

Dynamic Scheduled Queues

Masters Thesis

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	5
2	Literature Review	6
3	Dynamic Schedule	9
3.1	Aim	9
3.2	Assumptions	9
3.3	List of Variables	10
3.4	Objective Function	10
3.4.1	Base Case	11
3.5	Transition Probability	12
3.6	Expected Transition Cost	13
4	Conclusion	15
A	Optimal Cost Derivation	16

Chapter 1

Introduction

Introduction goes here.

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 Roshenshine'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 scheduled 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

Dynamic Schedule

3.1 Aim

The objective 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 dynamically. Immediately after a customer arrives and begins waiting for service, the scheduler chooses the arrival time of the next customer.

3.2 Assumptions

To simplify this problem, need to make several assumptions:

- Customer service times are independent and identically distributed (iid)
- Each customer 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

3.3 List of Variables

μ	: mean service time of each customer
c_W	: cost of customer's waiting time per unit time
c_I	: cost of server's idle time per unit time
k	: current number of customers waiting
j	: number of customers waiting 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 waiting and n customers remaining to be scheduled
$C_n(a, k)$: the expected cost of having k customers waiting, 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 waiting to j customers waiting 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 waiting to j customers waiting after a time units if the next customer is scheduled to arrive after a time units

3.4 Objective Function

The state (k, n) refers to k customers in the queue waiting for service and n customers remaining to be scheduled. The time the next customer is scheduled to arrive is a , and the number of customers waiting for service immediately after that customer's arrival is j . The expected cost at the current state is a function of the expected cost involved in transitioning to the next state, the expected cost at the next state and the probability of transitioning to the next state over all possible next states.

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

$$C_n^*(k) = \min_{a \geq 0} C_n(a, k) = \min_{a \geq 0} \left[\sum_{j=1}^{k+1} p_a(k, j) \left(R_a(k, j) + C_{n-1}^*(j) \right) \right] \quad (3.1)$$

Equation 3.1 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

$$C_n^*(k) = C_n(a^*, k) = \min_{a \geq 0} C_n(a, k) \quad (3.2)$$

It is reasonably intuitive that the minimum cost cannot occur at $a = \infty$. As $a \rightarrow \infty$, the probability that the server becomes idle converges to 1. In addition, the expected idle time of the server converges to ∞ . As the cost of the server's idle time c_I is strictly positive, the overall expected cost must also converge to ∞ as $a \rightarrow \infty$. Thus, $\lim_{a \rightarrow \infty} C_n(a, k) = \infty$.

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

$$\mathcal{A} = \{0\} \cup \left\{ a > 0 : \frac{\partial}{\partial a} C_n(a, k) = 0 \right\} \quad (3.3)$$

Solving Equation 3.2 involves solving a nonlinear optimisation problem over a left-closed interval. 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 optimal policy can be found by solving

$$C_n^*(k) = \min_{a \in \mathcal{A}} C_n(a, k) \quad (3.4)$$

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 (k, n) . Of course, this method becomes more computationally inefficient, the larger the number of states.

3.4.1 Base Case

Finding the solution iteratively requires a solution for the base case where $n = 0$. If $n = 0$, there are no customers remaining to be scheduled, which implies the server will not be idle for the remaining of service. The cost of state $(k, 0)$ (i.e., the base case) is thus the summation of the waiting cost of the k customers in the queue.

Let w_i be the expected waiting time of the customer that is currently in position i in the queue, and c_W be the cost of the customers' waiting time per unit time. The cost of the base case is thus given by

$$C_0(k) = \sum_{i=1}^k c_W w_i = c_W \sum_{i=1}^k \mu i = \frac{c_W \mu k(k+1)}{2} \quad (3.5)$$

3.5 Transition Probability

Let S_i be the service time of the customer that is currently in position i in the queue. Assume that the current state is (k, n) , then the service times S_1, \dots, S_k are iid (independent and identically distributed) random variables from an exponential distribution with mean μ .

The waiting time of the customer in the last position in the queue is

$$X = \sum_{i=1}^k S_i \sim \text{Erlang}(k, \mu) \quad (3.6)$$

which has the following pdf:

$$f(x; k) = \frac{1}{\mu \Gamma(k)} \left(\frac{x}{\mu} \right)^{k-1} e^{-\frac{x}{\mu}} \quad (3.7)$$

Let W_t be a Poisson Point Process such that $W_t \sim \text{Poisson}\left(\frac{t}{\mu}\right)$. Assume that the current state is (k, n) and there are no new arrivals to the queue. The number of customers served over the next t time units is given by $Y_t = \max(W_t, k)$. The probability that last customer in the queue waits longer than a time units is equal to the probability that Y_a is smaller than k , such that

$$\mathbb{P}(X > a) = \mathbb{P}(Y_a < k) \quad (3.8)$$

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:

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

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

Erlang distribution is defined as

$$F(a; k) = \begin{cases} 1 - \sum_{i=0}^{k-1} \frac{1}{i!} \left(\frac{a}{\mu}\right)^i e^{-\frac{a}{\mu}} & \text{where } a > 0 \\ 0 & \text{where } a = 0 \end{cases} \quad (3.9)$$

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.

$$p_a(k, j) = \begin{cases} \mathbb{1}(j = k + 1) & \text{where } a = 0 \\ \mathbb{1}(j = 1) & \text{where } a > 0, k = 0 \\ F(a; k) & \text{where } a > 0, k \geq 1, j = 1 \\ F(a; k - j + 1) - F(a; k - j + 2) & \text{where } a > 0, k \geq 1, 2 \leq j \leq k \\ 1 - F(a; 1) & \text{where } a > 0, k \geq 1, j = (k + 1) \\ 0 & \text{otherwise} \end{cases} \quad (3.10)$$

3.6 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 waiting time of all customers during the transition and the expected idle time of the server during the transition. Both the expected waiting and idle times depend on the conditional expectation of the Erlang distribution, which is derived here (for $a > 0$).

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

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

$$G(a; k) = \begin{cases} \mu k \cdot \frac{F(a; k+1)}{F(a; k)} & \text{where } a > 0 \\ 0 & \text{where } a = 0 \end{cases} \quad (3.11)$$

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

$$\lim_{a \rightarrow \infty} \frac{F(a; k+1)}{F(a; k)} = \frac{\lim_{a \rightarrow \infty} F(a; k+1)}{\lim_{a \rightarrow \infty} F(a; k)} = \frac{1}{1} = 1 \quad (3.12)$$

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

In a similar way to the transition probability, the expected transition cost is derived on a case by case basis. The per unit time cost of the customers' waiting time and the server's idle time are c_W and c_I respectively. The full derivation is included in the appendix and the final equation is presented here.

$$R_a(k, j) = \begin{cases} 0 & \text{where } a = 0 \\ c_I a + \frac{G(a; k)(c_W(k+1) - 2c_I)}{2} & \text{where } a > 0, k \geq 0, j = 1 \\ c_W a(j-1) + \frac{c_W G(a; k-j+1)(k-j+2)}{2} & \text{where } a > 0, k \geq 1, 2 \leq j \leq (k+1) \end{cases} \quad (3.13)$$

Chapter 4

Conclusion

Conclusion goes here.

Appendix A

Optimal Cost Derivation

Derivation of optimal cost goes here.

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.
- Huarng, Fenghueih 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 Insitute for Systems Research.