, $C_N^*(0) \leq \phi(\mathbf{x}^*)$ regardless of the number of customers to be scheduled. Equality is attained for the cases where $N \in \{1, 2\}$ as the schedules are identical in those cases.

Figure 5.1 plots the expected costs of both the static and dynamic schedules against the number of customers to be scheduled (N) assuming $\gamma = 0.5$ and $\mu = 1$. As expected, the costs are identical for $N \in \{1, 2\}$. For all other N values, the cost of the static schedule is greater than the cost of the dynamic schedule. As N increases, the difference between the costs increases.

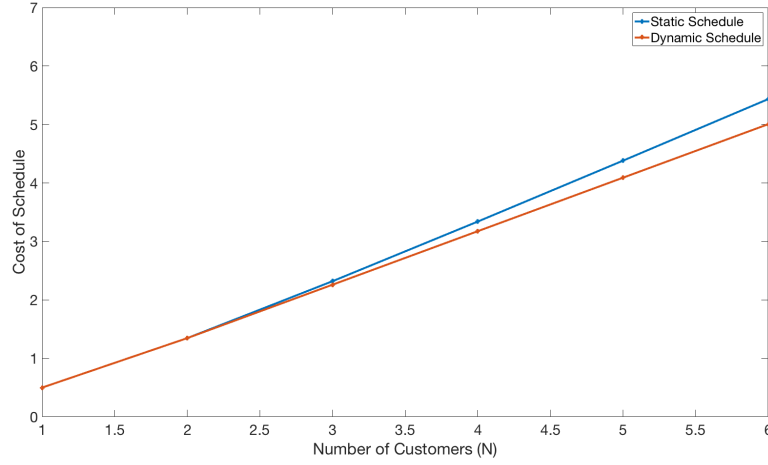


Figure 5.1: Plot of the expected cost of each schedule against the number of customers to be scheduled (N) for both the static and dynamic schedules where $N \in \{1, \dots, 6\}$, $\gamma = \frac{c_S}{c_S + c_W} = 0.5$ and $\mu = 1$.

5.2 Percentage Cost Saving

Define ΔC which measures the percentage difference between the cost of the static schedule and the dynamic schedule (i.e., the percentage cost saving).

$$\Delta C = 100 \times \frac{\phi(\mathbf{x}^*) - C_N^*(0)}{\phi(\mathbf{x}^*)} \quad (5.1)$$

Figure 5.2 plots ΔC against γ for various values of the number of customers to be scheduled (N). For $N \in \{1, 2\}$, $\Delta C = 0$ for all values of γ as the schedules are identical. As N increases with γ held constant, ΔC increases at a decreasing rate (i.e., the curves become closer together as N increases).

For the maximum considered value of γ (i.e., $\gamma = 0.95$), ΔC is at a minimum for all values of N . A large value of γ indicates that the server availability cost per unit time is significantly greater than the waiting cost per unit time (i.e., $c_S \gg c_W$). As the server availability cost is heavily prioritised, there is little difference between the static and dynamic schedules, thus ΔC is small.

For each value of $N \geq 3$, the peak value of ΔC occurs at a middle value of γ . As N increases, the peak occurs at a larger value of γ . For $N = 3, 4, 5$ and 6 , the

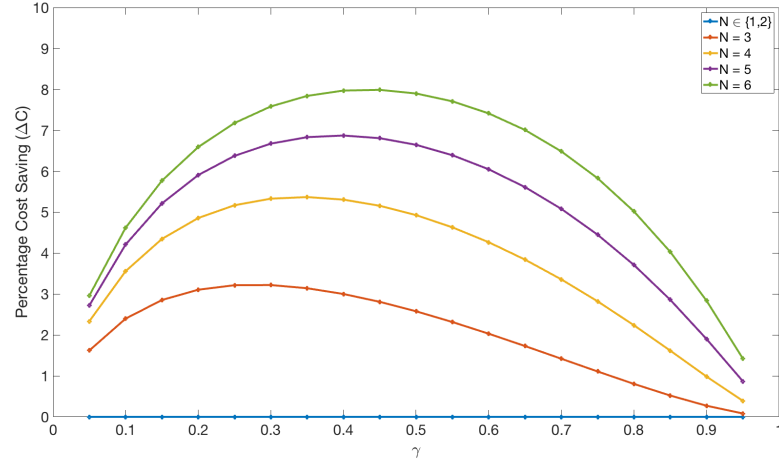


Figure 5.2: Plot of the percentage cost saving (ΔC) by using the dynamic schedule as opposed to the static schedule against $\gamma = \frac{c_S}{c_S + c_W}$ where $N = \{1, \dots, 6\}$ and $\mu = 1$.

peak occurs at $\gamma \approx 0.3, 0.35, 0.4$ and 0.45 respectively.

Chapter 6

Simulation Studies

Simulation studies goes here.

Chapter 7

Conclusion

Conclusion goes here.

Appendix A

Dynamic Cost Derivation

A.1 Transition Probability

A.1.1 Case 1 $k = 0$

$$p_a(k, j) = \mathbb{1}(j = 1)$$

A.1.2 Case 2 $k \geq 1, j = 1$

$$p_a(k, j) = \mathbb{P} \left(\sum_{i=1}^k S_i \leq a \right) = F(a; k)$$

A.1.3 Case 3 $k \geq 1, 2 \leq j \leq k$

$$\begin{aligned}
 & p_a(k, j) \\
 &= \mathbb{P} \left(\sum_{i=1}^{k-(j-1)} S_i \leq a, \sum_{i=1}^{k-(j-1)+1} S_i > a \right) \\
 &= \mathbb{P} \left(\sum_{i=1}^{k-(j-1)} S_i \leq a, \sum_{i=1}^{k-(j-1)} S_i + S_{k-(j-1)+1} > a \right) \\
 &= \mathbb{P} \left(\sum_{i=1}^{k-(j-1)} S_i \leq a, S_{k-(j-1)+1} > a - \sum_{i=1}^{k-(j-1)} S_i \right) \\
 &= \int \mathbb{P} \left(\sum_{i=1}^{k-(j-1)} S_i \leq a, S_{k-(j-1)+1} > a - \sum_{i=1}^{k-(j-1)} S_i \mid \sum_{i=1}^{k-(j-1)} S_i = z \right) \mathbb{P} \left(\sum_{i=1}^{k-(j-1)} S_i \in dz \right) \\
 &= \int_0^\infty \mathbb{P}(z \leq a, S_{k-(j-1)+1} > a - z) f(z; k - (j - 1)) dz \\
 &= \int_0^a \mathbb{P}(S_{k-(j-1)+1} > a - z) f(z; k - (j - 1)) dz \\
 &= \int_0^a f(z; k - (j - 1)) (1 - \mathbb{P}(S_{k-(j-1)+1} \leq a - z)) dz \\
 &= \int_0^a f(z; k - (j - 1)) (1 - F(a - z; 1)) dz \\
 &= \int_0^a \frac{1}{\mu \cdot (k - (j - 1) - 1)!} \left(\frac{z}{\mu} \right)^{k-(j-1)-1} \exp \left(\frac{-z}{\mu} \right) \cdot \exp \left(\frac{-(a - z)}{\mu} \right) dz \\
 &= \frac{1}{(k - (j - 1) - 1)!} \left(\frac{1}{\mu} \right)^{k-(j-1)} \exp \left(\frac{-a}{\mu} \right) \int_0^a z^{k-(j-1)-1} dz \\
 &= \frac{1}{(k - (j - 1) - 1)!} \left(\frac{1}{\mu} \right)^{k-(j-1)} \exp \left(\frac{-a}{\mu} \right) \cdot \frac{a^{k-(j-1)}}{k - (j - 1)} \\
 &= \frac{1}{(k - (j - 1))!} \left(\frac{a}{\mu} \right)^{k-(j-1)} \exp \left(\frac{-a}{\mu} \right) \\
 &= \left[1 - \sum_{i=0}^{k-(j-1)-1} \frac{1}{i!} \left(\frac{a}{\mu} \right)^i \exp \left(\frac{-a}{\mu} \right) \right] - \left[1 - \sum_{i=0}^{k-(j-1)} \frac{1}{i!} \left(\frac{a}{\mu} \right)^i \exp \left(\frac{-a}{\mu} \right) \right] \\
 &= F(a; k - (j - 1)) - F(a; k - (j - 1) + 1) \\
 &= F(a; k - j + 1) - F(a; k - j + 2)
 \end{aligned}$$

A.1.4 Case 4 $k \geq 1, j = (k + 1)$

$$p_a(k, j) = \mathbb{P}(S_1 > a) = 1 - \mathbb{P}(S_1 \leq a) = 1 - F(a; 1)$$

A.1.5 All Other Cases

$$p_a(k, j) = 0$$

A.1.6 Summary

These results can be summarised as:

$$p_a(k, j) = \begin{cases} \mathbb{1}(j = 1) & \text{where } k = 0 \\ F(a; k) & \text{where } k \geq 1, j = 1 \\ F(a; k - j + 1) - F(a; k - j + 2) & \text{where } k \geq 1, 2 \leq j \leq k \\ 1 - F(a; 1) & \text{where } k \geq 1, j = (k + 1) \\ 0 & \text{otherwise} \end{cases} \quad (\text{A.1})$$

A.2 Expected Transition Cost

A.2.1 Case 1 $k \in \{0, 1\}$

$$R_a(k, j) = c_S a = \gamma a$$

A.2.2 Case 2 $k \geq 2, j = 1$

$$\begin{aligned}
 R_a(k, j) &= c_W \sum_{i=2}^k \mathbb{E} \left[\sum_{l=1}^{i-1} S_l \mid \sum_{n=1}^k S_n \leq a \right] + c_S a \\
 &= c_W \sum_{i=2}^k \sum_{l=1}^{i-1} \mathbb{E} \left[S_l \mid \sum_{n=1}^k S_n \leq a \right] + c_S a \\
 &= c_W \mathbb{E} \left[S_1 \mid \sum_{n=1}^k S_n \leq a \right] \sum_{i=2}^k (i-1) + c_S a \\
 &= \frac{c_W k(k-1)}{2} \mathbb{E} \left[S_1 \mid \sum_{n=1}^k S_n \leq a \right] + c_S a \\
 &= \frac{c_W(k-1)}{2} \mathbb{E} \left[\sum_{n=1}^k S_n \mid \sum_{n=1}^k S_n \leq a \right] + c_S a \\
 &= \frac{c_W G(a; k)(k-1)}{2} + c_S a \\
 &= (1 - \gamma) \cdot \frac{G(a; k)(k-1)}{2} + \gamma a
 \end{aligned}$$

A.2.3 Case 3 $k \geq 2, 2 \leq j \leq k$

$$\begin{aligned}
 R_a(k, j) &= c_W \sum_{i=1}^{j-2} a + c_W \sum_{i=2}^{k-(j-2)} \mathbb{E} \left[\sum_{l=1}^{i-1} S_l \middle| \sum_{n=1}^{k-(j-1)} S_n \leq a \right] + c_S a \\
 &= c_W a(j-2) + c_W \sum_{i=2}^{k-(j-2)} \sum_{l=1}^{i-1} \mathbb{E} \left[S_l \middle| \sum_{n=1}^{k-(j-1)} S_n \leq a \right] + c_S a \\
 &= c_W a(j-2) + c_W \mathbb{E} \left[S_1 \middle| \sum_{n=1}^{k-(j-1)} S_n \leq a \right] \sum_{i=2}^{k-(j-2)} (i-1) + c_S a \\
 &= c_W a(j-2) + \frac{c_W [k - (j-1)] [k - (j-2)]}{2} \mathbb{E} \left[S_1 \middle| \sum_{n=1}^{k-(j-1)} S_n \leq a \right] + c_S a \\
 &= c_W a(j-2) + \frac{c_W [k - (j-2)]}{2} \mathbb{E} \left[\sum_{n=1}^{k-(j-1)} S_n \middle| \sum_{n=1}^{k-(j-1)} S_n \leq a \right] + c_S a \\
 &= c_W \left[a(j-2) + \frac{G(a; k - (j-1)) [k - (j-2)]}{2} \right] + c_S a \\
 &= (1 - \gamma) \left[a(j-2) + \frac{G(a; k - (j-1)) [k - (j-2)]}{2} \right] + \gamma a
 \end{aligned}$$

A.2.4 Case 4 $k \geq 2, j = (k+1)$

$$R_a(k, j) = c_W \sum_{i=1}^{j-2} a + c_S a = c_W a(j-2) + c_S a = (1 - \gamma) \cdot a(j-2) + \gamma a$$

A.2.5 Summary

These results can be summarised as:

$$\begin{aligned}
 & R_a(k, j) \\
 = & \begin{cases} \gamma a & \text{where } k \in \{0, 1\} \\ (1 - \gamma) \cdot \frac{G(a; k)(k-1)}{2} + \gamma a & \text{where } k \geq 2, j = 1 \\ (1 - \gamma) \left[a(j-2) + \frac{G(a; k-(j-1))[k-(j-2)]}{2} \right] + \gamma a & \text{where } k \geq 2, 2 \leq j \leq k \\ (1 - \gamma) \cdot a(j-2) + \gamma a & \text{where } k \geq 2, j = (k+1) \end{cases} \\
 & \hspace{25em} (\text{A.2})
 \end{aligned}$$

Bibliography

- Bailey, Norman TJ (1952). “A study of queues and appointment systems in hospital out-patient departments, with special reference to waiting-times”. In: *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 185–199.
- Rockart, John F and Paul B Hofmann (1969). “Physician and patient behavior under different scheduling systems in a hospital outpatient department”. In: *Medical Care* 7.6, pp. 463–470.
- Gupta, Ishwar, Juan Zoreda, and Nathan Kramer (1971). “Hospital manpower planning by use of queueing theory”. In: *Health Services Research* 6.1, pp. 76–82.
- Walter, SD (1973). “A comparison of appointment schedules in a hospital radiology department”. In: *British Journal of Preventive & Social Medicine* 27.3, pp. 160–167.
- Kao, Edward PC and Grace G Tung (1981). “Bed allocation in a public health care delivery system”. In: *Management Science* 27.5, pp. 507–520.
- O’Keefe, Robert M (1985). “Investigating outpatient departments: Implementable policies and qualitative approaches”. In: *Journal of the Operational Research Society* 36.8, pp. 705–712.
- Goldsmith, Jeff (1989). “A radical prescription for hospitals”. In: *Harvard Business Review*.
- Pegden, Claude Dennis and Matthew Rosenshine (1990). “Scheduling arrivals to queues”. In: *Computers & Operations Research* 17.4, pp. 343–348.
- Babes, Malika and GV Sarma (1991). “Out-patient queues at the Ibn-Rochd health centre”. In: *Journal of the Operational Research Society* 42.10, pp. 845–855.

- Ho, Chrwan-Jyh and Hon-Shiang Lau (1992). "Minimizing total cost in scheduling outpatient appointments". In: *Management Science* 38.12, pp. 1750–1764.
- Stein, William E and Murray J Côté (1994). "Scheduling arrivals to a queue". In: *Computers & Operations Research* 21.6, pp. 607–614.
- Huang, Fenghui and Mong Hou Lee (1996). "Using simulation in out-patient queues: A case study". In: *International Journal of Health Care Quality Assurance* 9.6, pp. 21–25.
- Bennett, Joanne C and DJ Worthington (1998). "An example of a good but partially successful OR engagement: Improving outpatient clinic operations". In: *Interfaces* 28.5, pp. 56–69.
- Cayirli, Tugba and Emre Veral (2003). "Outpatient scheduling in health care: A review of literature". In: *Production and Operations Management* 12.4, pp. 519–549.
- Mondschein, Susana V and Gabriel Y Weintraub (2003). "Appointment policies in service operations: A critical analysis of the economic framework". In: *Production and Operations Management* 12.2, pp. 266–286.
- DeLaurentis, Po-Ching et al. (2006). "Open access appointment scheduling - An experience at a community clinic". In: *IIE Annual Conference*. Institute of Industrial Engineers.
- Green, Linda (2006). "Queueing analysis in healthcare". In: *Patient flow: reducing delay in healthcare delivery*. Springer, pp. 281–307.
- Mendel, Sharon (2006). "Scheduling arrivals to queues: A model with no-shows". MA thesis. Tel-Aviv University.
- Fiems, Dieter, Ger Koole, and Philippe Nain (2007). "Waiting times of scheduled patients in the presence of emergency requests". In: *Technisch Rapport*.
- Fomundam, Samuel and Jeffrey W Herrmann (2007). *A survey of queuing theory applications in healthcare*. University of Maryland, The Institute for Systems Research.