

Dynamic Scheduling of Outpatient Appointments under Patient No-shows and Cancellations

Nan Liu • Serhan Ziya • Vidyadhar G. Kulkarni

*Department of Statistics and Operations Research, University of North Carolina, Hanes
Hall, CB #3260, Chapel Hill, NC 27599-3260, USA*

nliu@email.unc.edu • ziyas@email.unc.edu • vkulkarn@email.unc.edu

This paper develops a framework and proposes heuristic dynamic policies for scheduling patient appointments taking into account the fact that patients may cancel or not show up for their appointments. In a simulation study that considers a model clinic, which is created using data obtained from an actual clinic, we find that the heuristics proposed outperform all the other benchmark policies, particularly when the patient load is high compared with the regular capacity. We find that the Open Access policy, a popular scheduling paradigm proposed recently which calls for “meeting today’s demand today,” can be a reasonable choice when the patient load is relatively low, supporting earlier findings in the literature.

1. Introduction

Many firms strive to match demand and supply in the presence of uncertainty. Production systems deal with randomness in demand by keeping inventories. However, that is not an option for service systems since service capacity cannot be stored. Instead, service systems like patient clinics, repair shops, and hair salons, regulate demand through appointments. Appointment systems have two objectives: (i) provide a better service to customers by assigning them a very short time window during which they are guaranteed to get a service, and (ii) protect the system from daily fluctuations in demand which can lead to an inefficient system with low utilization levels in some days and overloads in others. However, appointment schedules do not resolve all the uncertainties in daily demand. Some service systems such as those in health care suffer from high no-show and cancellation rates, which introduce additional layers of uncertainty and can cause severe inefficiencies if not dealt with appropriately. The objective of this paper is to develop a framework and introduce a method

for the dynamic scheduling of patient appointments recognizing the possibility that patients can cancel their appointments or simply not show up for their appointments.

Past research shows that the longer the *appointment delay*, defined as the time between the day a patient requests an appointment and her actual appointment date, the higher the chances that she will cancel or not show up (Gallucci et al. (2005)). This suggests an obvious way of minimizing no-shows and cancellations: ask patients to come right away or make appointment requests on the day they want to be seen. This is called an *open access* (OA) or *advanced access* policy (see, e.g. Murray and Tantau (2000)) and of late it has become a popular paradigm in practice and the subject of active research. Several authors report on their experiences in implementing OA, both positive and negative. See Murray et al. (2003), Solberg et al. (2004), Belardi et al. (2004), and Dixon et al. (2006). Some practitioners strongly advocate OA (e.g., Murray and Tantau (2000)) while there are some who are strongly against it (e.g., Lamb (2002)). However, there seems to be an agreement on the fact that for OA to have at least a chance to work, demand and supply (capacity) need to be “in balance.” Simply ensuring that average demand is less than supply is not sufficient for OA to work. Because of the stochastic nature of the daily demand, if average patient demand is not sufficiently low relative to the capacity, this will result in a high frequency of days in which daily demand exceeds the regular daily capacity. The clinic has to deal with this excess demand in some way (e.g., by working overtime, squeezing in additional appointments during the day, delegating some work to nurses) that will cause the clinic to incur overtime costs and/or reduce the quality of the service provided to the patients. Therefore OA is unlikely to be sustainable for a clinic observing demand levels that are close to its capacity. Although it is in general difficult to define exactly what it means for supply and demand to be “in balance,” Green and Savin (2008) use a queueing model to provide some answers by developing a method to help determine the largest panel size sustainable for a physician using OA. They also find that panel sizes typically need to be much smaller than what is required for the queueing-theoretic stability of the system, i.e., long-run demand rate to be less than the long-run supply rate.

Even though the critics of OA have many concerns other than possible demand/supply imbalance, reported success stories strongly support the case in favor of OA since most clinics that implement OA report significant improvements in various measures of service quality when demand and supply are in balance. The problem, however, is that it appears that keeping the demand at a level that is necessary for an effective OA implementation may

not be possible for many clinics and the number of such clinics may even increase in the future. Some of the recent articles published in academic journals and in the news media discuss the physician shortage problem that is currently being felt especially in some of the rural areas within the United States (e.g., upstate New York), and this is expected to grow significantly worse in the next 10-15 years unless action is taken to increase physician supply (see, e.g. Blumenthal (2004), Cauchon (2005), York (2007), Arvantes (2007)). One of the more problematic states is Massachusetts, which, interestingly, is the state with the highest doctors per capita according to U.S. Census Bureau (2008). Yet, according to the Massachusetts Medical Society (2007), the state has severe or critical shortages in several specialty areas including family practice and internal medicine. The society's reports from the last five years have already indicated increasing levels of physician shortage, but since the middle of 2007, patient demand has started to increase at an even faster rate as a result of the state's mandate on its residents to have health insurance, which has been effective since July 1, 2007. A New York Times article published in April 2008 reports that since the law took effect, about 340,000 of the approximately 600,000 uninsured have gained health coverage, and as a result clinics around the state have started to admit more patients by stretching their regular capacities (Sack (2008)). One of the family physicians the article mentions has a panel size of 3,000, which is well above the number suggested by Green and Savin for a potentially successful OA implementation. It is no surprise that patients of this particular physician have to wait for more than a year for a physical at the time of the article's publication.

These reports and articles do not discredit OA as an efficient solution for clinics that can keep their demand and supply in balance. However, they clearly indicate that a strict implementation of OA cannot be a universal solution to the appointment scheduling problem since some of the clinics will be overloaded out of obligation and some by choice. For such systems, scheduling some of the appointments to a not-so-distant future can relieve some of the stress of having to keep up with the demand on a daily basis while not causing no-show rates to increase in any significant way. Distributing the demand over several days has the obvious benefit of having a more regular daily load on the system, thereby reducing the possibility and/or the severity of daily overloads. However, distributing demand over days will have the unavoidable consequence of increasing no-shows and cancellations. This is the basic trade-off we are dealing with in this paper. There are two policies at the two ends of the policy spectrum. On the one side is the OA, which leads to a minimal no-show rate at

the expense of frequent/severe daily overloads; on the other side is a policy that schedules appointments so that daily overloads are kept at a minimum, which however causes clinics to suffer from unavoidable no-shows. A clinic’s choice of a scheduling policy would depend on a variety of factors that determine its sensitivity towards no-shows, flexibility in adjusting its physician capacity, willingness to work overtime, and/or willingness to overbook and possibly cram in more patients depending on the daily load. For example, a physician working for his/her own private practice might not be bothered too much by working overtime and thus the overtime “cost” for such a clinic might be lower than what other clinics would typically face. Similarly, some clinics might prefer to see as many patients as possible and thus choose to overbook, while some others might refrain from that and choose to provide their patients with shorter waiting times and keep them more satisfied. The question is how to determine the scheduling policy in response to such diverse preferences. This paper provides some answers to this question.

More specifically, the objective of this paper is to develop dynamic methods that help assign an appointment date to each patient depending on the clinic’s appointment schedule at the time of the patient’s call. We first formulate the problem as a Markov decision process (MDP). We have chosen to use a model that makes it possible to estimate various parameters using data that are typically available for most clinics. The model takes the following as an input: the expected “net reward” of serving some x number of patients on a day with some z scheduled appointments at the beginning of the day, and the cancellation and no-show probability distributions. The objective of the MDP model is to maximize the long-run average net “reward” for the clinic.

In theory, one can solve this MDP to optimality using one of the standard solution methods. However, as we demonstrate in Section 3, that is not practical since the system state space is very large even for relatively small, toy systems. Thus, we propose heuristic methods instead. The heuristic methods are developed using a known technique that employs a single step of the policy improvement algorithm on a static (state-independent) policy. The static policy is ideally a “good” policy if not the optimal among the class of state-independent policies so that the dynamic policy to be obtained after the policy improvement is even better. Two static policies that we use are the OA policy and the two-day probabilistic scheduling policy. For either one of these two static policies, the policy improvement step gives an index policy that operates as follows: when an appointment request comes in, an index value is computed for each day in the scheduling horizon and the appointment is scheduled

for the day with the highest index value. We show in Section 4 that the indices can easily be computed.

We evaluate the performances of the heuristic policies in an extensive simulation study. The model clinic that is used in the simulation study is created by using data provided to us by the Department of Family Medicine at the University of North Carolina (UNC). These data help us estimate the no-show and cancellation probability distributions. The simulation results indicate that the proposed heuristics significantly outperform other benchmark heuristics especially when the system is highly loaded. As expected, the OA policy performs reasonably well when average demand is below daily regular capacity but it performs very poorly under high demand.

The remainder of this paper is organized as follows. Section 2 reviews the relevant literature. In Section 3, we introduce our MDP model for appointment scheduling, and in Section 4 we describe the general form of the proposed heuristics. Section 5 gives detailed descriptions of the heuristic policies that we propose as well as other benchmark heuristics. In Section 6, we give more precise descriptions of the heuristic methods for two special cases and prove that they are optimal under certain conditions. In Section 7, we describe our “model clinic,” which we use in our simulation study and discuss how we estimated various model parameters. We describe the simulation study and report the results in Section 8. Finally, we provide our concluding remarks in Section 9.

2. Literature Review

The Operations Research (OR) literature on outpatient appointment scheduling is extensive. A recent survey paper by Gupta and Denton (2008) discusses the main practical issues related to appointment scheduling, provides a review of the state of the art in modeling and optimization, and points to future directions. One classification that Gupta and Denton make regarding research on appointment scheduling is with respect to the type of waiting modeled: *direct* and *indirect*.

As Gupta and Denton indicate, most of the existing research has concentrated on *direct* waiting times, the times that patients spend waiting in the clinic on the day of their appointments from their arrival until their service. This work typically aims to minimize the expected “cost” for a day, which is a function of patients’ direct waiting times, and the physician’s idle time and overtime. In this body of work, typical decision variables are the

number of appointment intervals, the length of each interval, and the number of patients assigned to each interval etc. We refer the reader to Cayirli and Veral (2003) for a comprehensive review of this literature pre-2003. More recent work in this line of research includes Denton and Gupta (2003), Robinson and Chen (2003), and Klassen and Rohleder (2004). These articles focus on the determination of appointment times for a sequence of punctual patients (jobs) with random service times, and their objective is to balance server idling, customer waiting and tardiness (overtime) costs. On the other hand, LaGanga and Lawrence (2007) and Muthuraman and Lawley (2008) study how to use overbooking to compensate for patient no-shows in an appointment system so as to improve the overall performance of the clinic. Our formulation and that of Muthuraman and Lawley (2008) are similar in that they both consider sequential scheduling of the patients as they call for appointments. However, while we are interested in determining the appointment day for each incoming patient, Muthuraman and Lawley are interested in determining in which particular time slot in a service session the appointment should be scheduled. Furthermore, Muthuraman and Lawley assume that using past data the clinic has identified the correlation between various patient attributes and their no-show probabilities and use these patient attributes in making scheduling decisions. In our formulation, the clinic does not differentiate among the patients in terms of their attributes, but unlike the model of Muthuraman and Lawley (2008), it does take into account the fact that no-show probabilities depend on the appointment delays.

Few articles deal with *indirect* waiting times, which correspond to appointment delays in this paper, and which refer to the times between the days patients call for an appointment and the actual appointment dates. Gupta and Denton (2008) point to several difficulties of modeling indirect waiting, which might be part of the reason why work has been rather limited. Existing articles typically focus on the question of how many patients to admit or whether or not to admit a given patient for a particular day given the system state (i.e., current patient backlog). Patrick et al. (2008) study a dynamic multi-priority patient scheduling problem and develop cost-effective booking policies that meet waiting time targets for the patients. Gupta and Wang (2008) study a clinic capacity management problem using a model where patients' physician and time preferences are explicitly formulated, the decision is whether or not an appointment request should be accepted upon its arrival, and the objective is to maximize the revenue obtained on a given day. Our paper also belongs to this stream of research since we also deal with patients' indirect waiting times. However, unlike the articles above, our model takes into account the possibility that patients might cancel or

not show-up for their appointments. As Gupta and Denton (2008) discuss, “indirect patient waiting” and “late cancelations and no-shows” largely remain as open research challenges, and in particular no prior work has explicitly studied appointment scheduling decisions in a model that relates overbooking decisions to no-show and cancellation rates. To the best of our knowledge, our paper is the first such work that proposes dynamic appointment scheduling policies using a model that explicitly takes into account patient no-shows and cancellations.

Several recent articles have investigated the open access policy. Kopach et al. (2007) use discrete event simulation to investigate the performance of OA under various settings. One key finding is that for clinics which predominantly use open access, offering provider care groups and overbooking appointments can help maintain the continuity of care provided to the patients, which is one of the important performance measures. Qu et al. (2007) identify the optimal percentage of appointment slots that a clinic should keep open within a session so as to maximize the throughput using a model where patients may not show up for their appointments. They also investigate the sensitivity of this optimal percentage to the provider capacity, patient no-show rates, and demand distributions. Green and Savin (2008) use a single-server queueing system to carry out capacity analysis for a clinic that uses the open access policy. In their model, each patient can be a no-show with a probability that depends on the patient’s waiting time for the appointment. The authors provide a method to calculate the largest panel size that the clinic can handle; in a way, they “define” what it means for demand and supply to be in balance for a clinic using open access. Robinson and Chen (2008) compare the performance of OA with that of a traditional appointment scheduling system. Assuming deterministic service times and fixed and homogeneous patient no-show probabilities, the authors identify some of the structural properties of the optimal traditional scheduling policies and develop bounds for the system performance. Through numerical analysis, they find that in most cases - that is unless patient waiting times have marginal weights in the objective function and patient no-show rates are too small - OA is more preferable to traditional scheduling systems. In short, these four articles deal with the design of an OA system and the comparison of OA with traditional scheduling policies. In this paper, even though our main objective is not to investigate OA specifically, we also provide some support to some of the findings of this earlier work by identifying conditions under which OA performs reasonably well compared with the heuristic policies we propose.

Finally, a number of articles outside the OR literature investigate patient no-shows em-

pirically. For example, see Oppenheim et al. (1979), Pesata et al. (1999), Moore et al. (2001), and Gallucci et al. (2005). All of these articles point to patient no-shows as being a significant problem in appointment scheduling and find that no-show rates depend on a variety of factors including race, gender, socioeconomic status etc. In particular, Gallucci et al. (2005) find that no-show and cancellation rates increase with the appointment delay. As we discuss in Section 7.1, we also find a similar relationship using data from the UNC clinic.

3. The Model

We consider a clinic where patients call to make appointments for a visit in the future or some time during the day of the call. Given the current appointment schedule, the administrative staff schedules each incoming request for an appropriate day and updates the schedule accordingly. We assume that patients do not have a strong preference for the date they want to be seen and thus accept the first appointment date offered by the staff. In Section 9, we discuss how the heuristics we propose can be used in cases where patients do have time preferences.

Let A^t denote the number of appointment requests that arrive on day t . We assume that $\{A^t, t = 1, 2, \dots\}$ is a sequence of independent and identically distributed (i.i.d.) random variables. As defined earlier, the *appointment delay* for a patient is the time between the day the patient requests an appointment and her actual appointment date. The appointment delay for a patient is zero if the appointment is scheduled on the same day the patient calls. The clinic uses a scheduling horizon of length T so that no patient has an appointment delay that is larger than T . Hence, a patient requesting an appointment on day t will be scheduled on one of the days $t, t + 1, \dots, t + T$. For modeling convenience, we assume that all appointment requests are received at the beginning of the day so that all scheduling decisions on a given day t are made given the realization a^t of A^t . We shall see below that this assumption is needed to derive the heuristic policy, but is not needed when implementing the heuristic.

For any given day t , we define type (i, j) patients as those who called on day $t - i$, were given appointments for day $t + j$, and have not canceled their appointments by the beginning of day t . All patients of type (i, j) have an appointment delay of $i + j$ days and thus we must have $i + j \leq T$. Note that a given patient's type changes with time. For example, today's type (i, j) patient is tomorrow's type $(i + 1, j - 1)$ patient assuming she does not cancel.

Let X_{ij}^t denote the number of type (i, j) patients at the beginning of day t and define $\mathbf{X}^t = \{X_{ij}^t : 1 \leq i + j \leq T, i = 1, 2, \dots, T\}$ as the vector of the number of patients of each type in the schedule at the beginning of day t . In the rest of the paper, we refer to \mathbf{X}^t as the *backlog* or the *schedule* on day t , based on which the clinic schedules the incoming appointment requests.

There are three possible outcomes for each appointment made. The patient may show up for her appointment, cancel her appointment on or before the day of the appointment, or she may not cancel but simply not show up for her appointment. We assume that each patient's behavior is independent of that of the other patients and the arrival process $\{A^t, t = 1, 2, \dots\}$. Past research (e.g., Gallucci et al. (2005)) as well as our own analysis clearly show that the longer the appointment delay for a patient, the higher the chances that she will cancel or that she will be a no-show if she does not cancel. In order to capture this relationship, we formulate the cancellation/no-show behavior as follows: We assume that each patient cancels her appointment at some random time in the future, which can possibly be beyond the patient's appointment day. We use T_c to denote the generic random variable that represents the time between the day a patient calls for an appointment and the day she decides to cancel it. A patient who has not canceled on or before her appointment day (which happens if T_c for this particular patient is larger than the patient's appointment delay) may or may not show up for her appointment. Let "S" and "NS" represent the "show" and "no-show" events for a particular patient, respectively and define

$$\alpha_{ij} = \mathbf{P}(T_c \geq i + j + 1, S | T_c \geq i), \quad (1)$$

$$\beta_{ij} = \mathbf{P}(T_c \geq i + j | T_c \geq i). \quad (2)$$

Thus, α_{ij} is the probability that a patient who is currently of type (i, j) will show up for her appointment and β_{ij} is the probability that a patient who is currently of type (i, j) will not cancel her appointment before her appointment day. If there are n type (i, j) patients today, out of these n patients, the number of those who will not have canceled by the morning of their appointment day is a Binomial random variable with parameters n and β_{ij} . Similarly, the number of these patients who will show up for their appointment is a Binomial random variable with parameters n and α_{ij} .

We assume that events occur in the following order on each day. First, new patients call in and are given appointments. During the day, some patients cancel their appointments, and some do not show up for their appointments. At the end of the day, the clinic makes

an “expected net reward” of $r(x, z)$ if z patients were scheduled on that day and the clinic ended up serving x of these patients during the day. For generality, we do not specify any particular form for $r(\cdot, \cdot)$ and thus each clinic can choose the function that best captures the circumstances it is in and its own valuations of different situations (i.e., the cost of reduced quality of service given to the patients if they are served in overtime slots). One special case of $r(\cdot, \cdot)$ is

$$r(x, z) = \eta(x) - w(z) \tag{3}$$

where $\eta(x)$ can be seen as the “expected reward” and is possibly linear or more generally concave in x and $w(z)$ is the total “expected cost” for a day if there are z patients scheduled at the beginning of the day. Here, the idea of making the expected cost a function of z as opposed to x makes it possible to capture the possibility that the clinic will plan according to the number of scheduled appointments (e.g., staffing cost). Any “cost” that depends on the number of patients who show up can be included in the $\eta(\cdot)$ function. These reward and cost functions can be estimated using models that were developed for scheduling appointments within a day. For example, Denton and Gupta (2003) developed a model to schedule patients over a day so as to minimize the total cost of patient waiting, staff idling and overtime. Their work provides a way to estimate the cost based on the number of scheduled appointments.

Note that the general form of the function $r(\cdot, \cdot)$ allows for more sophisticated choices than the special case given in (3). For instance, as reported by Moore et al. (2001), not all no-show slots are wasted since walk-in patients fill in some of these empty slots. According to the authors’ study, only 12.2% of the slots are left unfilled by the end of the day, and on the average walk-in or triage patients help recover 89.5% of the costs of no-show patients. Thus, in some cases it might be reasonable to assume that the “reward” not only depends on x , the number of patients who show up, but also $z - x$, the number of no-show patients on a given day. This can be easily captured with our reward/cost formulation.

The objective of the clinic is to schedule arriving appointment requests so that the long-run average expected net reward is maximized. This problem can be modeled as an MDP, where the decision epochs are the times right after the appointment requests arrive every day and the system state at decision epoch t is given by (A^t, \mathbf{X}^t) . Let Y_j^t represent the number of patients who make their requests on day t and are given appointments for day $t + j$. (For example, Y_0^t is the number of patients who are given same-day appointments on day t .) Thus $\mathbf{Y}^t = \{Y_j^t : j = 0, 1, \dots, T\}$ is the scheduling “action” taken on day t . Given that there are A^t appointment requests, the set of actions available on day t is

$\{\mathbf{Y}^t : \sum_{i=0}^T Y_i^t = A^t, Y_i^t \in \mathbb{Z}_+, i = 0, 1, \dots, T\}$ where \mathbb{Z}_+ represents the set of all nonnegative integers. Note that it is straightforward to include “rejection” of the appointment request as another action. However, to keep the presentation simpler, we assume that rejection is not an option. Later, we discuss how our heuristic policies would change if rejection were an option.

Let $B(n, p)$ represent a Binomial random variable with parameters n and p . Given (A^t, \mathbf{X}^t) and the scheduling action \mathbf{Y}^t , \mathbf{X}^{t+1} can be characterized as follows:

$$X_{ij}^{t+1} \stackrel{d}{=} \begin{cases} B(Y_{j+1}^t, \beta_{01}), & i = 1, j = 0, \dots, T-1, \\ B(X_{i-1, j+1}^t, \beta_{i-1, 1}), & i \geq 2, 0 \leq j \leq T-i, \end{cases} \quad (4)$$

where $\stackrel{d}{=}$ denotes the equality in distribution. Let

$$b(n, k, p) = \binom{n}{k} p^k (1-p)^{n-k}$$

denote the probability mass function for a $B(n, p)$ random variable. Then, the transition probabilities given the action $\mathbf{Y}^t = \mathbf{y}^t$ can be expressed as follows:

$$\begin{aligned} & \mathbf{P}[(A^{t+1}, \mathbf{X}^{t+1}) = (a^{t+1}, \mathbf{x}^{t+1}) | (A^t, \mathbf{X}^t) = (a^t, \mathbf{x}^t), \mathbf{Y}^t = \mathbf{y}^t] \\ &= \mathbf{P}(A^{t+1} = a^{t+1}) \prod_{i=1, \dots, T, j=0, \dots, T-i} P_{ij}(\mathbf{x}^{t+1}, \mathbf{x}^t, \mathbf{y}^t), \end{aligned}$$

where

$$P_{ij}(\mathbf{x}^{t+1}, \mathbf{x}^t, \mathbf{y}^t) = \begin{cases} b(y_{j+1}^t, x_{ij}^{t+1}, \beta_{01}), & i = 1, j = 0, \dots, T-1, 0 \leq x_{ij}^{t+1} \leq y_{j+1}^t \\ b(x_{i-1, j+1}^t, x_{ij}^{t+1}, \beta_{i-1, 1}), & i \geq 2, 0 \leq j \leq T-i, 0 \leq x_{ij}^{t+1} \leq x_{i-1, j+1}^t. \end{cases}$$

Let $U_i^t, i = 0, 1, \dots, T$ denote the number of patients who call on day $t-i$ and show up for their appointments on day t . Then we have

$$U_i^t \stackrel{d}{=} \begin{cases} B(Y_0^t, \alpha_{00}), & i = 0, \\ B(X_{i0}^t, \alpha_{i0}), & i = 1, 2, \dots, T. \end{cases}$$

We define $c_0((A^t, \mathbf{X}^t), \mathbf{Y}^t)$ to be the net reward obtained on day t given A^t, \mathbf{X}^t , and \mathbf{Y}^t . Then,

$$c_0((A^t, \mathbf{X}^t), \mathbf{Y}^t) = r \left(\sum_{i=0}^T U_i^t, Y_0^t + \sum_{i=1}^T X_{i0}^t \right).$$

Now, consider a scheduling policy f , and let $\phi_f(a, \mathbf{x})$ be the long-run expected average net reward under policy f given the initial state $A^1 = a$ and $\mathbf{X}^1 = \mathbf{x}$, i.e.

$$\phi_f(a, \mathbf{x}) = \lim_{k \rightarrow \infty} \frac{\mathbf{E}_f[\sum_{t=1}^k c_0((A^t, \mathbf{X}^t), \mathbf{Y}^t) | (A^1, \mathbf{X}^1) = (a, \mathbf{x})]}{k}.$$

A scheduling policy f^* is said to be optimal if

$$\phi_{f^*}(a, \mathbf{x}) = \sup_f \phi_f(a, \mathbf{x}), \quad \forall a, \mathbf{x}.$$

In theory, one can solve this MDP problem using one of the standard procedures such as the policy improvement or value iteration algorithms. However, the formulation suffers significantly from the curse of dimensionality. To see that, suppose that the maximum number of appointment requests that can possibly be received on a single day is $N < \infty$. Then one can show that the number of states for the MDP formulation we have given above would be $(N + 1) \prod_{i=1}^T \sum_{k=0}^N \binom{k+i-1}{i-1}$. Note that even when $N = T = 5$, this number equals 1.34×10^9 and thus determining the optimal policy is not practically feasible for any realistically sized problem. Therefore it is of interest to develop heuristic scheduling methods that are efficient and perform well. We do this in the following section.

4. Policy Improvement Heuristics

In this section, we develop a heuristic dynamic appointment scheduling policy based on the idea of applying a single step of the policy improvement algorithm starting with a “good” initial policy. The idea is that since the policy improvement algorithm is generally believed to converge fast, applying a single step on an already “good” policy could give a policy that might be a reasonable substitute for the optimal dynamic policy. As we will see in this section, using a static policy as the initial policy makes it possible to fully characterize the policy obtained after the application of the policy improvement step so that the heuristic policy we propose will be easily implementable. In particular, it will not require the application of the policy improvement step for each instance of the appointment scheduling problem. This heuristic development technique has previously been used in such diverse areas as routing for parallel queues (see, e.g., Krishnan (1990), Opp et al. (2005) and Argon et al. (2009)) and dynamic kidney allocation (Zenios et al. (2000)), but to our knowledge, not within the context of appointment scheduling.

The procedure of developing policy improvement heuristics is as follows. Consider a state-independent policy that schedules each patient according to a stationary probability distribution $\mathbf{p} = (p_0, p_1, \dots, p_T)$ where p_i is the probability that the patient is given an appointment for the i th day from today. We call this policy *Probabilistic Static Policy* (PSP). Consider a policy $\pi_{\mathbf{y}}$ that takes action $\mathbf{y} \equiv \{y_j\}_{j=0}^T$ at the beginning of today, i.e., schedules

y_j patients on the j th day from today for $j \in \{0, 1, \dots, T\}$, but from tomorrow on uses PSP. Define $\Delta((a, \mathbf{x}), \mathbf{y}, \mathbf{p})$ to be the difference in total expected net rewards over an infinitely long period of time by following policy $\pi_{\mathbf{y}}$ rather than using PSP all along given the initial state (a, \mathbf{x}) . To conduct a one-step policy improvement, we maximize $\Delta((a, \mathbf{x}), \mathbf{y}, \mathbf{p})$ with respect to \mathbf{y} for each state (a, \mathbf{x}) , and the resulting optimal \mathbf{y} , denoted as $\mathbf{y}^*(a, \mathbf{x})$, specifies the heuristic scheduling policy for state (a, \mathbf{x}) . We call this policy *Heuristic Dynamic Policy* (HDP). In the following, we first give an expression for $\Delta((a, \mathbf{x}), \mathbf{y}, \mathbf{p})$. Then, we use this expression to determine HDP. For a simple example of how to develop policy improvement heuristics, see Example 3.6.2 of Tijms (1994).

4.1 Improving the Probabilistic Static Policy

Since patients can be given appointments at most T days in advance, the action \mathbf{y} taken today (say day 0) will only affect the schedule on days $0, 1, \dots, T$. Thus $\pi_{\mathbf{y}}$ and PSP are guaranteed to give stochastically the same appointment schedule from day $T + 1$ onwards. In order to compute $\Delta((a, \mathbf{x}), \mathbf{y}, \mathbf{p})$, we only need to find the difference between the total expected net rewards for days $0, 1, \dots, T$ under these two policies. Let $R_{\pi_{\mathbf{y}}}((a, \mathbf{x}), \mathbf{y}, \mathbf{p})$ and $R_{PSP}((a, \mathbf{x}), \mathbf{p})$ denote the total expected net reward accumulated over days $0, 1, \dots, T$ under policy $\pi_{\mathbf{y}}$ and PSP, respectively, given the initial state (a, \mathbf{x}) . Then,

$$\Delta((a, \mathbf{x}), \mathbf{y}, \mathbf{p}) = R_{\pi_{\mathbf{y}}}((a, \mathbf{x}), \mathbf{y}, \mathbf{p}) - R_{PSP}((a, \mathbf{x}), \mathbf{p}). \quad (5)$$

We first compute $R_{\pi_{\mathbf{y}}}((a, \mathbf{x}), \mathbf{y}, \mathbf{p})$. For $1 \leq i \leq T$ and $0 \leq j \leq T - i$, let $V_{ij}(\mathbf{x})$ denote the number of patients who called for appointments i days before day 0 and will not cancel by the morning of their appointment on day j , and $W_{ij}(\mathbf{x})$ denote the number of these patients who show up for their appointments. Similarly, for $0 \leq j \leq T$, let $\bar{V}_j(\mathbf{y})$ denote the number of patients who call for appointments on day 0, are scheduled for day j and will not cancel by the morning of their appointment, and $\bar{W}_j(\mathbf{y})$ denote the number of these patients who show up for their appointments. Finally, for $1 \leq k \leq T$ and $k \leq j \leq T$, we define $\hat{V}_{kj}(\mathbf{p})$ to be the number of patients who will call for appointments on day k and will not have canceled their appointment by the morning of their appointment on day j , and we define $\hat{W}_{kj}(\mathbf{p})$ to be the number of these patients who will show up for their appointments. Then, since the cancellation/no-show behaviors of the patients are independent of each other, we know that each one of these random variables has a Binomial distribution. More precisely,

$$V_{ij}(\mathbf{x}) \stackrel{d}{=} B(x_{ij}, \beta_{ij}), \bar{V}_j(\mathbf{y}) \stackrel{d}{=} B(y_j, \beta_{0j}), \hat{V}_{kj}(\mathbf{p}) \stackrel{d}{=} B(A_{(k)}, p_{j-k} \beta_{0, j-k}),$$

and

$$W_{ij}(\mathbf{x}) \stackrel{d}{=} B(x_{ij}, \alpha_{ij}), \bar{W}_j(\mathbf{y}) \stackrel{d}{=} B(y_j, \alpha_{0j}), \hat{W}_{kj}(\mathbf{p}) \stackrel{d}{=} B(A_{(k)}, p_{j-k} \alpha_{0,j-k}),$$

where $A_{(k)}$ denotes the number of appointment requests that will be made on day k , $1 \leq k \leq T$.

To see why the distribution of $\hat{V}_{kj}(\mathbf{p})$ and $\hat{W}_{kj}(\mathbf{p})$ are as given above, note that starting tomorrow, the policy $\pi_{\mathbf{y}}$ switches to implementing PSP, and therefore a patient requesting appointment on day k will be given an appointment on day j (where $k \leq j$) with probability p_{j-k} . From (1) and (2), this particular patient will not cancel the appointment by the morning of the appointment with probability $\beta_{0,j-k}$, and will show up for the appointment with probability $\alpha_{0,j-k}$.

We can then write

$$R_{\pi_{\mathbf{y}}}((a, \mathbf{x}), \mathbf{y}, \mathbf{p}) = \sum_{j=0}^T f_j(y_j, \mathbf{x}, \mathbf{p}),$$

where $f_j(y_j, \mathbf{x}, \mathbf{p})$ is the expected net reward accumulated on day j ($j = 0, 1, \dots, T$), and is given by

$$f_j(y_j, \mathbf{x}, \mathbf{p}) = \mathbf{E} \left[r \left(\sum_{i=1}^{T-j} W_{ij}(\mathbf{x}) + \bar{W}_j(\mathbf{y}) + \sum_{k=1}^j \hat{W}_{kj}(\mathbf{p}), \sum_{i=1}^{T-j} V_{ij}(\mathbf{x}) + \bar{V}_j(\mathbf{y}) + \sum_{k=1}^j \hat{V}_{kj}(\mathbf{p}) \right) \right] \quad (6)$$

where we let $\sum_{k=1}^j \hat{W}_{kj}(\mathbf{p}) = \sum_{k=1}^j \hat{V}_{kj}(\mathbf{p}) = 0$ for $j = 0$.

Now, we need to determine \mathbf{y} that maximizes (5) while ensuring that $\sum_{i=0}^T y_i = a$. But then since $R_{PSP}((a, \mathbf{x}), \mathbf{p})$ does not depend on \mathbf{y} , it can be dropped from the optimization problem altogether. Thus, HDP can be determined by solving the following resource allocation problem for given (a, \mathbf{x}) and \mathbf{p} :

$$\begin{aligned} \max \quad & R_{\pi_{\mathbf{y}}}((a, \mathbf{x}), \mathbf{y}, \mathbf{p}) = \sum_{j=0}^T f_j(y_j, \mathbf{x}, \mathbf{p}) \\ \text{s.t.} \quad & \sum_{j=0}^T y_j = a, \\ & y_j \in \mathbb{Z}_+, \quad j = 0, 1, \dots, T. \end{aligned} \quad (7)$$

4.2 Description of the Heuristic Dynamic Policy

In its most general form, the heuristic policy we propose can be described as follows: first, pick a distribution \mathbf{p} , and if the system state is (a, \mathbf{x}) at the beginning of the day, solve Problem (7). Clearly, solving (7) is significantly simpler than solving the original MDP problem in general, but under certain conditions this is especially the case. In particular, we make the following assumption:

Assumption 1. The net reward function $r(x, z)$ is (i) increasing in x for fixed z and decreasing in z for fixed x , (ii) submodular, and (iii) jointly concave in x and z .

This assumption enforces reasonable conditions on the net reward function. In particular, if $r(\cdot, \cdot)$ is in the form of (3), Assumption 1 holds if the function $\eta(\cdot)$ is increasing and concave, capturing the (possibly) diminishing returns for additional patients, and the function $w(\cdot)$ is increasing and convex, capturing the (possibly) increasing marginal cost of each additional appointment.

As the following proposition states, Assumption 1 ensures that the objective function of Problem (7) is well-behaved. The proof is given in Appendix A of the Online Companion.

Proposition 1. Under Assumption 1, for given \mathbf{x} and \mathbf{p} , the function $f_j(y_j, \mathbf{x}, \mathbf{p})$ is concave in y_j .

From Proposition 1 it follows that when Assumption 1 holds, Problem (7) is a resource allocation problem with a separable concave objective function, which has been well-studied in the literature. In particular, its optimal solution can be found by a simple algorithm given below (see, e.g. Section 4.2 of Ibaraki and Katoh (1988)). To avoid trivialities, we assume $a > 0$.

1. **Initialization:** Set $n := 1$ and $y_j := 0$, $j = 0, 1, \dots, T$.

2. **Scheduling the n th patient of the day:** For each $j \in \{0, 1, \dots, T\}$, compute

$$I_j = I_j(y_j, \mathbf{x}, \mathbf{p}) = f_j(y_j + 1, \mathbf{x}, \mathbf{p}) - f_j(y_j, \mathbf{x}, \mathbf{p}), \quad (8)$$

and determine

$$j^* = \arg \max_{j \in \{0, 1, \dots, T\}} I_j.$$

Set

$$y_{j^*} := y_{j^*} + 1.$$

3. **Termination test:**

- if $n < a$, let $n := n + 1$ and return to step 2;
- otherwise, terminate the algorithm with \mathbf{y} being the optimal solution.

At each iteration, the algorithm simply assigns a patient to the day j with the largest index value I_j , i.e., the day that will bring the largest improvement in the objective function of (7). Note that since patients are scheduled one by one and the index values I_j 's do not depend on a , i.e., the total number of appointment requests on that day, we can relax our assumption that the clinic knows the total number of requests when scheduling patients. As new appointments are made, the clinic updates the appointment schedule, and then when a request comes in, the indices are calculated based on the updated information. Thus the HDP works as follows:

Heuristic Dynamic Policy

When an appointment request comes in, first for each day $j \in \{0, 1, 2, \dots, T\}$ calculate the corresponding index I_j using (8), where y_j is the number of patients who called so far today and were scheduled for the j th day from today and \mathbf{x} is the appointment schedule upon the arrival of this request. Then, determine $j^ = \arg \max_{j \in \{0, 1, \dots, T\}} I_j$, schedule the new appointment for day j^* and set $y_{j^*} = y_{j^*} + 1$.*

If the clinic also had the option of rejecting appointment requests, the resulting policy would have an additional index for rejection, say the index I_R . It is then straightforward to show that this index would simply be equal to zero regardless of the system state. Therefore, with the rejection option available, the description of the Heuristic Dynamic Policy would need to be updated so that the appointment is rejected if all indices I_j where $j \in \{0, 1, 2, \dots, T\}$ are negative. Otherwise, the appointment is scheduled on day j^* as described above.

5. Picking the Distribution \mathbf{p} for HDPs and Alternative Heuristics

In Section 4.2, we described the heuristic policy that we propose for a given distribution \mathbf{p} , but we did not specify how this distribution should be picked. One can come up with different heuristic policies by choosing different distributions, but here we describe and propose two particular choices, which are both easy to determine and yield easily implementable HDPs. Then, we describe other benchmark policies that are not based on the policy described in Section 4.2.

5.1 Picking the Distribution \mathbf{p}

Since the policy improvement heuristic is guaranteed to improve upon the initial policy, ideally, one would like to pick the optimal \mathbf{p} , i.e., the one that maximizes the long-run average reward among all static policies. Although there is no guarantee that a better static policy will lead to a better dynamic policy, as long as there is no reason to believe otherwise, it appears to be the most reasonable choice. However, finding the optimal \mathbf{p} is a significant challenge by itself, requiring the solution of a maximization problem with $T + 1$ variables, and with an objective function that becomes increasingly more difficult to calculate as T increases. Therefore, in this paper, we propose two simple alternatives.

Policy I: Open Access Policy (OAP): Open access policy is a static policy where $p_0 = 1$ and $p_1 = p_2 = \dots = p_T = 0$.

Policy II: Optimal Two-day Probabilistic Static Policy (OTPSP): This is the optimal probabilistic static policy with the restriction that $p_2 = p_3 = \dots = p_T = 0$ (p_0 and p_1 are picked optimally).

OAP is readily available, but finding OTPSP needs some explanation. We first derive the long-run average net reward under a general two-day probabilistic static policy (TPSP) as a function of p_0 (since $p_1 = 1 - p_0$), which is to be maximized with respect to p_0 to determine OTPSP. The TPSP is state-independent. Thus we pick a random day and derive an expression for the expected net reward for that day. Now, since $p_i = 0$ for $i \geq 2$, patients seen on a given day either made appointment on that day or the day before. Let \tilde{A}_0 and \tilde{A}_1 denote the number of appointment requests that arrived today and yesterday, respectively. Let $\tilde{V}_0(p_0)$ and $\tilde{V}_1(p_0)$ denote the number of patients who called in today and yesterday, respectively, to make appointments for today have not canceled their appointments by the morning of today, and $\tilde{W}_0(p_0)$ and $\tilde{W}_1(p_0)$ respectively denote the number of these patients who will show up for their appointments. Then,

$$\tilde{V}_0(p_0) \stackrel{d}{=} B(\tilde{A}_0, p_0), \tilde{V}_1(p_0) \stackrel{d}{=} B(\tilde{A}_1, (1 - p_0)\beta_{01}),$$

and

$$\tilde{W}_0(p_0) \stackrel{d}{=} B(\tilde{A}_0, p_0\alpha_{00}), \tilde{W}_1(p_0) \stackrel{d}{=} B(\tilde{A}_1, (1 - p_0)\alpha_{01}).$$

It follows that the long-run average net reward under TPSP with a same-day scheduling probability p_0 can be written as

$$R_{TPSP}(p_0) = \mathbf{E} \left[r \left(\tilde{W}_0(p_0) + \tilde{W}_1(p_0), \tilde{V}_0(p_0) + \tilde{V}_1(p_0) \right) \right], \quad (9)$$

and the OTPSP is obtained by solving the following optimization problem

$$\max_{0 \leq p_0 \leq 1} R_{TPSP}(p_0), \quad (10)$$

where $R_{TPSP}(p_0)$ is given in (9). This is an optimization problem with a single decision variable, and thus determining the optimal solution is relatively straightforward given the probability distribution for the daily appointment requests, \tilde{A}_0 and \tilde{A}_1 . Furthermore, under certain conditions on this distribution and the reward/cost parameters, the objective function is well-behaved making the determination of the optimal solution to the problem easier. (See Appendix A of the Online Companion for the proof.)

Proposition 2. Suppose that the number of appointment requests that arrive on a single day has a Poisson distribution and the reward function $r(\cdot, \cdot)$ is as given in (3). Then, under Assumption 1, the objective function of Problem (10) is concave.

Although there is no strong evidence in support of the number of daily appointment requests having a Poisson distribution, it is a common assumption in the literature. As argued in Robinson and Chen (2008), if the panel size of a clinic is N and each patient independently makes an appointment for any given day with a small probability p , then the total number of appointment requests arriving on any given day has a Binomial distribution with parameters N and p , which converges to a Poisson distribution (with mean Np) as N gets large. For the heuristics we propose, although Poisson assumption is by no means necessary, it is nevertheless convenient.

Based on these two static policies, i.e., OAP and OTPSP, we propose two dynamic policies:

Policy III: Improved Open Access Policy (Imp-OAP): This policy works as described in Section 4.2 with distribution \mathbf{p} picked as specified for OAP, i.e., $p_0 = 1$ and $p_1 = p_2 = \dots = p_T = 0$.

Policy IV: Improved Optimal Two-day Probabilistic Static Policy (Imp-OTPSP): This policy works as described in Section 4.2 with distribution \mathbf{p} picked as in the description of OTPSP, i.e. $\mathbf{p} = [p_0, 1 - p_0, 0, \dots, 0]$ where p_0 is the optimal solution to (10).

5.2 Alternative Heuristic Policies

In Section 8, we investigate the performances of the heuristic policies we described above. Since determining the optimal policy is not practically feasible because of the large state

space, we compare the performances of the policies we propose with those of other heuristic policies that mimic some of the scheduling principles that are followed in practice and can be intuitively expected to perform well at least under certain conditions. Here, we describe these policies which will serve as benchmark policies.

Policy V: Threshold Policy (TP): This policy schedules each new appointment for the earliest day with less than M patients already scheduled. If there are no days within the scheduling horizon that has less than M already scheduled, the new appointment is scheduled for the day with the fewest appointments, and ties are broken in the favor of the earliest day. Threshold M is a policy parameter. One reasonable choice for M is the regular daily capacity of the clinic.

Policy VI: Balanced Scheduling Policy (BSP): This policy schedules each new appointment for the day with the fewest appointments, and ties are broken in favor of the earliest day.

Policy VII: Random Scheduling Policy (RSP): This policy schedules each new appointment randomly for one of the days in the scheduling horizon. More precisely, it is a probabilistic static policy with $p_i = 1/(T + 1), i = 0, 1, \dots, T$.

6. Two Examples and the Optimality of Imp-OAP

To give the reader a better idea about how the policies OTPSP, Imp-OTPSP and Imp-OAP look, we study two simple examples that lead to indices that can be expressed explicitly. For both examples, we assume that the function $r(\cdot, \cdot)$ is in the form of (3) with $\eta(x) = \tau x$ where τ is a positive constant, $\mathbf{E}[A^t] = \mu$ and $\mathbf{E}[(A^t)^2] = \xi$ for $t = 1, 2, \dots$. In the first example, the function $w(z)$ is increasing linearly with z , which would be a reasonable assumption in cases where the clinic does not have to deal with very high patient loads. In the second example, the function $w(z)$ is assumed to be quadratic, which means the “marginal cost” of a patient is increasing with each additional scheduled patient, which is reasonable in cases where the clinic is relatively understaffed so that additional patients bring increasingly more burden on the clinic.

6.1 Linear Rewards and Costs

Suppose that $w(z) = \nu_1 z$ where $\nu_1 \geq 0$ is the cost per scheduled appointment at the beginning of a day. Then, the optimal p_0 (denoted as p_0^*) which solves problem (10) is given by

$$p_0^* = \begin{cases} 0 & \text{if } \tau(\alpha_{00} - \alpha_{01}) \leq \nu_1(1 - \beta_{01}), \\ 1 & \text{otherwise.} \end{cases} \quad (11)$$

We can then show that for Imp-OTPSP, the indices (8) simplify to

$$I_j = I_j(y_j, \mathbf{x}, \mathbf{p}) = \tau\alpha_{0j} - \nu_1\beta_{0j}, \quad (12)$$

which corresponds to the expected net reward of scheduling one more patient for day j . Notice that in this case the index I_j does not depend on \mathbf{x} or y_j , and thus the HDP becomes a deterministic static policy which schedules all appointment requests received today for the j^* th day from today where $j^* = \arg \max_{j \in \{0, 1, \dots, T\}} I_j$. For this example, Imp-OAP is the same as Imp-OTPSP, and as we prove in Theorem 1, they are in fact optimal. (See Appendix A of the Online Companion for the proof.)

Theorem 1. *If the functions $\eta(\cdot)$ and $w(\cdot)$ are both linear, then the dynamic policy Imp-OAP, described by the indices given by (12), is optimal among all policies.*

6.2 Linear Rewards and Quadratic Costs

Suppose that $w(z) = \nu_2 z^2$ where $\nu_2 \geq 0$. Define $\kappa_0 = \mu(\tau\alpha_{01} - \nu_2\beta_{01}) - \nu_2(\xi - \mu)\beta_{01}^2$, $\kappa_1 = \mu[\tau(\alpha_{00} - \alpha_{01}) - \nu_2(1 - \beta_{01})] + \nu_2[2(\xi - \mu)\beta_{01}^2 - 2\mu^2\beta_{01}]$, and $\kappa_2 = \nu_2[2\mu^2\beta_{01} - (\xi - \mu)(1 + \beta_{01}^2)]$. Then, we can show that

$$R_{TPSP}(p_0) = \kappa_2 p_0^2 + \kappa_1 p_0 + \kappa_0. \quad (13)$$

Let p_0^* be the maximizer of (13) and \mathbf{p}^* be a $1 \times (T + 1)$ vector defined as $\mathbf{p}^* = [p_0^*, 1 - p_0^*, 0, 0, \dots]$. Then, index $I_j = I_j(y_j, \mathbf{x}, \mathbf{p}^*)$ for Imp-OTPSP can be shown to be

$$I_j = \begin{cases} \tau\alpha_{00} - \nu_2(1 + 2y_0 + 2\sum_{i=1}^T x_{i0}\alpha_{i0}) & \text{if } j = 0, \\ \tau\alpha_{01} - \nu_2\beta_{01}[1 + 2\beta_{01}y_1 + 2(\sum_{i=1}^{T-1} x_{i1}\alpha_{i1} + \mu p_0^*)] & \text{if } j = 1, \\ \tau\alpha_{0j} - \nu_2\beta_{0j}[1 + 2\beta_{0j}y_j + 2(\sum_{i=1}^{T-j} x_{ij}\alpha_{ij} + \mu((1 - p_0^*)\beta_{01} + p_0^*))] & \text{if } j = 2, 3, \dots, T. \end{cases} \quad (14)$$

Note that I_j decreases with y_j and x_{ij} . Thus under Imp-OTPSP, the more patients scheduled in one day, the smaller the chances that an additional patient will be scheduled on the same day, as we would expect for an intuitively “good” policy. This is markedly different from the linear-cost example since when costs are linear the marginal cost of an appointment is fixed

while in the quadratic-cost case, the marginal cost of an additional appointment for a given day increases with the total number of appointments already scheduled for that day.

Similarly, we can show that the index I_j for Imp-OAP can be shown to be

$$I_j = \begin{cases} \tau\alpha_{00} - \nu_2\beta_{00}(1 + 2\beta_{00}y_0 + 2\sum_{i=1}^T x_{i0}\alpha_{i0}) & \text{if } j = 0, \\ \tau\alpha_{0j} - \nu_2\beta_{0j}[1 + 2\beta_{0j}y_j + 2(\sum_{i=1}^{T-j} x_{ij}\alpha_{ij} + \mu)] & \text{if } j = 1, 2, \dots, T, \end{cases} \quad (15)$$

which is also decreasing in y_j and x_{ij} and thus carries similar characteristics as the index for Imp-OTPSP.

7. Estimating Model Parameters

In this section, we discuss how we estimated the parameters for the model clinic - which we use for our simulation study described in the next section - using data from an actual clinic. The data were obtained from the Department of Family Medicine of the School of Medicine at the University of North Carolina at Chapel Hill, and in consultation with Professor Samuel Weir, M.D. of the same organization. More specifically, the data came from the outpatient clinic of the Department of Family Medicine and consisted of logs from 7/1/2005 to 5/31/2007.

7.1 Cancellation and No-show Distributions

Implementation of the heuristics we propose requires the estimation of cancellation and no-show probabilities (α_{ij} and β_{ij} for integer i and j such that $i + j \leq T$). According to the data, for more than 99% of the patients, the appointment delay was less than 90 days. Therefore, estimates are determined for α_{ij} and β_{ij} for which $i + j \leq 90$. We determined the Maximum Likelihood Estimators (MLEs) for these probabilities. The details of this analysis are provided in Appendix B of the Online Companion.

Figure 1 compares the empirical results reported in Gallucci et al. (2005) with those obtained numerically from our statistical model. Gallucci et al. (2005) study a sample of around 6000 patients who were appointed to a psychiatry outpatient program, and estimate that the total rate of no-show and cancellation is 12%, 23%, 42%, and 44% corresponding to 0, 1, 7, and 13 days of appointment delay respectively. With our data and using our model, we found the same rates to be 17.99%, 18.48%, 21.37%, and 24.15%, respectively. Not surprisingly, the numbers we found are different from those found by Gallucci et al.

Nevertheless, they both point to the same relationship: the longer the appointment delay, the higher the chances of a no-show or cancellation.

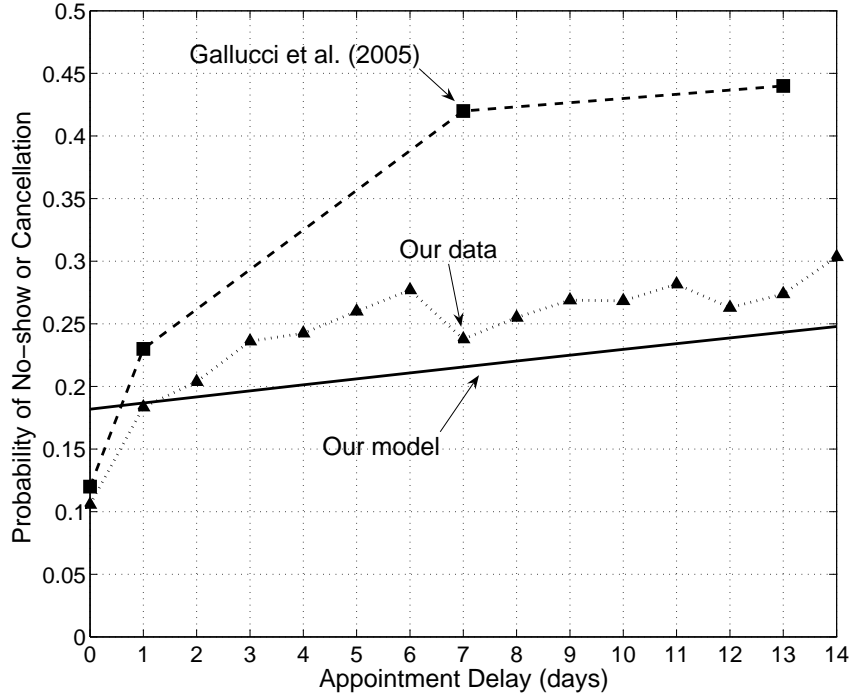


Figure 1: Probability of No-show or Cancellation vs. Appointment Delay

7.2 Daily Requests for Appointments

We estimated the average number of daily appointment requests to be approximately 50. However, there are no data that would make it possible to estimate its probability distribution. Therefore, in the simulation study, we assumed this distribution to be Poisson following other work in the literature, e.g. Patrick et al. (2008), Gupta and Wang (2008) and Robinson and Chen (2008).

7.3 Cost and Reward Functions

The reward function $r(\cdot, \cdot)$ for a clinic can be estimated in dollar terms given the relevant data. In most cases, however, the function will be used to reflect the preferences of a clinic regarding various trade-offs that are in play, which will no doubt be highly influenced by financial concerns. While profitability would be a major concern for many private clinics, this

might be less of an issue at university hospitals because their missions mainly lie elsewhere: research, education, and public service. For such hospitals, it is difficult to quantify the cost of overtime for physicians since it depends on time spent away from other important activities such as research and teaching. In fact, regardless of the type of the clinic even if the only concern was maximizing profits, not everything could be easily quantified. For example, what is the cost of the reduced quality of service given to the patients if these patients are seen on a day when the clinic exceeded its regular capacity by 15%? On an overloaded day patients will experience longer waits and they will be seen by physicians who are under pressure to clear up the high load. Therefore, the reward function should mostly be seen as user inputs that will differ depending on the circumstances different clinics are operating in as well as their preferences.

In our simulation study, we used the reward function that we determined in consultation with Samuel Weir, Professor and Co-Director of the Family Medicine Center at the University of North Carolina. Although values of specific parameters might be different for different clinics, we believe that the structure of these functions capture the basic trade-off that will be faced by many clinics. Specifically, we assumed that the reward function is in the form of (3) with $\eta(x) = x$ (implying one nominal reward for each patient served) and

$$w(z) = \begin{cases} K + h_1 z, & z \leq M, \\ K + h_1 M + h_2(z - M), & z > M, \end{cases} \quad (16)$$

In (16), $K \geq 0$ can be seen as the daily fixed cost, $M \geq 0$ as the regular daily capacity of the clinic, $h_1 \geq 0$ as the regular time cost of one scheduled patient, $h_2 \geq h_1$ as the overtime cost of one scheduled patient. This cost formulation makes it possible for the clinic to make its preference regarding daily patient overloads. If h_2 is set arbitrarily large, the policy will not schedule more than M patients on a single day unless all days are full. On the other hand, if h_1 and h_2 are set equal to each other, that implies the clinic does not mind overloading the clinic. Most clinics will prefer being somewhere in between, which they can determine by setting h_1 and h_2 accordingly.

8. Performance Comparison of the Heuristics: A Simulation Study

This section summarizes the results of the simulation study we carried out in order to investigate the performances of the scheduling policies proposed. As we have discussed in

the previous section, the model clinic we used in our simulation study was created using data from the Department of Family Medicine at UNC. We first derive OTPSP, and indices for Imp-OTPSP, and Imp-OAP for the “model clinic,” which has a linear reward function with one nominal reward for each patient served and a cost function given by (16). Then, we report and discuss the findings of the simulation study.

8.1 Derivation of the Heuristics for the Simulation Study

Suppose that λ denotes the mean number of daily appointment requests. Note that for our “model clinic,” λ is estimated to be 50. Let $\tilde{\lambda}_1(p_0) = \lambda p_0(\alpha_{00} - \alpha_{01}) + \lambda \alpha_{01}$ and $\tilde{\lambda}_2(p_0) = \lambda p_0(1 - \beta_{01}) + \lambda \beta_{01}$. Then, we can show that

$$R_{TPSP}(p_0) = \tilde{\lambda}_1(p_0) - \left[K + h_1 \tilde{\lambda}_2(p_0) + (h_2 - h_1) \left(\sum_{i=M}^{\infty} (i - M) e^{-\tilde{\lambda}_2(p_0)} \frac{[\tilde{\lambda}_2(p_0)]^i}{i!} \right) \right]. \quad (17)$$

Recall that OTPSP is obtained by finding p_0 that maximizes (17). Let p_0^* be the maximizer and $\mathbf{p}^* = [p_0^*, 1 - p_0^*, 0, 0, \dots, 0]$, which is a $(T + 1)$ -dimensional vector. Then, OTPSP is the policy of scheduling appointments independently of the system state using probability vector \mathbf{p}^* .

Now, let $\Gamma_j \stackrel{d}{=} \sum_{i=1}^{T-j} V_{ij}(\mathbf{x}) + \bar{V}_j(\mathbf{y}) + \sum_{k=\max\{j-1,1\}}^j \hat{V}_{kj}(\mathbf{p}^*)$ where $\sum_{k=\max\{j-1,1\}}^j \hat{V}_{kj}(\mathbf{p}^*) = 0$ if $j = 0$. Then, using (6), (8), and (16) we can show that the indices for Imp-OTPSP are given by

$$I_j = \alpha_{0j} - \beta_{0j} \{ \mathbf{E}[w(\Gamma_j + 1)] - \mathbf{E}[w(\Gamma_j)] \} = \alpha_{0j} - \beta_{0j} [h_1 + \mathbf{P}(\Gamma_j \geq M)(h_2 - h_1)]. \quad (18)$$

The indices for Imp-OAP are also given by (18) except that for Imp-OAP, $\mathbf{p} = [1, 0, 0, \dots]$, $\Gamma_j \stackrel{d}{=} \sum_{i=1}^{T-j} V_{ij}(\mathbf{x}) + \bar{V}_j(\mathbf{y}) + \hat{V}_{jj}(\mathbf{p})$, and $\hat{V}_{00}(\mathbf{p}) = 0$.

8.2 Findings of the Simulation Study

The objective of the simulation study was to compare the performances of the heuristic policies we propose (OTPSP, Imp-OAP, and Imp-OTPSP) with those of the benchmark heuristics (OAP, TP, BSP, RSP) introduced in Section 5. For TP, we set the threshold to be the regular daily capacity of the clinic, which is denoted by M as defined in Section 5.2.

Recall that we picked $\eta(x) = x$ and $w(\cdot)$ as in (16). As for the values of the parameters, we set $K = 0$ without loss of generality and $h_2 = 0.95$. We simulated various scenarios by considering different combinations of values for M and h_1 . Specifically, M took values

in $\{40, 45, 50, 55\}$ and h_1 took values in $\{0.0, 0.2, 0.5\}$. Since daily arrivals were set to 50, picking different values for M allowed us to test the performances of the policies under various conditions of system load. Clinics in scenarios where $M = 55$ can be seen as underloaded while those in scenarios where $M = 40$ or $M = 45$ are overloaded. On the other hand, assigning different values for h_1 helps us capture different preferences that a clinic might have for admitting more patients than the regular daily capacity. Our preliminary analysis indicated that, the policies we propose very rarely scheduled appointments more than 15 days in advance, and the times it took the simulation runs to complete quickly increased with the value of T , the maximum value for the appointment delay. Therefore, we set $T = 15$ in all scenarios.

We used the batch-means method (see, e.g. Section 9.5 in Law and Kelton (2000)). For each scenario, we ran 11 batches, each batch consisting of 200 consecutive workdays. The first batch was used as the warm-up period. A total of 12 different scenarios were considered each with a different pair of values for M and h_1 . Each scenario was simulated under each one of the seven scheduling policies, i.e. OAP, Imp-OAP, OTPSP, Imp-OTSP, TP, BSP and RSP, and the long-run average net reward was computed. For all scenarios considered, OTPSP turns out to be a deterministic static policy which always assigns the arriving appointment request to the next day, i.e. $p_0 = 0$ and $p_1 = 1$. It is difficult to give a clear description of the dynamic heuristics Imp-OAP and Imp-OTSP, but our numerical observations suggest that these policies generally resemble BSP, but unlike BSP, which schedules patients to the day with the fewest scheduled appointments, these two heuristics “balance” the load according to the indices. In order to facilitate comparison, we chose OAP (Open Access Policy) as the benchmark policy and for every other policy computed the percentage improvement that would be obtained by using the policy as opposed to OAP. Finally, we determined the 95% confidence interval for the mean percentage improvement. The results are given in Table 1.

In Table 1, the first number for each scenario-policy pair is the mean percentage improvement while the second number is the half width of the 95% confidence interval for the mean. Therefore, if the first number is larger than the second number, that indicates the corresponding policy is superior than OAP at the 5% significance level. On the other hand, if the first number is negative and it is larger than the second number in absolute value, that implies the superiority of OAP. The cases where the numbers indicate superiority of one policy over the other (in either direction) at the 5% significance level are shown in bold face. Note that the comparison is inconclusive in only two cases.

		Imp-OTPSP	OTPSP	Imp-OAP
$M = 55$	$h_1 = 0$	2.11% \pm 0.46%	0.78% \pm 0.32%	2.18% \pm 0.49%
	$h_1 = 0.2$	4.10% \pm 0.95%	3.23% \pm 0.75%	3.08% \pm 0.63%
	$h_1 = 0.5$	12.74% \pm 1.05%	12.14% \pm 1.10%	3.72% \pm 1.34%
$M = 50$	$h_1 = 0$	6.77% \pm 0.76%	2.75% \pm 0.37%	5.42% \pm 0.70%
	$h_1 = 0.2$	8.28% \pm 0.97%	5.48% \pm 0.86%	6.96% \pm 0.41%
	$h_1 = 0.5$	18.56% \pm 1.30%	15.29% \pm 1.31%	9.25% \pm 1.26%
$M = 45$	$h_1 = 0$	10.63% \pm 0.52%	6.23% \pm 0.60%	9.25% \pm 0.51%
	$h_1 = 0.2$	13.35% \pm 0.77%	9.13% \pm 0.76%	11.53% \pm 0.72%
	$h_1 = 0.5$	25.01% \pm 2.10%	20.32% \pm 1.41%	21.78% \pm 1.57%
$M = 40$	$h_1 = 0$	9.84% \pm 0.67%	9.12% \pm 0.44%	10.21% \pm 0.39%
	$h_1 = 0.2$	13.03% \pm 0.66%	12.48% \pm 0.68%	13.69% \pm 0.90%
	$h_1 = 0.5$	27.41% \pm 1.87%	26.82% \pm 1.49%	28.13% \pm 1.59%
		TP	BSP	RSP
$M = 55$	$h_1 = 0$	2.11% \pm 0.46%	-6.30% \pm 0.53%	-3.28% \pm 0.41%
	$h_1 = 0.2$	3.25% \pm 0.61%	-5.48% \pm 0.72%	-1.53% \pm 0.49%
	$h_1 = 0.5$	4.39% \pm 1.08%	-4.53% \pm 1.21%	2.68% \pm 1.02%
$M = 50$	$h_1 = 0$	6.45% \pm 0.73%	-2.22% \pm 0.73%	-1.20% \pm 0.51%
	$h_1 = 0.2$	8.21% \pm 0.92%	-1.09% \pm 0.89%	0.50% \pm 0.66%
	$h_1 = 0.5$	12.11% \pm 1.68%	0.72% \pm 1.47%	5.31% \pm 1.28%
$M = 45$	$h_1 = 0$	5.24% \pm 0.80%	4.11% \pm 0.54%	1.81% \pm 0.70%
	$h_1 = 0.2$	6.28% \pm 0.10%	4.91% \pm 0.69%	3.28% \pm 0.88%
	$h_1 = 0.5$	10.40% \pm 1.97%	8.10% \pm 1.38%	9.16% \pm 1.65%
$M = 40$	$h_1 = 0$	2.79% \pm 0.70%	2.99% \pm 0.55%	4.23% \pm 0.73%
	$h_1 = 0.2$	3.57% \pm 0.97%	3.83% \pm 0.75%	6.05% \pm 0.95%
	$h_1 = 0.5$	6.79% \pm 2.12%	7.32% \pm 1.62%	13.58% \pm 1.96%

Table 1: Results of the Simulation Study - (The first number indicates the difference between the mean performances of the corresponding policy and OAP, and the second number indicates the half width for the 95% confidence interval.)

A quick look at Table 1 reveals that Imp-OTPSP, OTPSP, Imp-OAP, and TP are the “best” policies under a variety of conditions. Although the Open Access Policy (OAP) does not perform as well as these policies, it does perform better as the regular capacity M gets larger. This is not surprising since as we also discussed in Section 1, OAP is an ideal policy when the system is not overloaded. It would thus be reasonable to expect that OAP would be among the best policies if M were larger.

In order to better compare the four best policies Imp-OTPSP, OTPSP, Imp-OAP, and TP among each other we also determined the best policies for each scenario separately, which are listed in Table 2. More specifically, for each scenario we conducted paired t-tests for every possible pair of policies and determined whether or not there is a statistical difference

between their performances at a significance level of 0.05. For every scenario, we listed the policies whose performances are better than those of the others. Note that for some of the scenarios, there is more than one policy. That is because in those cases, paired t-tests are inconclusive meaning that the performances of the policies are not statistically different.

Scenarios		Best Policies		
$M = 55$	$h_1 = 0$	Imp-OTPSP	Imp-OAP	TP
	$h_1 = 0.2$	Imp-OTPSP	Imp-OAP	TP
	$h_1 = 0.5$	Imp-OTPSP	OTPSP	
$M = 50$	$h_1 = 0$	Imp-OTPSP		
	$h_1 = 0.2$	Imp-OTPSP		TP
	$h_1 = 0.5$	Imp-OTPSP		
$M = 45$	$h_1 = 0$	Imp-OTPSP		
	$h_1 = 0.2$	Imp-OTPSP		
	$h_1 = 0.5$	Imp-OTPSP		
$M = 40$	$h_1 = 0$	Imp-OTPSP	Imp-OAP	
	$h_1 = 0.2$	Imp-OTPSP	Imp-OAP	
	$h_1 = 0.5$	Imp-OTPSP	Imp-OAP	

Table 2: “Best” Policies for Each Scenario at the 5% Significance Level

Table 2 clearly shows that the policies we propose particularly Imp-OTPSP perform well. The superiority of Imp-OTPSP over OTPSP is actually guaranteed given that Imp-OTPSP is obtained by applying a policy improvement step over OTPSP. It is also not surprising (though not guaranteed) that overall Imp-OTPSP performs better than Imp-OAP since the static policy that Imp-OTPSP improves upon (OTPSP) is superior than the static policy Imp-OAP improves upon (OAP).

TP performs quite well when the regular capacity M is large, but as it can be observed from Table 1, when M is small it is no better than BSP or even RSP, which schedules appointments randomly. On the other hand, all three policies we propose perform consistently well across all scenarios, suggesting that they are more robust than TP. Note that, one can in fact show that TP is an optimal policy if the clinic ignores patient no-shows and cancellations. Therefore, the poor performance of TP when the regular capacity is small also shows how damaging such omissions in formulation can be.

The superiority of Imp-OTPSP is more pronounced for mid values of the regular capacity ($M = 45$ or $M = 50$). This might be a consequence of the fact that there is no room for Imp-OTPSP to make a difference when the capacity is large or small. When the regular capacity is large, most policies manage to keep scheduled appointments below daily capacity

without delaying appointments too long. As a result most policies perform reasonably well and there is not much to gain from more sophisticated policies. On the other hand, when the regular capacity is small, under most policies the regular capacity is exceeded, which again limits the opportunity for sophisticated policies to make a difference. For mid values of the regular capacity, there are more choices that an intelligent scheduling policy can make because unlike the low or high capacity cases, some days will be overloaded and some days will be underloaded. An intelligent policy can work to smooth the daily loads so that overloads and no-shows are avoided to the extent possible.

So, what are the implications of our findings for practice? One should seek simplicity if complexity does not bring any significant advantages. Therefore, it appears that as suggested by Green and Savin (2008), if demand and supply are in balance, Open Access is quite reasonable especially if this policy is believed to have additional benefits that were not quantified in our formulation. However, if demand and supply are not in balance, if average demand is close to regular capacity or somewhat higher, then the policies we propose appear to be much better than the Open Access policy or other alternatives.

9. Concluding Remarks

In this paper we developed a model for the dynamic appointment scheduling decisions for a clinic that explicitly takes into account the possibility that patients may cancel or not show up for their appointments with probabilities that increase with their appointment delays. The model parameters can easily be estimated from data typically available for most clinics; in particular, the model does not assume the existence of data that are impossible or very difficult to obtain in practice, and the model allows us to develop dynamic heuristic policies that are easy to implement.

We tested the performances of the heuristic policies in a simulation study where we considered a model clinic generated by data provided to us by the Department of Family Medicine at the University of North Carolina. The results of the study clearly indicate that the heuristic policies we propose perform quite well particularly when the patient load is high. The policies outperform the Open Access policy even when the load is low, but the differences in the performances are not as significant and thus the Open Access policy might be a more reasonable choice if one wishes to keep scheduling simple. However, for high patient loads, the heuristics we propose appear to be much better than the Open Access

policy as well as other benchmark policies.

It is important to note that in this paper we used the term “Open Access” (OA) strictly to refer to the policy of scheduling all the appointments for the day when the requests are made. In practice, OA is sometimes interpreted or implemented in a more flexible way. For example, some clinics that use OA allow a certain portion of their patients to make appointments for the future. Some clinics use open access mostly as a guiding principle rather than a specific prescription of how appointments should be scheduled. On the other hand, some others implement OA but maintain more flexibility by scheduling appointments over a period of 2-5 days rather than a single day. Such flexible implementations of OA fit into our framework. The heuristics we propose can actually be seen as not necessarily alternatives to open access in the broader meaning of the term, but also as policies that can be used by clinics that implement some of the flexible versions of OA.

There are many ways to extend our model. First, our heuristics specify on which day to schedule an incoming appointment request, but not the time of the appointment on that day. It will be interesting to develop a sequential decision model to jointly decide on which day and at what time to schedule appointments. One possibility is to merge the existing slot scheduling algorithms already developed in the literature with our heuristic methods. More precisely, one can use the methods we propose to determine the day of the appointment and use a slot scheduling algorithm (such as the one proposed by Muthuraman and Lawley (2008)) to determine the time of the appointment. Another interesting extension would be to consider the heterogeneity of no-show behavior among the patients. One way of capturing this heterogeneity would be - as in the model of Muthuraman and Lawley (2008) - to group patients into different classes according to characteristics that correlate with their no-show behavior, and assign them appointments accordingly.

In this paper, we assumed that patients accept the first appointment date offered to them, which can be a reasonable assumption in many cases. However, in practice not all patients will be satisfied with the day offered to them and some will ask for another day. In fact, the heuristics we propose can be used even when the patients have preferences since they do not simply tell us on which day the next patient should be scheduled for. They also provide a ranking of alternative appointment dates since the index that the heuristics calculate for each day can be seen as a heuristic score of how much the clinic should rather assign that day instead of the others. Thus, if a patient is not happy with the day offered, the scheduler can offer the day with the next highest score. However, the performance of such a policy

should be investigated carefully since the heuristics are in fact derived under the assumption that patients accept the first appointment time offered to them. One interesting avenue for future work is to explicitly incorporate patient preferences into an appointment scheduling framework, which appears to be a significant modeling challenge. One possibility is to use a probabilistic discrete choice model as in Gupta and Wang (2008).

Finally, it is important to note that although this paper has been motivated by the scheduling of outpatient appointments for clinics, our framework and the policies we propose are in fact relevant to any system that schedules jobs in a similar fashion and experiences customer no-shows and cancellations (e.g., hair salons, mechanics, computer support services) that depend on appointment delays.

Acknowledgment

We are grateful to Sam Weir, M.D., director of the University of North Carolina Family Medicine Center, for his help in building cost and reward functions used in the simulation study, and providing us with the data which we used to estimate cancellation and no-show probabilities. We also thank Mr. Ratcliffe for his help in collecting and organizing the data. This research is partially supported by the National Science Foundation Grant #0620736.

References

- Argon, N. T., L. Ding, K. D. Glazebrook, S. Ziya. 2009. Dynamic routing of customers with general delay costs in a multiserver queuing system. *Probability in the Engineering and Informational Sciences* **23**(2) 175–203.
- Arvantes, J. 2007. Survey confirms growing demand for primary care physicians. *American Academy of Family Physicians News Now*. Oct 16, 2007.
- Belardi, F. G., S. Weir, F. W. Craig. 2004. A controlled trial of an advanced access appointment system in a residency family medicine center. *Fam. Med.* **36**(5) 341–345.
- Blumenthal, D. 2004. New steam from an old cauldron — The physician-supply debate. *N. Engl. J. Med.* **350**(17) 1780–1787.
- Cauchon, D. 2005. Medical miscalculation creates doctor shortage. *USA Today*. Mar 2, 2005.
- Cayirli, T., E. Veral. 2003. Outpatient scheduling in health care: A review of literature. *Production and Operations Management* **12**(4) 519–549.

- Denton, B., D. Gupta. 2003. A sequential bounding approach for optimal appointment scheduling. *IIE Transactions* **35**(11) 1003–1016.
- Dixon, S., F. C. Sampson, A. O’Cathain, M. Pickin. 2006. Advanced access: More than just GP waiting times? *Fam. Pract.* **23**(2) 233–239.
- Gallucci, G., W. Swartz, F. Hackerman. 2005. Brief reports: Impact of the wait for an initial appointment on the rate of kept appointments at a mental health center. *Psychiatr. Serv.* **56**(3) 344–346.
- Green, L. V., S. Savin. 2008. Reducing delays for medical appointments: A queueing approach. *Operations Research* **56**(6) 1526–1538.
- Gupta, D., B. Denton. 2008. Appointment scheduling in health care: Challenges and opportunities. *IIE Transactions* **40**(9) 800–819.
- Gupta, D., L. Wang. 2008. Revenue management for a primary-care clinic in the presence of patient choice. *Operations Research* **56**(3) 576–592.
- Ibaraki, T., N. Katoh. 1988. *Resource Allocation Problems : Algorithmic Approaches*. M.I.T. press, Cambridge, MA.
- Klassen, K. J., T. R. Rohleder. 2004. Outpatient appointment scheduling with urgent clients in a dynamic, multi-period environment. *International Journal of Service Industry Management* **15**(2) 167–186.
- Kopach, R., P. C. DeLaurentis, M. Lawley, K. Muthuraman, L. Ozsen, R. Rardin, H. Wan, P. Intrevado, X. Qu, D. Willis. 2007. Effects of clinical characteristics on successful open access scheduling. *Health Care Management Science* **10**(2) 111–124.
- Krishnan, K. R. 1990. Joining the right queue: A state-dependent decision rule. *IEEE Transactions on Automatic Control* **35**(1) 104–108.
- LaGanga, L. R., S. R. Lawrence. 2007. Clinic overbooking to improve patient access and increase provider productivity. *Decision Sciences* **38**(2) 251–276.
- Lamb, A. 2002. Why advanced access is a retrograde step. *Br. J. Gen. Pract.* **52**(485) 1035.
- Law, A. M., W. D. Kelton. 2000. *Simulation Modeling and Analysis*. New York: McGraw-Hill.
- Massachusetts Medical Society. 2007. MMS physician workforce study. Retrieved on May 14, 2009, <http://www.massmed.org>.

- Moore, C. G., P. Wilson-Witherspoon, J. C. Probst. 2001. Time and money: Effects of no-shows at a family practice residency clinic. *Fam. Med.* **33**(7) 522–527.
- Murray, M., T. Bodenheimer, D. Rittenhouse, K. Grumbach. 2003. Improving timely access to primary care: Case studies of the advanced access model. *J. of the American Medical Association* **289**(8) 1042–1046.
- Murray, M., C. Tantau. 2000. Same-day appointments: Exploding the access paradigm. *Family Practice Management*. September, 2000.
- Muthuraman, K., M. Lawley. 2008. A stochastic overbooking model for outpatient clinical scheduling with no-shows. *IIE Transactions* **40**(9) 820–837.
- Opp, M., K. Glazebrook, V. G. Kulkarni. 2005. Outsourcing warranty repairs: Dynamic allocation. *Naval Research Logistics* **52**(5) 381–398.
- Oppenheim, G. L., J. J. Bergman, E. C. English. 1979. Failed appointments: A review. *J. Fam. Pract.* **8**(4) 789–96.
- Patrick, J., M. L. Puterman, M. Queyranne. 2008. Dynamic multi-priority patient scheduling for a diagnostic resource. *Operations Research* **56**(6) 1507–1525.
- Pesata, V., G. Pallija, A. A. Webb. 1999. A descriptive study of missed appointments: Families’ perceptions of barriers to care. *J. Pediatr. Health Care* **13**(4) 178–182.
- Qu, X., R. L. Rardin, J. A. S. Williams, D. R. Willis. 2007. Matching daily healthcare provider capacity to demand in advanced access scheduling systems. *European Journal of Operational Research* **183**(2) 812–826.
- Robinson, L. W., R. R. Chen. 2003. Scheduling doctors’ appointments: Optimal and empirically-based heuristic policies. *IIE Transactions* **35**(3) 295–307.
- Robinson, L. W., R. R. Chen. 2008. The effects of patient no-shows on appointment scheduling policies. Working Paper, Graduate School of Management, University of California, Davis, CA, USA.
- Sack, K. 2008. In Massachusetts, universal coverage strains care. *The New York Times*. April 5, 2008.
- Shaked, M., J. G. Shanthikumar. 1994. *Stochastic Orders and Their Applications*. Academic Press, San Diego, CA.
- Solberg, L. I., M. C. Hroschikoski, J. M. Sperl-Hillen, P. J. O’Connor, B. F. Crabtree. 2004.

- Key issues in transforming healthcare organizations for quality: The case of advanced access. *Jt. Comm. J. Qual. Saf.* **30**(1) 15–24.
- Tijms, H. C. 1994. *Stochastic Models : An Algorithmic Approach*. John Wiley & Sons, Chichester, England.
- U.S. Census Bureau. 2008. The 2008 statistical abstract. Retrieved on May 14, 2009, <http://www.census.gov/prod/2007pubs/08abstract/health.pdf>.
- York, M. 2007. Few young doctors step in as upstate population ages. *The New York Times*. July 23, 2007.
- Zenios, S. A., G. M. Chertow, L. M. Wein. 2000. Dynamic allocation of kidneys to candidates on the transplant waiting list. *Operations Research* **48**(4) 549–569.

Online Supplement

Appendix A - Proofs of the Results

Proof of Proposition 1: To show that $f_j(y_j, \mathbf{x}, \mathbf{p})$ is concave in y_j , it suffices to show that $\mathbf{E}[r(w + \bar{W}_j(\mathbf{y}), v + \bar{V}_j(\mathbf{y}))]$ is concave in y_j for any real numbers w and v . Now, notice that $\bar{W}_j(\mathbf{y})$ and $\bar{V}_j(\mathbf{y}) = i$ are dependent. Furthermore, conditioning on $\bar{V}_j(\mathbf{y}) = i$, $\bar{W}_j(\mathbf{y})$ is a binomial random variable with parameters i and α_{0j}/β_{0j} . A direct application of Lemma 1 yields the result. ■.

Lemma 1. Define $B(i, p)$ to be a Binomial random variable with parameters i and p . Let $\{W_j, j = 0, 1, 2, \dots\}$ and $\{V_j, j = 0, 1, 2, \dots\}$ be two sequences of random variables such that (1) $W_j \stackrel{d}{=} B(j, \alpha)$, (2) $V_j \stackrel{d}{=} B(j, \beta)$ and (3) conditioning on $V_j = i$, $W_j \stackrel{d}{=} B(i, \alpha/\beta)$ for $0 \leq i \leq j$, where $0 \leq \alpha \leq \beta \leq 1$. Suppose that $r(x, z)$ is submodular and jointly concave in x and z . For any arbitrary real numbers w and v , let $g(j) = \mathbf{E}\{r(W_j + w, V_j + v)\}$. Then $g(j)$ is a concave function over the set of non-negative integers, i.e., $g(j+2) + g(j) \leq 2g(j+1)$, $\forall j = 0, 1, 2, \dots$.

Proof: The proof uses a coupling argument. Let $\{Z_{i,k}, k = 0, 1, 2, \dots\}$ be a sequence of i.i.d. Bernoulli random variables with parameter β for each $k \in \{1, 2, 3, 4\}$. First, we write

$$\begin{aligned} & g(j+2) + g(j) - 2g(j+1) \\ = & \mathbf{E}\left\{r\left(B\left(\sum_{k=1}^{j+2} Z_{1,k}, \alpha/\beta\right) + w, \sum_{k=1}^{j+2} Z_{1,k} + v\right) + r\left(B\left(\sum_{k=1}^j Z_{2,k}, \alpha/\beta\right) + w, \sum_{k=1}^j Z_{2,k} + v\right) \right. \\ & \left. - r\left(B\left(\sum_{k=1}^{j+1} Z_{3,k}, \alpha/\beta\right) + w, \sum_{k=1}^{j+1} Z_{3,k} + v\right) - r\left(B\left(\sum_{k=1}^{j+1} Z_{4,k}, \alpha/\beta\right) + w, \sum_{k=1}^{j+1} Z_{4,k} + v\right)\right\}. \quad (\text{A-2}) \end{aligned}$$

Now we couple the random variables so that $Z_{1,k} = Z_{2,k} = Z_{3,k} = Z_{4,k} = z_k$ for $k = 0, 1, 2, \dots, j$, $Z_{1,j+1} = Z_{3,j+1} = \hat{z}$ and $Z_{1,j+2} = Z_{4,j+1} = \tilde{z}$. There are four possible cases for the pair (\hat{z}, \tilde{z}) : (i) $\hat{z} = \tilde{z} = 0$, (ii) $\hat{z} = 1$ and $\tilde{z} = 0$, (iii) $\hat{z} = 0$ and $\tilde{z} = 1$ and (iv) $\hat{z} = \tilde{z} = 1$. For case (i), (ii) and (iii), (A-2) reduces to 0. In the rest of the proof, we show that (A-2) is also non-positive under case (iv).

Let $z = \sum_{k=1}^j z_k$. Then, under case (iv), we have

$$\begin{aligned} & g(j+2) + g(j) - 2g(j+1) \\ &= \mathbf{E}\{r(B(z+2, \alpha/\beta) + w, z+2+v) + r(B(z, \alpha/\beta) + w, z+v) \\ &\quad - r(B(z+1, \alpha/\beta) + w, z+1+v) - r(B(z+1, \alpha/\beta) + w, z+1+v)\}. \end{aligned} \quad (\text{A-3})$$

Following as above, we define $\{U_{i,k}, k = 0, 1, 2, \dots\}$ be a sequence of i.i.d. Bernoulli random variables with parameter α/β for each $k \in \{1, 2, 3, 4\}$. Then we write equation (A-3) as

$$\begin{aligned} & g(j+2) + g(j) - 2g(j+1) \\ &= \mathbf{E}\{r(\sum_{k=1}^{z+2} U_{1,k} + w, z+2+v) + r(\sum_{k=1}^z U_{2,k} + w, z+v) \\ &\quad - r(\sum_{k=1}^{z+1} U_{3,k} + w, z+1+v) - r(\sum_{k=1}^{z+1} U_{4,k} + w, z+1+v)\}. \end{aligned} \quad (\text{A-4})$$

Now we couple the random variables so that $U_{1,k} = U_{2,k} = U_{3,k} = U_{4,k} = u_k$ for $k = 0, 1, 2, \dots, z$, $U_{1,z+1} = U_{3,z+1} = \hat{u}$ and $U_{1,z+2} = U_{4,z+1} = \tilde{u}$. There are four possible cases for the pair (\hat{u}, \tilde{u}) :

Case 1: $\hat{u} = \tilde{u} = 0$. Then the term inside the expectation in (A-4) becomes

$$r(\sum_{k=1}^z u_k + w, z+2+v) + r(\sum_{k=1}^z u_k + w, z+v) - r(\sum_{k=1}^z u_k + w, z+1+v) - r(\sum_{k=1}^z u_k + w, z+1+v) \leq 0$$

due to the concavity of $r(\cdot, \cdot)$.

Case 2: $\hat{u} = 1$ and $\tilde{u} = 0$. Then the term inside the expectation in (A-4) becomes

$$r(\sum_{k=1}^z u_k + 1 + w, z+2+v) + r(\sum_{k=1}^z u_k + w, z+v) - r(\sum_{k=1}^z u_k + 1 + w, z+1+v) - r(\sum_{k=1}^z u_k + w, z+1+v) \leq 0$$

since

$$\begin{aligned} & r(\sum_{k=1}^z u_k + 1 + w, z+2+v) - r(\sum_{k=1}^z u_k + 1 + w, z+1+v) \\ & \leq r(\sum_{k=1}^z u_k + 1 + w, z+1+v) - r(\sum_{k=1}^z u_k + 1 + w, z+v) \quad (\text{due to the concavity of } r(\cdot)) \\ & \leq r(\sum_{k=1}^z u_k + w, z+1+v) - r(\sum_{k=1}^z u_k + w, z+v). \quad (\text{due to the submodularity of } r(\cdot, \cdot)) \end{aligned}$$

Case 3: $\hat{u} = 0$ and $\tilde{u} = 1$. Then the term inside the expectation in (A-4) becomes

$$r(\sum_{k=1}^z u_k + 1 + w, z+2+v) + r(\sum_{k=1}^z u_k + w, z+v) - r(\sum_{k=1}^z u_k + w, z+1+v) - r(\sum_{k=1}^z u_k + 1 + w, z+1+v) \leq 0,$$

following the same proof in Case (2) above.

Case 4: $\hat{u} = \tilde{u} = 1$. Then the term inside the expectation in (A-4) becomes

$$r\left(\sum_{k=1}^z u_k + 2 + w, z + 2 + v\right) + r\left(\sum_{k=1}^z u_k + w, z + v\right) - r\left(\sum_{k=1}^z u_k + 1 + w, z + 1 + v\right) - r\left(\sum_{k=1}^z u_k + 1 + w, z + 1 + v\right) \leq 0$$

due to the concavity of $r(\cdot, \cdot)$. Hence, the result follows. ■

Proof of Proposition 2: Let λ denote the mean number of daily arrivals and $PO(\alpha)$ denote a Poisson random variable with mean α . Then,

$$\tilde{V}_0(p_0) \stackrel{d}{=} B(\tilde{A}_0, p_0) \stackrel{d}{=} PO(\lambda p_0), \tilde{V}_1(p_0) \stackrel{d}{=} B(\tilde{A}_1, (1 - p_0)\beta_{01}) \stackrel{d}{=} PO(\lambda(1 - p_0)\beta_{01}),$$

and

$$\tilde{W}_0(p_0) \stackrel{d}{=} B(\tilde{A}_0, p_0\alpha_{00}) \stackrel{d}{=} PO(\lambda p_0\alpha_{00}), \tilde{W}_1(p_0) \stackrel{d}{=} B(\tilde{A}_1, (1 - p_0)\alpha_{01}) \stackrel{d}{=} PO(\lambda(1 - p_0)\alpha_{01}).$$

Note that $\tilde{V}_0(p_0)$ is independent of $\tilde{V}_1(p_0)$ and $\tilde{W}_0(p_0)$ is independent of $\tilde{W}_1(p_0)$ since \tilde{A}_0 is independent of \tilde{A}_1 . Let $C_1(p_0)$ be a random variable such that $C_1(p_0) \stackrel{d}{=} PO(\lambda p_0(\alpha_{00} - \alpha_{01}))$ (note that $\alpha_{00} - \alpha_{01} \geq 0$), and $D_1 \stackrel{d}{=} PO(\lambda\alpha_{01})$ be a random variable which is independent of $C_1(p_0)$. Also, let $C_2(p_0)$ be a random variable such that $C_2(p_0) \stackrel{d}{=} PO(\lambda p_0(1 - \beta_{01}))$, and $D_2 \stackrel{d}{=} PO(\lambda\beta_{01})$ be a random variable which is independent of $C_2(p_0)$. Then, using the fact that the sum of two independent Poisson random variables is another Poisson random variable, we have

$$\tilde{W}_0(p_0) + \tilde{W}_1(p_0) \stackrel{d}{=} PO(\lambda p_0(\alpha_{00} - \alpha_{01}) + \lambda\alpha_{01}) \stackrel{d}{=} C_1(p_0) + D_1$$

and

$$\tilde{V}_0(p_0) + \tilde{V}_1(p_0) \stackrel{d}{=} PO(\lambda p_0(1 - \beta_{01}) + \lambda\beta_{01}) \stackrel{d}{=} C_2(p_0) + D_2.$$

Assumption 1 implies that $w(\cdot)$ is an increasing convex function and $r(\cdot)$ is an increasing concave function. Let $z^+ = \max\{z, 0\}$. Since the class of functions $g_s(z) = (z - s)^+$ for all $s \in \mathbb{R}$ generates all the increasing convex functions and the class of functions $h_s(z) = -(s - z)^+$ for all $s \in \mathbb{R}$ generates all the increasing concave functions (see, e.g. page 173 in Shaked and Shanthikumar (1994)), in order to show

$$\mathbf{E} \left[r \left(\tilde{W}_0(p_0) + \tilde{W}_1(p_0) \right) - w \left(\tilde{V}_0(p_0) + \tilde{V}_1(p_0) \right) \right]$$

is concave in p_0 it suffices to show that

$$\mathbf{E}\{-[s - (C_1(p_0) + D_1)]^+\}$$

is concave in p_0 for all $s \in \mathbb{R}$ and

$$\mathbf{E}\{[(C_2(p_0) + D_2) - s]^+\}$$

is convex in p_0 for all $s \in \mathbb{R}$, both of which immediately follow from Lemma 2 shown and proved below. ■.

Lemma 2. Let $z^+ = \max\{z, 0\}$. Suppose that $\tilde{C}(p)$ is a Poisson random variable with mean $p\theta$ where $p, \theta > 0$ and A is a random variable that is independent of $\tilde{C}(p)$. Then $\mathbf{E}\{[(\tilde{C}(p) + A) - s]^+\}$ and $\mathbf{E}\{[s - (\tilde{C}(p) + A)]^+\}$ are both convex in p for any given $s \in \mathbb{R}$.

Proof: It is sufficient to show that for all $a \in \mathbb{R}$, $\mathbf{E}\{[(\tilde{C}(p) + A) - s]^+ | A = a\}$ and $\mathbf{E}\{[s - (\tilde{C}(p) + A)]^+ | A = a\}$ are both convex in p for an arbitrary $s \in \mathbb{R}$. We first show that $\mathbf{E}\{[(\tilde{C}(p) + A) - s]^+ | A = a\}$ is convex in p for any given $a, s \in \mathbb{R}$. This is equivalent to showing that $\mathbf{E}\{[\tilde{C}(p) - s]^+\}$ is convex in p for any given $s \in \mathbb{R}$.

Now, to show that $\mathbf{E}\{[\tilde{C}(p) - s]^+\}$ is convex in p , first note that the proof is trivial if $s \leq 0$, and thus we only need to consider $s > 0$. Denote $\lceil s \rceil$ to be the smallest integer that is larger than or equal to s , and let $g(\cdot | \lambda)$ be the probability mass function of a Poisson random variable with mean λ , i.e.

$$g(i | \lambda) = \begin{cases} e^{-\lambda} \frac{\lambda^i}{i!}, & \text{if } i = 0, 1, 2, \dots, \\ 0, & \text{otherwise.} \end{cases}$$

After some straightforward but tedious calculus, one can show that for any $s > 0$,

$$\frac{d^2}{dp^2} \mathbf{E}\{[\tilde{C}(p) - s]^+\} = \theta^2 [(s + 1 - \lceil s \rceil)g(\lceil s \rceil - 1 | \theta p) + (\lceil s \rceil - s)g(\lceil s \rceil - 2 | \theta p)] \geq 0,$$

which establishes that $\mathbf{E}\{[\tilde{C}(p) - s]^+\}$ is convex in p for any given $s \in \mathbb{R}$.

Similarly, to show $\mathbf{E}\{[s - (\tilde{C}(p) + A)]^+\}$ is convex in p for any given $s \in \mathbb{R}$, it is sufficient to show that $\mathbf{E}\{[s - \tilde{C}(p)]^+\}$ is convex in p for any given $s \in \mathbb{R}$. But then since $\mathbf{E}\{[s - \tilde{C}(p)]^+\} = \mathbf{E}\{[\tilde{C}(p) - s]^+ + [s - \tilde{C}(p)]\} = \mathbf{E}\{[\tilde{C}(p) - s]^+\} + s - \theta p$, the result immediately follows. ■.

Proof of Theorem 1: Let f^* denote the policy that schedules all appointment requests received today for the j^* th day from today where $j^* = \arg \max_{j \in \{0, 1, \dots, T\}} I_j$ where I_j is as described in (12).

First, the long-run average net reward under policy f^* is $\phi_{f^*} = (\tau\alpha_{0j^*} - \nu_1\beta_{0j^*})\mu$. For almost all sample paths w of $\{A^1, A^2, A^3, \dots\}$, we show that, for any arbitrary policy f , the long-run average net reward along this path, denoted by $\phi_f(w)$, is no larger than

ϕ_{f^*} . Without loss of generality, we assume that the initial backlog is empty, i.e. $\mathbf{X}^0 = \mathbf{0}$. Let $N_j(t, \omega)$ be the total number of patients scheduled with appointment delay j days up to day t along sample path ω under policy f . Now, since the cost and reward are both linear, we can view the total net reward as the sum of net rewards contributed by individual patients. Let $R_{ij}(t, \omega)$ be the net reward generated by the i th patient of those $N_j(t, \omega)$ patients whose appointment delay is j days along sample path ω up to day t , $i = 1, 2, \dots, N_j(t, \omega)$. Notice that $\{R_{ij}(t, \omega), i = 1, 2, \dots, N_j(t, \omega)\}$ is the realization of a sequence of i.i.d. random variables with mean $\tau\alpha_{0j} - \nu_1\beta_{0j}$ along sample path ω . Let $\mathcal{J} = \{j : \lim_{t \rightarrow \infty} N_j(t, \omega) = \infty\}$. Then,

$$\begin{aligned}
\phi_f(\omega) &= \lim_{t \rightarrow \infty} \frac{\sum_{j=0}^T \sum_{i=1}^{N_j(t, \omega)} R_{ij}(t, \omega)}{t} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j=0}^T N_j(t, \omega) \frac{\sum_{i=1}^{N_j(t, \omega)} R_{ij}(t, \omega)}{N_j(t, \omega)} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j \in \mathcal{J}} N_j(t, \omega) \frac{\sum_{i=1}^{N_j(t, \omega)} R_{ij}(t, \omega)}{N_j(t, \omega)} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j \in \mathcal{J}} N_j(t, \omega) \lim_{t \rightarrow \infty} \frac{\sum_{i=1}^{N_j(t, \omega)} R_{ij}(t, \omega)}{N_j(t, \omega)} \\
&= \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j \in \mathcal{J}} N_j(t, \omega) (\tau\alpha_{0j} - \nu_1\beta_{0j}) \quad (\text{by Strong Law of Large Numbers}) \\
&\leq (\tau\alpha_{0j^*} - \nu_1\beta_{0j^*}) \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j \in \mathcal{J}} N_j(t, \omega) \quad (\text{by the definition of } j^*) \\
&\leq (\tau\alpha_{0j^*} - \nu_1\beta_{0j^*}) \mu \quad (\text{since } \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j \in \mathcal{J}} N_j(t, \omega) \leq \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{j=0}^T N_j(t, \omega) = \mu) \\
&= \phi_{f^*}.
\end{aligned}$$

This completes the proof. ■.

Appendix B - Derivations of MLEs for α_{ij} and β_{ij}

From the data obtained from the Department of Family Medicine at the University of North Carolina, we extracted the following information:

C_i : Number of patients who had an appointment delay of i days but canceled their appointments on or before their appointment days.

S_i : Number of patients who had an appointment delay of i days and showed up for their appointments.

M_i : Number of patients who had an appointment delay of i days, did not cancel in advance but missed their appointments.

First, define the following:

$$\begin{aligned} q_i &= \mathbf{P}(\text{NS}|T_c \geq i + 1), \\ r_i &= \mathbf{P}(\text{S}|T_c \geq i + 1), \\ u_i &= \mathbf{P}(T_c \leq i). \end{aligned} \tag{A-5}$$

Note that q_i and r_i are respectively the probabilities of the “patient no-show” and “patient show” events given that the patient does not cancel in the first $i + 1$ days after the day she calls for an appointment; u_i is the probability that a patient will cancel no later than i days after she calls for an appointment.

Clearly, we must have $q_i + r_i = 1$ and $u_i \leq u_{i+1}$, $i \in \{0, 1, \dots, T\}$. Recall that each patient’s cancellation and no-show behaviors are independent of those of other patients and for any appointment made, there are three possible outcomes: the patient cancels any time on or before the appointment day, the patient misses the appointment without cancellation, and the patient shows up for the appointment. Let $\mathbf{q} = \{q_i\}_{i=0}^T$, $\mathbf{r} = \{r_i\}_{i=0}^T$, $\mathbf{u} = \{u_i\}_{i=0}^T$. Then the MLEs for q_i , r_i , and u_i can be obtained by solving the following maximization problem

$$\begin{aligned} \max \quad & L(\mathbf{q}, \mathbf{r}, \mathbf{u}) = \prod_{i=0}^T u_i^{C_i} [(1 - u_i)q_i]^{M_i} [(1 - u_i)r_i]^{S_i} \\ \text{s.t.} \quad & q_i + r_i = 1, \quad i \in \{0, 1, \dots, T\}, \\ & u_i \leq u_{i+1}, \quad i \in \{0, 1, \dots, T - 1\}, \\ & u_i \geq 0, q_i \geq 0, r_i \geq 0, \quad i \in \{0, 1, \dots, T\}. \end{aligned} \tag{A-6}$$

Suppose that we solve the optimization problem (A-6) and obtain the MLEs \hat{q}_i , \hat{r}_i , and \hat{u}_i for q_i , r_i , and u_i , respectively. Then, we can get the MLEs $\hat{\alpha}_{ij}$ and $\hat{\beta}_{ij}$ for α_{ij} and β_{ij} as follows:

$$\begin{aligned} \hat{\alpha}_{ij} &= \frac{\hat{r}_{i+j}(1-\hat{u}_{i+j})}{1-\hat{u}_{i-1}}, \\ \hat{\beta}_{ij} &= \frac{1-\hat{u}_{i+j-1}}{1-\hat{u}_{i-1}}. \end{aligned} \tag{A-7}$$

where $\hat{u}_{-1} = 0$ by definition.

However, solving Problem (A-6) is difficult especially since it has too many decision variables, i.e., there are too many parameters to estimate. Hence, we propose a parsimonious parametric model, which requires estimating only four parameters. Specifically, assume that

q_i , r_i , and u_i take the following form:

$$\begin{aligned} q_i &= 1 - \theta b^{i+1}, & i \geq -1, \\ r_i &= \theta b^{i+1}, & i \geq -1, \\ u_i &= \begin{cases} 0, & i = -1, \\ 1 - \gamma a^i, & i \geq 0. \end{cases} \end{aligned} \tag{A-8}$$

This model is appealing not only because it is sufficiently simple but also because it has an interpretation that is quite fitting in the appointment scheduling context. First, note that u_i is the cumulative distribution function for the random variable T_c , i.e. the time between the patient's call and the day she cancels (or the day she would cancel if her appointment was not earlier). Under the parametric form we describe in (A-8), the probability mass function of T_c is a mixture of two distributions, one being a constant and the other being a geometric distribution. To be more precise, we have

$$\mathbf{P}(T_c = i) = (1 - \gamma)\mathbf{1}_{\{i=0\}} + \gamma\mathbf{P}(Y_c = i)\mathbf{1}_{\{i \geq 1\}},$$

where $\mathbf{1}_A$ is the indicator function and Y_c is a geometric random variable with parameter $1 - a$. One way of interpreting this mixture structure is that there are two different types of patients: those who cancel on the same day they make their appointments, which constitute $1 - \gamma$ fraction of the whole patient population, and those who cancel later, which constitute γ fraction of the whole patient population. Furthermore, the model also implies that for those who make appointments at least one day before their appointment day, the probability of cancelling on each day is $1 - a$ independently of everything else. (Note that it is possible to make similar interpretations for q_i and r_i as well.)

Notice that once the probabilities are restricted to be in the form given above, the only additional condition needed for all the constraints of Problem (A-6) to hold is that $0 \leq \gamma, a, \theta, b \leq 1$. Let $\hat{\gamma}$, \hat{a} , $\hat{\theta}$, and \hat{b} denote the MLEs for γ, a, θ , and b , respectively. Then, the first order optimality condition yields

$$\begin{aligned} \sum_{i=0}^T (C_i + M_i + S_i) &= \sum_{i=0}^T \frac{C_i}{1 - \hat{\gamma}\hat{a}^i}, \\ \sum_{i=0}^T i(C_i + M_i + S_i) &= \sum_{i=0}^T \frac{iC_i}{1 - \hat{\gamma}\hat{a}^i}, \\ \sum_{i=0}^T (M_i + S_i) &= \sum_{i=0}^T \frac{M_i}{1 - \hat{\theta}\hat{b}^{i+1}}, \\ \sum_{i=0}^T (i+1)(M_i + S_i) &= \sum_{i=0}^T \frac{(i+1)M_i}{1 - \hat{\theta}\hat{b}^{i+1}}. \end{aligned}$$

This system of nonlinear equations can be solved by one of the standard algorithms such as Gauss-Newton method or Trust-Region method. Once this solution is found, the estimates for α_{ij} and β_{ij} can be determined using the following equations, which are obtained by simply substituting (A-8) into (A-7):

$$\hat{\alpha}_{ij} = \begin{cases} \hat{\theta}\hat{b}^{j+1}\hat{\gamma}\hat{a}^j & \text{if } i = 0, \\ \hat{\theta}\hat{b}^{i+j+1}\hat{a}^{j+1} & \text{if } i \geq 1, \end{cases}$$

$$\hat{\beta}_{ij} = \begin{cases} 1 & \text{if } i = 0, j = 0, \\ \hat{\gamma}\hat{a}^{j-1} & \text{if } i = 0, j \geq 1, \\ \hat{a}^j & \text{if } i \geq 1, j \geq 1. \end{cases}$$

Using our data, we found that $\hat{\gamma} = 0.9297$, $\hat{a} = 0.9987$, $\hat{\theta} = 0.8863$, and $\hat{b} = 0.9953$. Then, we numerically verified that this solution is indeed a maximizer. In order to build confidence for our statistical model, we have also conducted a Chi-square goodness-of-fit test, and we found that the distribution we proposed can not be rejected at the significance level of 0.01 (the p-value is 0.042).