

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

A Data-Driven Model of an Appointment-Generated Arrival Process at an Outpatient Clinic

Song-Hee Kim

Data Sciences and Operations, USC Marshall School of Business, songheek@marshall.usc.edu, <http://www-bcf.usc.edu/songheek>

Ward Whitt

Industrial Engineering and Operations Research, Columbia University, ww2040@columbia.edu, <http://www.columbia.edu/ww2040>

Won Chul Cha

Department of Emergency Medicine, Samsung Medical Center, docchaster@gmail.com

We develop a high-fidelity simulation model of the patient arrival process to an endocrinology clinic by carefully examining appointment and arrival data from that clinic. The data includes the time that the appointment was originally made as well as the time that the patient actually arrived, or if the patient did not arrive at all, in addition to the scheduled appointment time. We take a data-based approach, specifying the schedule for each day by its value at the end of the previous day. This data-based approach shows that the schedule for a given day evolves randomly over time. Indeed, in addition to three recognized sources of variability, namely, (i) no-shows, (ii) extra unscheduled arrivals and (iii) deviations in the actual arrival times from the scheduled times, we find that the primary source of variability in the arrival process is variability in the daily schedule itself. Even though service systems with arrivals by appointment can differ in many ways, we think that our data-based approach to modeling the clinic arrival process can be a guideline or template for constructing high-fidelity simulation models for other arrival processes generated by appointments.

Key words: simulation stochastic input modeling, simulating appointment-generated arrival processes, scheduled arrivals in service systems, outpatient clinics, data-driven modeling, stochastic models in healthcare

History: June 20, 2016

1. Introduction

In this paper we aim to contribute to simulation stochastic input modeling. In particular, we develop an approach for creating high-fidelity stochastic models of arrival processes generated by appointments. We do that so that the arrival process model can be part of a full simulation model used as to improve operations (e.g., to improve throughput, control individual workloads, set staffing levels and allocate other resources), with the goal of efficiently providing good service in a service system with arrivals by appointment.

We carefully examine data from the endocrinology outpatient clinic of the Samsung Medical Center in South Korea. The data were collected over a 13-week period from July 2013 to September 2013. Included in the data are the day and time of each appointment and when the appointment was made as well as the final disposition, i.e., whether or not the scheduled arrival actually came and, if so, what was the time of arrival, and if the arrival did not come, if there was a cancellation; if not, it is regarded as a no-show.

1.1. A Long History of Modeling and Analysis

There is a long history of modeling and analysis of outpatient clinics and other healthcare systems, with notable early work Bailey (1952), Welch and Bailey (1952), Fetter and Thompson (1965) and Swartzman (1970); surveys Cayirli and Veral (2003), Gupta and Denton (2008), Jacobson et al. (2006) and Jun et al. (1999); and edited reviews Hall (2006) and Hall (2012). The large literature can be divided roughly into three types of analyses, depending on their focus: (i) conducting a full analysis of an outpatient clinic to make operational improvements, (ii) designing an effective appointment system and (iii) conducting a performance analysis of queueing models based on assumed properties of clinic arrival processes.

As illustrated by the seminal paper by Fetter and Thompson (1965), it is recognized that outpatient clinics can be usefully represented as a complex network of queues associated with the reception area, nurses, labs and doctors. Patients often follow different paths through the clinic, depending on many factors, such as the doctor whom they are scheduled to see, their medical condition and the results of medical tests. Thus, to analyze and improve the performance of a clinic, it is important to construct careful process maps or work-flow diagrams, e.g., as in Figure 1 in Chand et al. (2009) and Figure 2 in Harper and Gamlin (2003). The system complexity has made simulation the dominant choice for detailed analysis of a clinic. Many successful simulation studies have been conducted, as can be seen from Chand et al. (2009), Chakraborty et al. (2010), Guo et al. (2004), Harper and Gamlin (2003), Swisher et al. (2001). There is also potential for analytical queueing network models such as in Whitt (1983), as discussed by Zonderland and Boucherie (2012).

Most outpatient clinics have a substantial portion of their arrivals scheduled in advance, i.e., generated by an appointment system. Thus, as expected, a large part of the literature is devoted to designing an effective appointment system, as can be seen from surveys Cayirli and Veral (2003) and Gupta and Denton (2008) and other recent works Liu et al. (2010), Luo et al. (2012) and Liu

and Ziya (2014). It is also recognized that appointment-generated arrival processes are very different from arrival processes where customers independently decide when to arrive. It is well known that the arrival process often tends to have a nearly periodic structure determined by appointment time slots but that it also can be significantly variable because of no-shows, unscheduled arrivals and earliness or lateness. Empirical studies of patient no-shows and non-punctual arrivals have been conducted, and they indicate that no-show rates vary across different services and patient populations: the reported no-show rates are as low as 4.2% at a general practice outpatient clinic in the United Kingdom (Neal et al. 2001) and as high as 31% at a family practice clinic (Moore et al. 2001). Thus, ever since the seminal papers of Bailey (Bailey 1952, Welch and Bailey 1952), work has been done to analyze queueing models that reflect key structural properties of appointment-generated arrival processes, e.g., see Feldman et al. (2014), Hassin and Mendel (2008), Jouini and Benjaafar (2012), Kaandorp and Koole (2007), Wang et al. (2010, 2014) and Zacharias and Pinedo (2014).

1.2. Probing Deeply into One Clinic Arrival Process

In this paper, we do not follow any of the three time-tested approaches discussed in Section 1.1. Instead, we devote this entire paper to carefully examining arrival data from the outpatient clinic appointment system. In doing so, we aim to construct a high-fidelity stochastic arrival process model that can be part of a simulation model or analytic queueing model that can be used to improve the performance of the clinic. We want to understand the consequence of existing appointment schedules; we do not consider alternative scheduling algorithms.

Sixteen doctors work in this clinic, but patients make an appointment to see a particular doctor, so each arriving patient knows which doctor he or she will meet. Hence, each doctor operates as a single-server system. Each doctor works within a subset of available days and shifts, with three shifts available: morning (am) shifts, roughly from 8:30 am to 12:30 pm; afternoon (pm) shifts, roughly from 12:30 pm to 4:30 pm; and full-day shifts. During the weekdays of the 13-week study period, the 16 doctors worked for a total of 228 am shifts, 220 pm shifts and 25 full-day shifts. The shifts were not evenly distributed among the doctors, with the number of shifts per doctor ranging from 11 to 46.

We have studied the data for all 16 doctors, but in this paper, we primarily focus on patient arrivals during the am shifts of one doctor. This doctor was selected among the 16 candidate doctors because of his relatively high volume of patients: he worked for a total of 22 am shifts (12 on Tuesdays and 10 on Fridays) and 22 pm shifts (11 on Mondays, 2 on Wednesdays and 9 on

Thursdays) during our study period. The results of the analysis of the other doctors are presented in our longer, more detailed study Kim et al. (2015a). We emphasize that analysis of the arrival process for each doctor is important because patients with appointments for different doctors tend to follow different paths through the clinic. The overall arrival process for all doctors is of course important for studying the congestion in the waiting room, but the total arrival process can be directly modeled as the superposition of the arrival processes for the individual doctors.

1.3. Randomness in the Appointment Schedule

For the clinic arrival process, even though we find the customary deviations from an ideal deterministic pattern of arrivals, including no-shows, extra unscheduled arrivals and non-punctuality, *we find that the main deviation from a regular deterministic arrival pattern is variability in the schedule itself.* However, we first need to define what we mean by “the schedule” because it can be defined in different ways.

Overall, we take a data-based approach. Instead of relying on the framework of the appointment system or management expressed intentions, we rely on the data, which includes the day and time when the appointment was made as well as when the arrival was scheduled and actually arrived. *By “the schedule,” we specifically mean both the number and the arrival times of all arrivals scheduled for a given day, as determined by the appointment system by the end of the previous day.*

Like most appointment schedules, the schedule we consider has a general framework based on time slots over the day or a portion of the day. There is a separate schedule for each doctor. The clinic has multiple arrivals scheduled in each time slot. As a consequence, our appointment system has relatively high volume: The doctor we focus on sees about 60 patients per day.

1.4. Organization

The paper analyzes the appointment-generated arrival process in steps, leading up to a full stochastic process model. We do not immediately present the final model because we regard the process leading up to the model as more important than the resulting model for the arrival process.

We start by focusing on what we regard as the most novel and important step: developing the model for the random schedule. In §2, we first examine the observed schedules to infer an underlying deterministic framework. Afterward, we view the actual schedule as a random modification of that framework. We find that the main deviation from a regular deterministic arrival pattern is variability in the schedule itself. Next, in §3, we view the actual arrivals as a random modification of the schedule and examine to what extent the actual arrivals adhere to the schedule. In §4, we

study the pattern of arrivals over each day and directly compare the arrivals to the schedule. In §5, we provide mathematical representations of the stochastic counting processes for the schedules and actual arrivals as well as a simple parsimonious model that may be a convenient substitute for mathematical analysis. We provide a classification for appointment-generated arrival processes in §6. This provides a basis for comparing the different doctors in this clinic with each other and with doctors in similar clinics. The classification scheme should also help to compare alternative appointment systems. We conclude in §7.

2. Defining and Modeling the Daily Schedule

We now examine the schedule and arrival data for one clinic doctor over his 22 morning (am) shifts. The arrivals planned for each day are given in a daily schedule, which has a specified number of arrivals in each of several evenly spaced ten-minute time intervals. Our schedule data are the 22 observed schedules for the doctor during his am shifts. Even though much can be learned from consulting the appointment manager, we try to see what can be learned directly from the data. While conducting the data analysis, we confirmed our observations with clinic doctors and administrators.

2.1. The Evolution of a Schedule

The actual schedule for a given day evolves over time, typically starting many weeks before the specified day. We do not consider the scheduling process; instead, we consider the evolution of the resulting schedule, we regard the evolution of the schedule as a stochastic process, with additions and cancellations occurring randomly over time. For each day, we define the final schedule as its value at the end of the previous day.

In the left-hand plot in Figure 1, we illustrate the evolution of the daily total number of patients scheduled over the previous year for the 22 days in the data set for 2013. The plot shows the specific appointment days as well, which are spread out between July and October.

The right-hand plot in Figure 1 presents a useful alternative view, showing the percentage of the final schedule reached k days before the appointment data as a function of k . For all 22 days, 100% of the schedule is filled at $k = 0$. We see much less variability in the right-hand plot than in the left-hand plot. The percentage of the schedule reached 30 days before appears at $k = -30$. Especially revealing is the average of the 22 sets of percentage data, which is shown by the single thick line. From this average plot, we see jumps at regular intervals, especially around 90 days (3 months) before the appointment date. The right-hand plot in Figure 1 shows that about 24% of

all appointments are made more than 93 days in advance, while about 30% are made between 93 and 84 days in advance (about 3 months). Moreover, about 30% are made in the last month, while about 13% are made in the last week.

2.2. New and Repeat Visits

There is increasing interest in the delays from request to appointment, including how to determine panel sizes (the pools of potential patients) for doctors; see Green et al. (2007), Liu and Ziya (2014), and Liu et al. (2010) and references therein. Unfortunately, our data set does not include a measure of the urgency or time sensitivity of each appointment, so we cannot determine whether patients are unable to get urgent appointments quickly enough. Fortunately, the data set does specify whether or not each scheduled arrival is a repeat visit or a new visit. Since 78% of all appointments are repeat visits, we conclude that the long intervals between the scheduling date and the appointment date do not imply that patients are failing to get urgent needs addressed promptly.

Figure 2 separately displays the evolution of the schedules for new and repeat visits, expanding upon the view in Figure 1. The figure panels show that this classification is very important. Figure 2 specifically shows that only about half of new patients wait for more than a week for an appointment. The median number of days between the appointment scheduling date and the actual appointment date is 93 for repeat visits and 24 for new patients.

2.3. Inferring a Deterministic Framework

From the perspective of the eventual arrival process over each day, the evolution of the schedule should not matter much if the final schedule reaches a nearly deterministic, regular form, which

Figure 1 The evolution of the daily number of patients scheduled over the previous year for 22 appointment days (left) and the percentage of patients who are scheduled k days in advance for each of the 22 appointment days (right). The thick line indicates the average over the 22 appointment days.

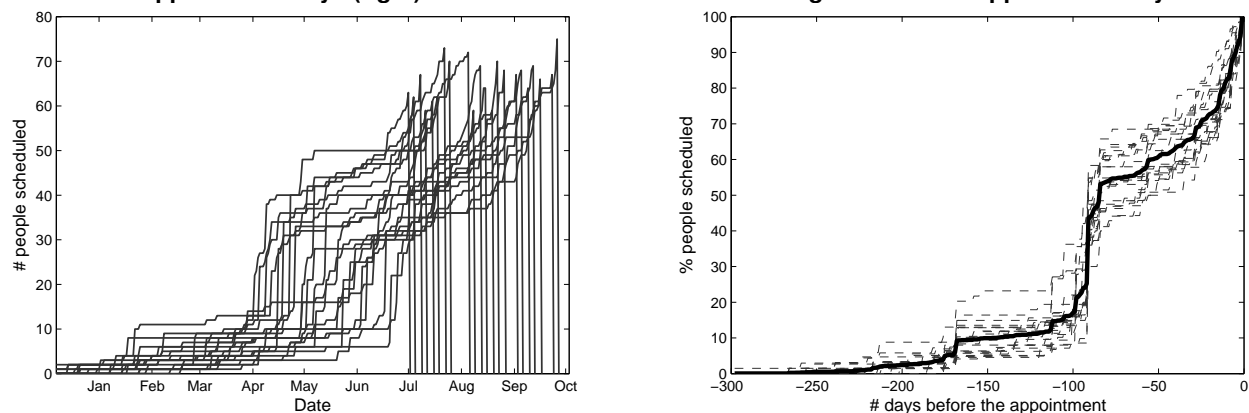
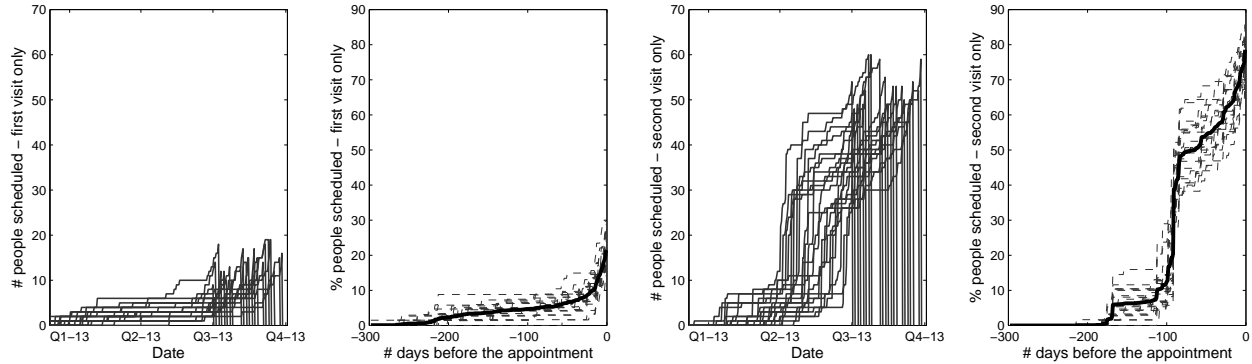


Figure 2 The evolution of the daily number of patients scheduled and the percentage of patients who are scheduled k days in advance for each of the 22 appointment days for new patients (left two panels) and repeat visits (right two panels). The thick line indicates the average over the 22 appointment days.



varies little from day to day. However, for the clinic, there is considerable variability in the schedules, so the evolution may matter.

We first define the schedule as the daily total plus the actual scheduled arrival times of all these patients. In particular, we define the schedule as its value at the end of the previous day, and we define arrivals on the same day as unscheduled arrivals. Given that definition, we next look for an underlying deterministic framework. The starting point for our data analysis is the 22 observed daily schedules. These are displayed in Table 1. Table 1 shows the number of patients scheduled for different ten-minute time slots (displayed vertically) over the am shifts of 22 days (displayed horizontally). Each ten-minute time slot is specified by its start time.

Most appointment schedules today are designed and managed to fit into a deterministic framework, usually using a computerized appointment management system. However, it seems prudent to look at the actual schedules and infer the realized framework from the data. Not all variability occurs because of adherence to the schedule; rather, the schedules show that there is substantial variability in the schedule itself.

We next define what we mean by a deterministic framework for the appointment schedule. A general deterministic framework has batches of size β_j customers arriving at intervals τ_j after an initial time 0 for $1 \leq j \leq \nu$. Thus, the associated arrival times are

$$\psi_j \equiv \sum_{i=1}^{j-1} \tau_i \quad \text{for } 1 \leq j \leq \nu \quad \text{and} \quad \psi_1 \equiv 0. \quad (1)$$

The framework has a total targeted number N_F and time T_F defined by

$$N_F = \sum_{j=1}^{\nu} \beta_j \quad \text{and} \quad T_F = \sum_{j=1}^{\nu-1} \tau_j = \psi_{\nu-1}. \quad (2)$$

Table 1 The number of patients scheduled in each 10-minute time slot (displayed vertically) during 22 morning shifts (displayed horizontally).

time slot	22 days in July-October 2013																						Avg	Var	Var/Avg		
7:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00			
8:00	0	0	0	0	0	0	1	1	1	0	1	0	0	0	1	0	0	0	0	0	0	0	1	1	0.32	0.23	0.71
8:10	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00
8:20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:50	3	4	5	4	4	4	4	4	1	3	2	1	4	2	4	4	2	4	5	4	3			3.41	1.30	0.38	
9:00	3	4	2	3	3	2	3	3	3	3	2	2	2	3	4	3	2	3	4	2	3	4	2		2.77	0.47	0.17
9:10	3	3	3	2	2	2	4	2	2	3	2	3	2	3	3	3	2	2	3	2	3	3	3		2.59	0.35	0.13
9:20	2	2	4	2	3	2	3	2	2	3	3	2	3	2	3	3	3	3	2	3	2	3	2		2.59	0.35	0.13
9:30	3	2	3	4	3	3	4	3	3	3	3	1	3	2	2	2	2	3	3	3	3	3	3		2.77	0.47	0.17
9:40	3	3	3	2	2	2	2	3	3	2	2	3	2	2	2	2	2	3	2	2	2	2	2		2.36	0.24	0.10
9:50	3	3	3	3	2	3	3	3	3	3	2	2	3	3	3	3	3	2	2	3	3	3	3		2.77	0.18	0.07
10:00	3	2	3	3	2	3	2	3	2	3	3	3	3	3	3	3	4	4	3	3	3	3	3		2.91	0.28	0.10
10:10	3	3	3	3	3	3	3	2	2	3	3	3	3	3	3	3	3	3	3	3	3	3	3		2.91	0.09	0.03
10:20	2	3	3	3	3	3	3	2	3	3	2	3	3	3	3	2	3	2	3	4	3	3	3		2.82	0.25	0.09
10:30	3	2	3	3	3	2	4	2	3	2	3	3	3	3	2	3	3	2	4	3	3	3	3		2.82	0.35	0.12
10:40	3	1	3	3	3	1	3	2	3	2	3	2	3	2	1	3	2	3	3	3	2	3	2		2.45	0.55	0.22
10:50	2	3	3	3	1	2	3	2	3	3	2	3	3	3	3	3	3	2	3	3	3	3	3		2.68	0.32	0.12
11:00	3	2	3	2	3	2	3	2	2	4	4	2	3	3	3	3	3	3	4	3	4	3	4		2.95	0.52	0.18
11:10	3	3	3	1	3	3	3	3	2	3	3	2	3	2	1	3	2	3	3	3	3	3	3		2.64	0.43	0.16
11:20	2	3	3	3	3	3	3	3	3	3	3	3	2	2	3	3	3	3	3	3	3	3	4		2.91	0.18	0.06
11:30	3	3	2	3	3	3	3	3	3	3	3	3	3	3	3	3	3	2	2	2	3	2	2		2.77	0.18	0.07
11:40	3	2	3	3	2	3	3	3	3	1	2	3	3	2	3	3	3	3	3	3	2	3	2		2.68	0.32	0.12
11:50	3	3	3	3	3	2	2	3	3	2	3	2	4	3	3	3	2	2	3	3	1	3	3		2.68	0.42	0.16
12:00	2	3	3	2	3	3	4	3	3	2	3	3	3	3	3	3	3	2	2	3	4	3	4		2.86	0.31	0.11
12:10	3	3	3	2	3	2	3	2	3	2	3	2	3	3	4	3	1	2	3	2	3	3	3		2.68	0.42	0.16
12:20	2	4	3	2	3	3	3	4	3	3	3	2	2	3	1	3	1	4	3	3				2.77	0.66	0.24	
12:30	2	1	0	0	0	3	3	3	3	2	2	2	3	3	3	2	4	3	1	2	3	3		2.14	1.27	0.59	
12:40	0	0	0	0	0	2	2	4	3	0	3	2	1	2	3	3	4	2	3	0	0	3		1.68	2.13	1.27	
12:50	0	0	0	0	0	0	1	4	0	0	0	0	3	4	0	2	0	4	0	0	4			1.00	2.67	2.67	
13:00	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		0.09	0.09	0.95	
Daily Total	63	62	67	59	61	62	73	70	72	59	69	64	59	70	68	67	68	64	69	66	67	75		66.09	21.32	0.32	
[8:50, 12:20] Total	60	61	67	59	60	57	67	60	61	57	63	60	56	62	57	61	60	58	59	65	64	64		60.82	9.77	0.16	
All slot avg	2.0	2.0	2.2	1.9	2.0	2.0	2.4	2.3	2.3	1.9	2.2	2.1	1.9	2.3	2.2	2.2	2.2	2.1	2.2	2.1	2.2	2.4		2.07	1.73	0.84	
All slot var	1.5	1.9	2.2	1.9	1.8	1.5	1.8	1.3	1.5	1.7	1.5	1.5	1.6	1.5	1.3	1.7	1.8	1.6	1.6	2.2	1.8	1.6					
All slot var/avg	0.7	1.0	1.0	1.0	0.9	0.8	0.8	0.6	0.6	0.9	0.7	0.7	0.8	0.6	0.6	0.8	0.8	0.8	0.7	1.1	0.8	0.7					
[8:50, 12:20] avg	2.7	2.8	3.0	2.7	2.7	2.6	3.0	2.7	2.8	2.6	2.9	2.7	2.5	2.8	2.6	2.8	2.7	2.6	2.7	3.0	2.9	2.9		2.76	0.42	0.15	
[8:50, 12:20] var	0.2	0.6	0.3	0.5	0.4	0.4	0.4	0.3	0.4	0.5	0.2	0.3	0.5	0.3	0.4	0.4	0.7	0.3	0.4	0.7	0.4	0.4					
[8:50, 12:20] var/avg	0.1	0.2	0.1	0.2	0.1	0.2	0.1	0.1	0.1	0.2	0.1	0.1	0.2	0.1	0.2	0.1	0.3	0.1	0.2	0.2	0.1	0.1					

A principal case is the *stationary framework*, with $\beta_j = \beta$ and $\tau_j = \tau$ for all j , which makes $N_F = \beta\nu$ and $T_F = (\nu - 1)\tau$, leaving the target parameter triple (β, τ, ν) , but there often are variations in practice. In the more general model, it is important to consider alternative nonstationary schedules that might be used or contemplated to improve various measures of performance.

From Table 1, we infer that the deterministic framework above is approximately valid with $\tau = 10$ minutes. However, the scheduled arrivals in each time slot are not constant over different days or over different times on each day. Table 1 indicates that for the am shifts of the doctor in the endocrinology clinic, the stationary framework is roughly valid as an idealized model, with $\beta = 3$, $\tau = 10$ minutes and $\nu = 22$ and starting at 8:50 and ending at 12:20 (including the intervals [8:50, 9:00) and [12:20, 12:30), closed on the left and open on the right), which we refer to as the interval [8:50, 12:20]. However, some shifts start as early as 8:00, while some shifts end as late as 13:00 (including the interval [13:00, 13:10)). The daily total for the stationary framework is $22 \times 3 = 66$, which matches the average daily total for the 22 days, even though the schedule is otherwise more variable.

Upon closer examination, we can see consistent structure in the schedule variability. First, we see that some days have higher daily totals, evidently because an effort is being made to respond

Table 2 The estimated distribution of the batch sizes (B_s) within the main interval [8:50, 12:20] and the estimated distribution of the total number of scheduled arrivals after the main interval (N_o) on the 10 at-capacity (AC) days, on the 12 overloaded (OL) days and on all days.

number k	$\hat{P}(B_s = k)$					$\hat{P}(N_o = k)$										
	1	2	3	4	5	0	1	2	3	4	5	6	7	8	9	10
10 at-capacity days	0.04	0.25	0.63	0.07	0.01	0.30	0.20	0.30	0.10	0.10						
12 overloaded days	0.02	0.27	0.63	0.08							0.25	0.17		0.25		0.33
All days	0.03	0.26	0.63	0.08	0.004	0.14	0.09	0.14	0.05	0.05	0.14	0.09		0.14		0.18

to high demand. Second, we see random batch sizes in the slots over the entire shift. We discuss each of these features in turn.

2.4. High-Demand Service with Overloaded and At-Capacity Schedules

In general, it seems useful to classify service systems with arrivals by appointment into two categories. First, there are the low-demand service systems, for which it is challenging to fill a target schedule. For such service systems, the randomness in the schedule is due to the random level of demand. We then might focus on the extent to which demand is adequate to fill the target schedule.

Second, there are the high-demand service systems, for which there is almost always ample demand, and often, there is excess demand. In that case, the system may or may not actually respond to the excess demand, i.e., it may or may not schedule more than the normal workload in order to meet that excess demand. Of course, there can be more complicated scenarios in which a service system oscillates between the low-demand and high-demand modes.

From Table 1, we infer first that the doctor operates as a high-demand service system and that indeed he responds to excess demand on some days, but not all days. We deduce from Table 1 that the appointment schedule is at capacity (AC) on some days but overloaded (OL) on other days. If we identify the deterministic framework for the am shifts to be the 22 ten-minute time slots in the interval [8:50, 12:20], then we observe that the daily totals within this interval are remarkably stable, having mean 60.82 and variance 9.77. In contrast, the full daily totals for the entire am shifts are much more highly variable, having a variance of 21.19.

This conclusion is further confirmed by the observation that the extra patients tend to be scheduled outside (after) the main am shift interval [8:50, 12:20]. In particular, we regard days with 5 or more appointments outside of (after the) the main interval as OL. By this definition, we see 12 OL days and 10 AC days among the 22 am shifts. Table 2 shows the distribution of the number of scheduled patients in these outside intervals, N_o , among AC and OL days.

Table 1 shows that overflows happen without any empty slots in between on OL days. Furthermore, we observe the possibility of interdependence over successive appointment days because OL

days are often followed by more OL days. As further confirmation of the idea that overload appears outside the main time interval, we also see higher numbers in the first shift, at 8:50 (the interval [8:50, 9:00)); this suggests that at least some of the patients scheduled in the first interval, at 8:50, are scheduled in response to pressure to provide service to more patients than the usual number. We note that this interval might be regarded as an overload period as well, though we choose not to do so. Moreover, the data show that the appointments in the OL portion of the schedule (at the beginning and the end) were consistently made far earlier than the other appointments.

When we next consider random batch sizes for the slots, we see that the batch sizes are very consistent inside the main interval on AC and OL days, further supporting the inference that arrivals outside the main interval primarily occur because of an effort to respond to excess demand.

2.5. Random Batch Sizes

Table 1 clearly indicates that the number of patients scheduled for each 10-minute time slot is variable. This distribution becomes quite consistent over the days and the time slots if we focus on the main time interval [8:50, 12:20]. Table 2 shows the distribution of the schedule within each time slot within the main interval for the AC days, the OL days and all days. The conclusions of §2.4 are supported by the fact that the estimated distributions for AC and OL days are very similar; those days tend to differ only to the extent that they have extra arrivals scheduled outside the main interval.

From Table 2, we conclude that it is reasonable to assume that the batch sizes in each of the time slots of the main time interval can be regarded as realizations of a random variable B_s , assuming values in the set $\{1, 2, 3, 4\}$ for any j . (We omit the value 5 because the frequency is so low, and we could also possibly omit the value 1 for the same reason.) In particular, we estimate the distribution as

$$P(B_s = k) = 0.03, 0.26, 0.63, 0.08, \quad 1 \leq k \leq 4, \quad (3)$$

respectively, so that

$$\begin{aligned} E[B_s] &= 2.76, & E[B_s^2] &= 8.02, & Var(B_s) &= 0.402 \\ & \text{and } SD(B_s) &= 0.634, \end{aligned} \quad (4)$$

respectively, for all j . The variance is considerably less than the mean, so we can conclude that the distribution of B_s is much less variable than Poisson. The squared coefficient of variation (scv, or the variance divided by the square of the mean) is remarkably low as well, being $c_B^2 = 0.053$.

2.6. Independence or Dependence among Batch Sizes

In §2.5, we focused on the distribution of the batch size of the scheduled arrivals in any time slot within the main time interval on any day. We now consider the joint distribution of the batch sizes over successive time slots on the same day.

Let $B_{s,j}$ be the scheduled batch size in slot j , $1 \leq j \leq 22$, on a given day. For simplicity from a stochastic modeling perspective, it is natural to assume that the batch variables $B_{s,j}$ in successive slots j are independent, which corresponds to appointments being made independently for specific slots. However, it may be more realistic to assume that the appointments are primarily made with a specific day in mind and that the actual appointments are distributed approximated evenly over the day, with the person or system creating the schedule only partly in response to patient requests regarding specific time slots. Alternatively, appointments may overflow into nearby slots, which should also create positive correlation. Therefore, in any context, it is interesting to explore the dependence among the scheduled batch sizes $B_{s,j}$ on each day.

To illustrate, let N_S be the daily total of the schedule (focusing on the main interval [8:50, 12:20] with $\nu = 22$ slots) and consider the case in which the distribution of B_s is independent of j . If the batch sizes are mutually independent, then

$$\text{Var}(N_S) = \nu \text{Var}(B_s). \quad (5)$$

In contrast, if we assume that the daily total is random and if we distribute it evenly among the slots, then we might have

$$B_s \approx \frac{N_S}{\nu} \quad \text{so that} \quad \text{Var}(N_S) = \nu^2 \text{Var}(B_s). \quad (6)$$

More generally, the dependence among the batch sizes might be usefully summarized by the correlations

$$\rho_{j_1, j_2} \equiv \text{corr}(B_{s, j_1}, B_{s, j_2}) = \frac{\text{cov}(B_{s, j_1}, B_{s, j_2})}{\sqrt{\text{Var}(B_{s, j_1}) \text{Var}(B_{s, j_2})}}. \quad (7)$$

We propose a model that enables us to incorporate a range of possibilities in a parsimonious manner. We assume that

$$\rho_{j_1, j_2} = \rho_S, \quad -1 \leq \rho_S \leq 1 \quad \text{for all} \quad j_1 \neq j_2. \quad (8)$$

We can then estimate the single pairwise correlation parameter ρ_S empirically in any given appointment setting.

Under assumption (8), we have

$$\begin{aligned}\sigma_S^2 \equiv \text{Var}(N_S) &= \sum_{j=1}^{\nu} \sum_{k=1}^{\nu} \text{Cov}(B_{s,j}, B_{s,k}) \\ &= \nu \text{Var}(B_s)[1 + (\nu - 1)\rho_S].\end{aligned}\quad (9)$$

We thus estimate the correlation ρ_S in (8) by

$$\rho_S \equiv \frac{\text{Var}(N_S) - \nu \text{Var}(B_s)}{\nu \text{Var}(B_s)(\nu - 1)}, \quad (10)$$

where we use our estimates of $\text{Var}(N_S)$ and $\text{Var}(B_s)$. From Table 1, our estimate of $\text{Var}(N_S)$ is 9.77; from (4), our estimate of $\nu \text{Var}(B_s)$ is $22 \times 0.402 = 8.80$. We thus estimate that ρ_S is $0.97/185 = 0.0052$, which is quite small. In fact, it is sufficiently small that we consider the i.i.d. model reasonable.

2.7. Outside the Main Time Interval

It remains to specify arrivals scheduled outside the main time interval. Since the average total outside is only about 10% of the full daily total and since we do not have a great amount of data overall, we will not try to develop a high-fidelity model. Based on the limited data provided by Tables 1 and 2, we classify a typical day as AC with probability $10/22$ and as OL with probability $12/22$. For days of each type, we allocate the total number of scheduled arrivals outside (after) the main interval according to the appropriate distributions specified for each day in Table 2. If the total number is 7 or fewer, then we divide the number into two parts, putting the larger or equal number in the first slot and the smaller or equal number in the second slot. If the total number is 8 or more, we divide the total into three parts, as evenly as possible, and put the numbers in decreasing order in the first three slots after the main interval.

2.8. Summary of the Schedule Model

In summary, the clinic data clearly indicate a well-defined structured framework, provided that we focus on a main time interval $[8:50, 12:20]$ containing 22 slots. The scheduled numbers in these slots can be regarded as i.i.d. random variables distributed as the random variable B_s , as in (3). Our analysis in §2.6 supports regarding these slot numbers as mutually independent.

Our doctor evidently experiences high demand. The data indicate that some days are OL, while others are AC. Based on the empirical data, we would say that a typical day is OL with probability $12/22$ but AC with probability $10/22$. The distributions of batch sizes within the slots are the same for these two kinds of days. In contrast, the number of extra arrivals scheduled after the

main interval does depend on this classification. As stipulated in §2.7, we allocate the totals randomly according to the distributions in Table 2, and we distribute them in a balanced, decreasing order over the outside intervals. Since the numbers outside are smaller, we devote less effort to developing a high-fidelity model for that part.

Only about 10% of the mean of the daily totals (66) is due to the arrivals scheduled outside the main interval (the mean inside is 60.8), while the variance 21.2 in the daily totals is primarily due to the random occurrence of arrivals scheduled outside the main interval because the variance inside is 9.77. (See equation (6) for a more precise statement.) Thus, we tentatively conclude that the greatest contributor to the overall variability of the schedules for the doctor in our study is the inconsistent response to extra demand. By examining both the scheduled and the realized arrivals for the other 15 doctors in the clinic, we find that this conclusion applies to all the other doctors as well: see Figures 1-3 and Figures 4-11 in our longer, more detailed study (Kim et al. 2015a). To draw a firm conclusion, we would need to consider data on the original demand, i.e., requests for appointments, including ones that were not satisfied or that were moved to another day.

3. Adherence to the Schedule

We now come to the question of adherence to the schedule. The level of adherence converts the schedule into the actual arrival process. We identify three familiar forms of additional randomness in the model: (i) no-shows, (ii) extra unscheduled arrivals and (iii) lateness or earliness. We first focus on the no-shows and the unscheduled arrivals, which together determine how the scheduled daily number of arrivals is translated into the actual daily total number of arrivals. In §4, we focus on lateness or earliness, which each have a significant impact on the pattern of actual arrivals over the day.

3.1. No-Shows

The no-shows are the scheduled arrivals that do not actually arrive. Instead of the number of actual arrivals in time slot j on a given day, which we denote by $B_{a,j}$, we now begin by focusing on the number from among the $B_{s,j}$ arrivals that were scheduled to arrive in slot j on that day that arrived *at some time on that appointment day*, which we denote by $B_{a|s,j}$, which necessarily satisfies the inequalities

$$0 \leq B_{a|s,j} \leq B_{s,j} \quad \text{for all } j. \quad (1)$$

The no-shows in slot j are thus defined as

$$B_{n,j} \equiv B_{s,j} - B_{a|s,j}. \quad (2)$$

Table 3 The number of no-shows ($B_{n,j} \equiv B_{s,j} - B_{a|s,j}$) for each 10-minute time slot j (displayed vertically) during 22 morning shifts (displayed horizontally).

time slot	22 days in July-October 2013																						Avg	Var	Var/Avg	
7:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:00	0	0	0	0	0	0	1	1	1	0	1	0	0	0	1	0	0	0	0	0	0	1	1	0.32	0.23	0.71
8:10	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00
8:20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00		
8:50	1	0	0	1	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0.23	0.28	1.23
9:00	0	1	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0.18	0.16	0.86
9:10	0	0	0	1	0	0	2	0	0	0	1	0	1	0	0	0	0	0	1	0	0	0	0	0.27	0.30	1.11
9:20	0	0	0	0	0	0	0	0	0	0	1	0	0	2	0	0	1	0	0	0	0	0	0	0.18	0.25	1.38
9:30	0	0	0	2	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.18	0.25	1.38
9:40	0	0	2	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0.18	0.25	1.38	
9:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00			
10:00	0	0	0	0	0	0	0	0	0	1	0	2	1	0	0	0	1	0	1	0	0	0	0.27	0.30	1.11	
10:10	0	0	0	1	0	0	0	1	0	0	1	1	0	0	0	0	1	0	0	0	2	0	0.32	0.32	1.01	
10:20	0	0	1	1	0	0	0	0	1	0	0	1	0	0	0	0	0	1	0	0	0	0	0.23	0.18	0.81	
10:30	0	0	0	0	0	0	0	1	1	0	0	0	0	0	1	0	0	0	0	0	1	0	0.18	0.16	0.86	
10:40	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	1	0	1	0.23	0.18	0.81	
10:50	0	0	1	0	0	1	1	0	0	0	0	0	1	0	0	0	0	0	0	0	2	0	0.27	0.30	1.11	
11:00	0	0	0	1	1	0	0	0	0	1	1	0	1	0	0	1	0	0	0	1	0	0	0.32	0.23	0.71	
11:10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	0	0.18	0.16	0.86	
11:20	0	0	0	0	0	1	0	0	0	0	0	1	0	0	1	0	1	0	0	1	0	0	0.23	0.18	0.81	
11:30	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	1	1	0	0	0	0	0	0.23	0.18	0.81	
11:40	0	0	1	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0.23	0.18	0.81	
11:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	0	0	0.14	0.12	0.90	
12:00	0	0	1	0	0	0	1	0	1	1	0	1	2	0	0	1	2	0	0	0	1	1	0.55	0.45	0.83	
12:10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0.05	0.05	1.00	
12:20	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0.09	0.09	0.95	
12:30	0	0	0	0	0	0	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0.14	0.12	0.90	
12:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0.09	0.09	0.95	
12:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0.05	0.05	1.00	
13:00	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.09	0.09	0.95	
Daily Total	3	2	6	8	2	2	10	7	5	4	6	10	6	2	5	5	6	5	4	7	5	10	5.45	6.35	1.17	
[8:50, 12:20] Total	2	2	6	8	1	2	8	4	4	4	4	10	6	1	4	5	6	5	4	7	4	7	4.73	5.64	1.19	
All slot avg	2.0	2.0	2.2	1.9	2.0	2.0	2.4	2.3	2.3	1.9	2.2	2.1	1.9	2.3	2.2	2.2	2.2	2.1	2.2	2.1	2.2	2.4	0.17	0.17	1.00	
All slot var	1.5	1.9	2.2	1.9	1.8	1.5	1.8	1.3	1.5	1.7	1.5	1.5	1.6	1.5	1.3	1.7	1.8	1.6	1.6	2.2	1.8	1.6	(across all days)			
All slot var/avg	0.7	1.0	1.0	1.0	0.9	0.8	0.8	0.6	0.6	0.9	0.7	0.7	0.8	0.6	0.6	0.8	0.8	0.8	0.8	0.7	1.1	0.8	0.7	(across all days)		
[8:50, 12:20] avg	2.7	2.8	3.0	2.7	2.7	2.6	3.0	2.7	2.8	2.6	2.9	2.7	2.5	2.8	2.6	2.8	2.7	2.6	2.7	3.0	2.9	2.9	0.21	0.21	0.98	
[8:50, 12:20] var	0.2	0.6	0.3	0.5	0.4	0.4	0.4	0.3	0.4	0.5	0.2	0.3	0.5	0.3	0.4	0.4	0.7	0.3	0.4	0.7	0.4	0.4	(across all days)			
[8:50, 12:20] var/avg	0.1	0.2	0.1	0.2	0.1	0.2	0.1	0.1	0.1	0.2	0.1	0.1	0.2	0.1	0.2	0.1	0.3	0.1	0.2	0.2	0.1	0.1	0.1	(across all days)		

These are shown in Table 3.

Table 3 shows that no-shows are rarer than in many other appointment systems: the number of no-shows ranges from 2 to 10 per day, with an average of 5.45 per day. The overall proportion of no-shows is 5.45/66.09, or 8.2%.

In general, we might try to model the no-shows quite carefully, as we did the schedule batch sizes $B_{s,j}$, but here, we simply assume that each scheduled patient fails to arrive in each slot on each day with probability $\delta = 0.082$, independently of all other patients. Overall, in the model, the total number of no-shows would have a binomial distribution with parameters equal to the total number, say n , of scheduled patients over all days and with probability $p = \delta = 0.082$, which would make the distribution approximately Poisson, with variance slightly less than the mean. Table 3 shows that the observed sample variance of the average number of no-shows is 6.35, which is only slightly greater than the overall average of 5.45. Hence, we conclude that the model with i.i.d. Bernoulli no-shows is quite well supported by the data.

3.2. Additional Unscheduled Arrivals

Some medical services have significant proportions of both unscheduled and scheduled arrivals. However, there are relatively few unscheduled arrivals at the clinic that we studied. As indicated

before, these are defined as scheduled arrivals that are scheduled on the same day (after the end of the previous day). On average, there are 2.18 unscheduled patients per day, among whom 1.95 arrived. In the Appendix, Table 7 shows all additional unscheduled arrivals, while Table 8 shows the additional unscheduled arrivals that actually arrived. The total number of unscheduled arrivals (that arrived) on all 22 days is 43. Table 8 shows that the total daily number of unscheduled arrivals exceeds 3 on only two days, with values of 4 and 7. The one exceptional day is evidently responsible for the variance for all days, 2.43, being somewhat larger than the mean. The unscheduled arrivals are somewhat more likely to be outside the main time interval, but that is consistent with our interpretation of outside the main time interval being a time for overload.

Paralleling our previous modeling, we could represent the daily total number of unscheduled arrivals within the main time interval as Poisson with mean 1.55 and those outside the main interval as Poisson with mean 0.40. We could then distribute those arrivals randomly (uniformly) within the respective time periods. With larger numbers, we might try more careful modeling. However, in general, some sort of Poisson process is natural for unscheduled arrivals because they are likely to be a result of individual people making decisions independently.

3.3. Daily Totals

We now examine the impact of no-shows and unscheduled arrivals on the actual daily totals of arrivals. Let N_A , N_S , N_N and N_U be the random daily total numbers of actual arrivals, scheduled arrivals, no-shows and unscheduled arrivals, respectively. In general, we have the basic flow conservation formula

$$N_A = N_S - N_N + N_U. \quad (3)$$

Combining the summary data from Tables 1, 3 and 8, we see that the means are

$$\begin{aligned} E[N_A] &= E[N_S] - E[N_N] + E[N_U] \\ &= 66.1 - 5.5 + 2.0 = 62.6 \end{aligned} \quad (4)$$

We see that the final mean daily number of arrivals $E[N_A] = 62.6$ is only about 5% less than the mean scheduled daily number $E[N_S] = 66.1$. Hence, from the perspective of the daily totals, there is strong adherence to the schedule.

Moreover, we see that the variability of the daily number of arrivals N_A is primarily due to the variability of the schedule. Indeed, the sample variances of the four daily numbers were

$$\begin{aligned} \text{Var}(N_A) &= 17.4, \quad \text{Var}(N_S) = 21.3, \quad \text{Var}(N_N) = 6.4 \\ \text{and } \text{Var}(N_U) &= 2.4. \end{aligned} \quad (5)$$

Note that the estimated variances are ordered by $Var(N_A) < Var(N_S)$. The dispersions (sample variance divided by the sample mean) are ordered as well:

$$\begin{aligned} Var(N_A)/E[N_A] &= 17.4/62.6 = 0.278 \\ &< 0.322 = 21.3/66.1 = Var(N_S)/E[N_S]. \end{aligned} \quad (6)$$

We also find that the data show a significant correlation between N_S and N_N . From further analysis, we find that $Var(N_S - N_N) = Var(N_A - N_U) = 18.91$. Since we necessarily have

$$\begin{aligned} Var(N_S - N_N) &= Var(N_S) + Var(N_N) - 2Cov(N_S, N_N) \\ &= 21.32 + 6.35 - 2Cov(N_S, N_N) = 18.91, \end{aligned}$$

we estimate the covariance $Cov(N_S, N_N)$ and the associated correlation $Cor(N_S, N_N)$ by

$$\begin{aligned} Cov(N_S, N_N) &= (27.67 - 18.91)/2 = 4.38 \quad \text{and} \\ Cor(N_S, N_N) &\equiv \frac{Cov(N_S, N_N)}{\sqrt{Var(N_S)Var(N_N)}} = \frac{4.38}{11.64} = 0.376, \end{aligned} \quad (7)$$

which is quite high. However, notice that two of the three largest no-show values (10) occur on days 7 and 22, which have the two largest daily totals, or 73 and 75, respectively. It remains to determine why the variables N_S and N_N should be positively correlated.

4. The Arrival Pattern over the Day

We now shift our attention to the pattern of arrivals over each day, given the daily totals. Here, ‘‘pattern’’ primarily means whether each patient arrives before or after the appointment time (earliness or lateness), but it might also mean systematic time dependence of the schedule, the no-shows or the unscheduled arrivals over the day.

4.1. The Big Picture of the Daily Pattern

Table 4 provides the details of the big picture for the time interval [8:50,12:20]. The first four columns of Table 4 show the average numbers scheduled, the percentage of no-shows, the percentage late and the percentage late by more than 15 minutes by half-hour intervals over the am shift, while the first four columns of Table 5 separately show the same summary statistics for new and repeat patients; these statistics are significantly different. Table 4 shows that the scheduled numbers and the no-shows are remarkably stable over time. As we have observed in previous sections, the main irregularity in the schedule occurs due to occasional overload scheduled outside these time intervals.

Table 4 Average numbers of scheduled arrivals for each 30-minute interval within the main 3.5-hour time interval as well as the proportions of no-shows and lateness and the average earliness (X^-), lateness (X^+) and overall deviation (X), plus 95% confidence intervals.

Interval	Avg # Scheduled	% No-show	% Late	% (Late>15 min)	Avg(X^+)	Avg(X^-)	Avg(X)
[8:50, 9:20)	8.8±0.7	7.9±4.8	21.2±6.9	12.3±5.5	35.8±18.7	-25.8±2.7	-11.4±6.9
[9:20, 9:50)	7.7±0.5	6.9±4.6	16.7±6.1	4.8±3.4	24.1±25.4	-35.7±5.6	-25.7±5.2
[9:50, 10:20)	8.6±0.4	6.8±4.4	15.0±6.7	6.4±3.4	20.3±10.9	-38.8±5.2	-30.2±6.5
[10:20, 10:50)	8.1±0.6	7.9±3.2	17.6±5.0	3.3±2.9	9.7±4.9	-45.0±7.3	-34.5±6.0
[10:50, 11:20)	8.3±0.5	9.0±3.9	13.6±4.4	5.4±3.9	18.4±11.2	-48.6±9.1	-39.2±8.7
[11:20, 11:50)	8.4±0.3	7.9±3.7	10.4±4.7	3.9±3.0	16.0±6.4	-61.2±9.1	-53.3±9.4
[11:50, 12:20)	8.2±0.5	9.3±4.1	9.5±5.4	3.8±3.5	12.7±6.6	-58.2±9.5	-51.7±9.8
[8:50, 12:20)	58.0±1.3	8.0±1.7	15.0±1.5	5.8±1.6	21.3±5.6	-44.9±3.0	-34.9±2.9

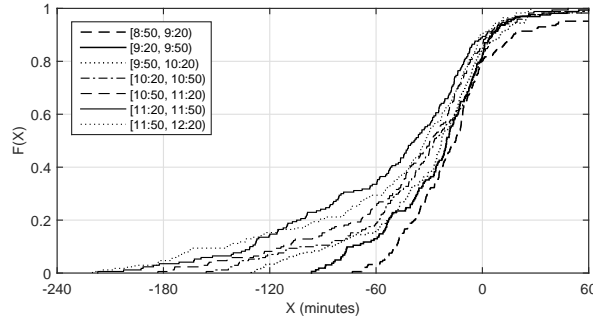
Table 5 Average numbers for new and repeat patients for the main interval and outside of the interval as well as the proportions of no-shows and lateness and the average earliness (X^-), lateness (X^+) and overall deviation (X), plus 95% confidence intervals.

Interval	Avg # Scheduled	% No-show	% Late	% (Late>15 min)	Avg(X^+)	Avg(X^-)	Avg(X)
New	14.2±1.3	5.5±2.4	22.2±4.4	7.7±3.4	23.2±9.7	-34.2±4.4	-21.2±3.9
New - [8:50, 12:30)	13.7±1.3	5.7±2.5	23.1±4.6	8.0±3.5	23.2±9.7	-33.9±4.7	-20.5±4.1
New - outside	0.5±0.3	0	0	0		-42.0±27.9	-42.0±27.9
Repeat	51.9±2.1	8.8±1.8	11.8±1.7	4.7±1.8	18.7±5.8	-49.3±3.3	-41.2±3.6
Repeat - [8:50, 12:30)	47.1±1.7	8.3±1.9	12.1±1.7	4.7±1.8	18.8±5.8	-48.7±2.9	-40.4±3.1
Repeat - outside	4.8±1.5	16.9±11.6	7.0±6.8	4.0±5.7	14.6±12.9	-62.6±24.1	-59.9±25.3

However, we see a different pattern in the lateness or earliness, as shown in the last four columns of Table 4. Specifically, Table 4 shows the percentage of patients arriving late, the average of the lateness X^+ among those patients arriving late, the average of the earliness X^- among those patients arriving early and the overall average lateness X (whose values are negative when the patient is early). Table 4 shows that the likelihood of lateness and the expected value of lateness tend to decrease over the day. In particular, we see that on average, 15% of the patients are late (arrive after the appointment time) each day, with an average lateness of $E[X^+] = 21$ minutes, but the percentage decreases over the day, from 21.2% in the first half hour to 9.5% in the last half hour. Meanwhile the average amount of lateness among these late patients, $E[X^+]$, decreases from 35.8 minutes to 12.7 minutes. In general, Table 4 shows that patients tend to arrive early, rather than late. This again reflects strong adherence to the schedule.

4.2. Toward a Model of the Deviations

We now look closer into the deviations of the actual arrival times from the scheduled arrival times. Figure 3 shows the *empirical cumulative distribution functions* (ecdfs) for the lateness for each of the half-hour time slots in Table 4. Figure 3 shows that the lateness consistently decreases over the day in the sense that each successive ecdf is stochastically less than the one before; see §9.1 of Ross (1996). (One ecdf is stochastically less than or equal to another if the entire ecdf lies *above*

Figure 3 The lateness ecdfs in each of the 30-minute intervals.

the other, e.g., the stochastically largest ecdf (with the most lateness) falls below all others and occurs in the first half hour.)

We now create a model of patient lateness (or earliness). The model has each scheduled arrival arrive at a random deviation from its scheduled arrival time. Let the arrivals scheduled to arrive at each time be labeled in some determined order, independent of the actual arrival time. We let the k^{th} arrival among the scheduled arrivals in time slot j (at time ψ_j in (1)) actually occur at time

$$A_{j,k} = \psi_j + X_{j,k} = \sum_{i=1}^{j-1} \tau_i + X_{j,k}, \quad (1)$$

where $X_{j,k}$ are mutually independent random variables, independent of the schedule (assuming arrivals are acting independently), and where $X_{j,k}$ is distributed as the random variable X_j with *cumulative distribution function* (cdf)

$$F_j(x) \equiv P(X_j \leq x), \quad -\infty < x < +\infty. \quad (2)$$

We allow X_j to assume both positive and negative values, representing arriving late and arriving early, respectively.

The ecdfs in Figure 3 can be regarded as estimates of the cdf F_j , and we use the same cdf F_j for all three ten-minute time slots j in the specified half hour. For a simple model, we might want a single cdf F , but Table 4 and Figure 3 present strong evidence that F_j should be allowed to depend on j , at least to some extent.

Finally, we note that it may be deemed useful to incorporate constraints on the arrival times at the beginning and the end of the time period. We might replace $A_{j,k}$ with the constrained version

$$A_{j,k}^c \equiv \max\{0, \min\{T_F, A_{j,k}\}\}. \quad (3)$$

To generate concrete stochastic models, we suggest fitting $P(X_j > 0)$ to the observed proportion of lateness in the half hour containing j and then fitting distributions to the observed values of

lateness X^+ or earliness X^- separately. The lateness probability estimates are given directly in Table 4. Similarly, we can use the ecdfs, denoted by $\hat{F}(x)$, to generate the model cdfs of X^+ and X^- , letting

$$F_{X^+}(x) \equiv P(X \leq x | X > 0) \equiv \frac{\hat{F}_j(x) - \hat{F}_j(0)}{1 - \hat{F}_j(0)} \quad \text{and}$$

$$F_{X^-}(x) \equiv P(X \leq -x | X \leq 0) \equiv \frac{\hat{F}_j(-x)}{\hat{F}_j(0)}, \quad x \geq 0. \quad (4)$$

Based on Figure 3, it appears that it should also be possible to use more elementary parametric models. We show the results of fitting exponential cdfs to X^+ and X^- over hours to the sample means in Figure 4. Figure 4 specifically shows that the estimated scv c^2 is less than 1 for X^- and greater than 1 for X^+ . Given the limited data, the exponential fit for X^- might be judged adequate, but we might also want to allow for greater variability in the lateness. We provide for that by considering a two-moment hyperexponential (mixture of two exponentials, with $c^2 > 1$ and balanced means, as on p. 137 of Whitt (1982)) in Figure 5.

Given that we have specified the cdf's F_j , we have completed construction of a full stochastic model of the arrival process that can be used to simulate arrivals to the clinic.

4.3. Comparing the Arrivals to the Schedule

We now directly compare the realized arrivals to the schedule. Table 9 in the appendix shows the difference between the numbers scheduled for the slot and the numbers that arrived in that slot for each time slot during the 22 days. The difference is often large, which we have seen must be primarily due to deviations from the scheduled arrival times, and especially earliness. Figures 6 and 7 provide summary views.

Let $S(t)$ and $A(t)$ count the number of scheduled and actual arrivals up to time t in the am shift. Figure 6 shows the histograms of the 22 observed values of the counting processes $S(t)$ and $A(t)$ for a few values of t : 10 am, 11 am, 12 pm and 1 pm. In particular, Figure 6 exposes systematic effects and shows the variability. Based on the figure, the arrivals scheduled for 10 am and 11 am tend to arrive early, but then they are about on time at 12 pm and fall slightly behind at 1 pm. We have seen that this is caused by the earliness of patient arrivals.

Figure 7 summarizes the data by plotting the average numbers of scheduled and actual arrivals for each of the ten-minute time slots within the 22 am shifts. Figure 7 also shows linear rate functions fit by least squares to the 22 averages of the scheduled and actual arrivals for each of the 22 ten-minute time slots within the main time interval (the solid lines). As should be expected,

Figure 4 Earliness (X^-) and lateness (X^+) histograms and associated exponential fits. Top to bottom: scheduled arrivals in [9,10), [10,11), and [11,12).

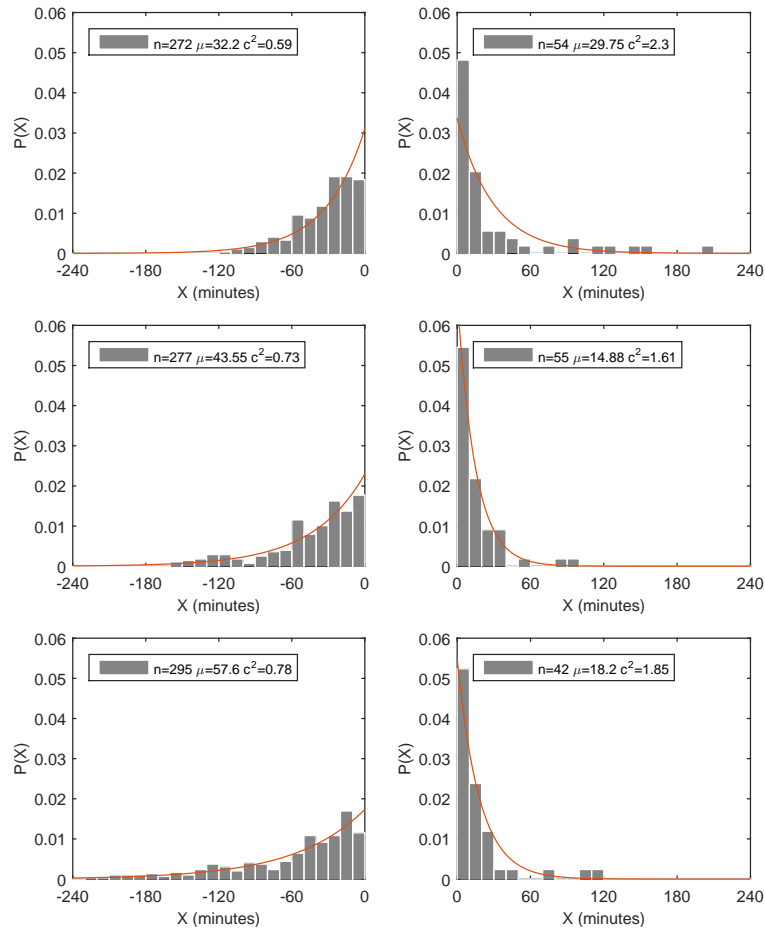
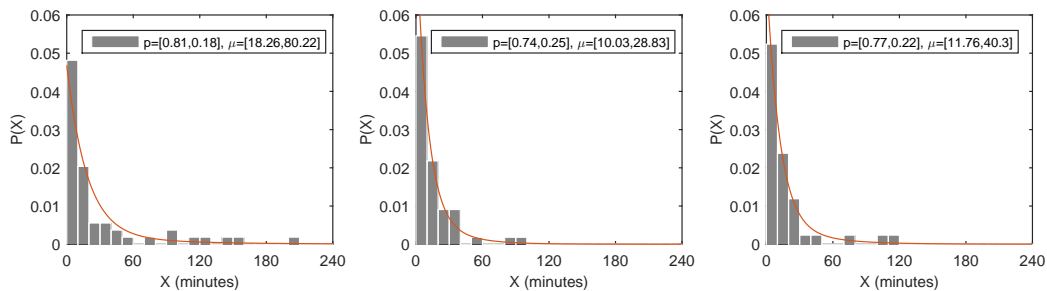


Figure 5 Lateness (X^+) histograms and associated hyperexponential (H_2) fits. Left to right: scheduled arrivals in [9,10), [10,11) and [11,12).



we see that the estimated rate function for the schedule within the main time interval is constant but that the estimated rate function of the actual arrivals is decreasing because of the tendency for patients to arrive early.

Finally, Figure 7 shows an additional continuous piecewise-linear estimated arrival rate function (the dotted lines) for the arrivals over the three intervals of the am shift. This dotted line has an

Figure 6 Histograms of the counting processes $S(t)$ and $A(t)$ at four different times. From left to right: 10 am, 11 am, 12 pm and 1 pm.

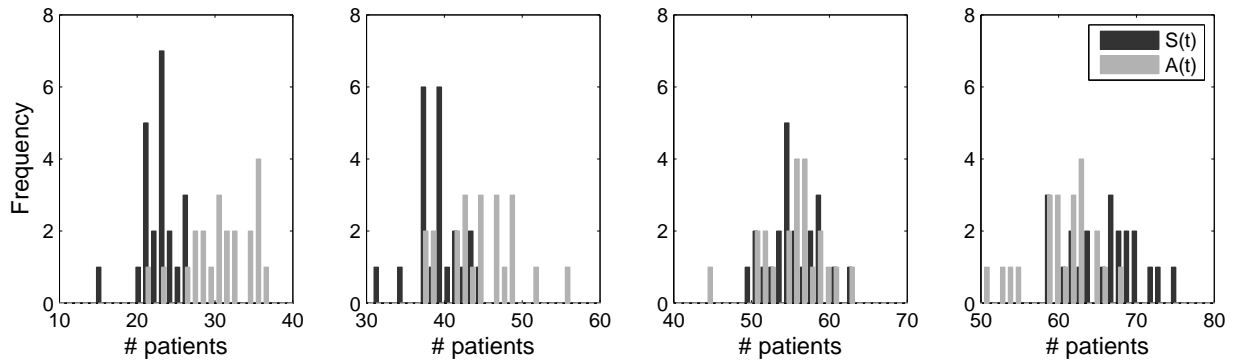
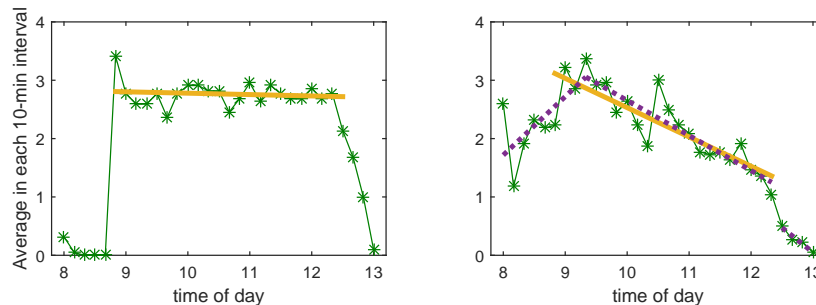


Figure 7 Plots of the average numbers of scheduled (left) and actual (right) arrivals in each of the 22 10-



extra linear piece before the main interval to account for the earliness. We will use this construction as the arrival rate resulting from the schedule in the main interval in the simple model constructed in §5.3.

5. Mathematical Models

In this section, we give concise mathematical representations of the stochastic counting processes $S(t)$ and $A(t)$, counting the number of scheduled and actual arrivals up to time t in the am shift, defined in terms of the model elements developed in previous sections.

The number of scheduled arrivals up to time t can first be expressed as the sum

$$S(t) = \sum_{j=1}^k B_{s,j}, \quad \psi_k \leq t < \psi_{k+1}, \quad k \geq 0, \quad (1)$$

for all t , for ψ in (1) and the batch sizes $B_{s,j}$. According to the model in §2, $B_{s,j}$ should be i.i.d. random variables with distribution in (3) inside the main time interval and distributed outside according to §2.7.

Let $A_S(t)$ count the number of scheduled arrivals that actually arrive up to time t . To define it, let the scheduled arrivals in each arrival epoch j (at time ψ_j) be ordered in some definite manner

not having to do with their actual arrival time. Let $I_{j,k} = 1$ if scheduled arrival k at time ψ_j actually arrives on that day and let $X_{j,k}$ be the deviation of the actual arrival time from the scheduled time. If $X_{j,k} > 0$, the arrival is late; otherwise, the arrival is early. (For simplicity in labeling, we have variables $X_{j,k}$ even when $I_{j,k} = 0$, but they will play no role.) We combine these two random features with the indicator random variable $I_{j,k}(t)$, defined by

$$I_{j,k}(t) \equiv 1_{\{I_{j,k}=1, X_{j,k} \leq t\}}, \quad -\infty < t < \infty, \\ 1 \leq k \leq B_{s,j}, \quad j \geq 0. \quad (2)$$

Given the definitions above, we can write

$$A_S(t) = \sum_{j=1}^{\infty} \sum_{k=1}^{B_{s,j}} 1_{\{I_{j,k}=1, X_{j,k} \leq t - \psi_j\}} \\ = \sum_{j=1}^{\infty} \sum_{k=1}^{B_{s,j}} I_{j,k}(t - \psi_j), \quad (3)$$

for $-\infty < t < +\infty$, where ψ_j is defined in (1). We may have $A_S(t) > 0$ for $t < 0$ because of early arrivals.

Let $A_U(t)$ ($A(t)$) count the number of unscheduled (all) arrivals by time t . Then we have

$$A(t) = A_S(t) + A_U(t), \quad \text{for all } t. \quad (4)$$

From §3.2, for the clinic, A_U would be independent of A_S , having two independent Poisson-based components, one for inside the main time interval and another for outside.

5.1. Conditional Means and Variances

Now suppose that the schedule is known, i.e., we know $B_{s,j}$ for all j , as would be the case at the end of the previous day in the clinic. Let the information about the schedule ($B_{s,j}$ for all j) be denoted by \mathcal{S} .

Since the ordering on k for each j is totally arbitrary, it is natural to assume that the joint distribution of $(I_{j,k}, X_{j,k})$ is independent of k for each j , and we make that assumption. The conditional cumulative arrival rate function for the scheduled arrivals given the schedule is then simply the conditional expected value

$$\Lambda_S(t|\mathcal{S}) \equiv E[A_S(t)|\mathcal{S}] = \sum_{j=1}^{\infty} B_{s,j} p_j(t), \quad (5)$$

$-\infty < t < +\infty$, where

$$\begin{aligned} p_j(t) &\equiv E[I_{j,k}(t - \psi_j)] = P(I_{j,k} = 1, X_{j,k} \leq t - \psi_j) \\ &= (1 - \delta)F_j(t - \psi_j), \end{aligned} \quad (6)$$

with $F_j(t) \equiv P(X_{j,k} \leq t)$, which is independent of k . As usual, the associated arrival rate function $\lambda_S(t|\mathcal{S})$ is the derivative with respect to t of the cumulative arrival rate function $\Lambda_S(t|\mathcal{S})$, i.e.,

$$\lambda_S(t|\mathcal{S}) = \sum_{j=1}^{\infty} B_{s,j}(1 - \delta)f_j(t - \psi_j), \quad (7)$$

where f_j is the probability density function (pdf) associated with the cdf F_j . The associated conditional variance is

$$V_S(t|\mathcal{S}) \equiv Var(A_S(t)|\mathcal{S}) = \sum_{j=1}^{\infty} B_{s,j}^2 p_j(t)(1 - p_j(t))$$

for $p_j(t)$ in (6).

5.2. The Total Mean and Variance of $A(t)$

The total arrival rate function is then

$$\begin{aligned} \Lambda(t) &\equiv E[A(t)] = E[A_S(t)] + E[A_U(t)] \\ &= E[\Lambda_S(t|\mathcal{S})] + E[A_U(t)] \\ &= \sum_{j=1}^{\infty} E[B_{s,j}](1 - \delta)f_j(t - \psi_j) + E[A_U(t)]. \end{aligned} \quad (8)$$

Applying the conditional variance formula, assuming that the random variables $B_{s,j}$ are mutually independent, the associated variance is

$$\begin{aligned} Var(A(t)) &= Var(A_S(t)) + Var(A_U(t)) \\ &= Var(E[A_S(t|\mathcal{S})]) + E[Var(A_S(t)|\mathcal{S})] + Var(A_U(t)) \\ &= \sum_{j=1}^{\infty} Var(B_{s,j})[(1 - \delta)f_j(t - \psi_j)]^2 \\ &\quad + \sum_{j=1}^{\infty} B_{s,j}^2 p_j(t)(1 - p_j(t)) + Var(A_U(t)). \end{aligned} \quad (9)$$

5.3. A Parsimonious Simplified Arrival Process Model

As before, we classify each day as AC or OL, but in our model, we make the days random, with OL days occurring with probability $12/22$ and AC days occurring otherwise, coinciding with the observed frequencies among the 22 days. We divide the overall time interval $[8:00,13:00]$ into two parts: before and after 12:30. We let D_F be the daily total during the final interval $[12:30,13:00]$, and we let it be conditional on whether the day is OL or AC. For each kind of day, we let the daily totals be distributed as in Table 2; that makes the mean number of OL days 7.60 and the mean number of AC days 1.50. (Thus, the overall mean number in $[12:30,13:00]$ is 4.82.) We then distribute the D_F arrivals among the intervals, as indicated in §2.7.

We let D_I be the random daily total for the initial interval $[8:00,12:20]$, and we treat all days the same. We let $E[D_I] = 66.1 - 4.8 = 61.3$, making it coincide with the observed average total of 66.1 in Table 1. We let the variance coincide roughly with the variance of the schedule inside the interval in Table 1, so that $Var(D_I) = 10.0$. We can use a Gaussian distribution (rounded to the nearest integer) with this estimated mean and variance. Alternatively, we can fit a binomial distribution with parameter pair (n, p) to this mean and variance, yielding two equations with two unknowns: $E[D_I] = np = 61.3$ and $Var(D_I) = np(1-p) = 10$, so that $(1-p) = 10/61.3 = 0.163$ and $n = 61.3/0.837 = 73.2$, rounded to 73. Hence, we regard D_I as binomial: $(n, p) = (73, 0.837)$.

Given D_I , the daily total in the initial interval, we let these arrivals be i.i.d. over the initial interval $[8:00,12:20]$, with a pdf proportional to the continuous two-piece arrival rate function in Figure 7, i.e., with a pdf equal to the arrival rate function divided by its integral over the interval.

In Kim et al. (2015b), binomial-uniform and Gaussian-uniform models were proposed. Our model here differs in two respects. First, we treat the final subinterval $[12:30,13:00]$ separately, accounting for whether or not the day is OL or AC. Second, we treat the initial interval similarly, but our more careful analysis here suggests a non-uniform density for the individual arrivals. We propose a scaled version of the continuous piecewise-linear curve on the right in Figure 7, which should better fit the actual arrival rate.

6. Classification of Appointment-Generated Arrival Processes

While diverse appointment systems should have much in common, there also can be important differences. A useful first step when considering appointment systems and appointment-generated arrival processes is to classify the system. Our analysis of the clinic, summarized in Table 6, helps to show how that can be done. There are three main steps, which we explain in detail below.

Table 6 Steps to classify an appointment-generated arrival process and the steps' application to the arrivals for the doctor at the clinic.

Category	Issue	For the doctor at the endocrinology outpatient clinic
General	time frame for arrivals time from scheduling to appointment time sensitivity (urgency) of appointment repeat versus new scale variability of the arrival process	one morning shift on a single day mostly 1-4 months not known 78% of visits are repeat moderately large, average daily total of 66 significant but less than Poisson, dispersion $V/M = 0.3$ for daily totals
Schedule	variability of the schedule deterministic framework primary deviation from the framework high or low demand extent of overload manifestation of overload distribution of the main schedule	significant but less than Poisson, dispersion $V/M = 0.3$ for daily totals identifiable as 22 ten-minute intervals with batches of size 3 extra arrivals scheduled outside the main interval high demand 12 of 22 days overloaded, with overload producing 10% of daily totals overload occurs outside, usually after, the main interval the data support i.i.d. batches with mean 2.76 in all time slots
Adherence	no-shows unscheduled arrivals deviations (lateness or earliness)	relatively few no-shows, or about 8.5% relatively few unscheduled arrivals, or about 2 per day (3%) significant deviations of about 60 minutes, but mostly early; about 15% late, with average conditional lateness of about 20 minutes

Step 1. General Classification We first identify the *time frame*, which we take to be a day. However, there are two different perspectives: first, the times when the arrivals occur, and second, the times when the appointments are actually made. We primarily focus on the times when the arrivals occur, aiming to understand variability over the day.

However, as in the clinic studied here, the appointments may have been made over a much longer time frame, weeks or even months before the appointment day, so the delay in getting an appointment may occur over a longer time scale. With such long delays between the date that the appointment is scheduled and the appointment date, we have observed that it is important to consider whether arrivals represent, perhaps routinely, *repeat visits* or are new requests. Especially in healthcare, an important question is whether the system can respond well to urgent requests for service. Unfortunately, time sensitivity or urgency was not part of the clinic arrival data analyzed here, but we were able to identify repeat visits, which accounted for 78% of all visits. However, it is important to recognize that long delays do not necessarily mean that patients with urgent problems are experiencing excessive delays before their needs can be addressed. In general, for healthcare appointment systems, it would be useful to have information on the *delay sensitivity or urgency* of the service to be provided.

We next focus on the *scale*, determined by the typical daily totals. Is the scale large or small? The clinic doctors considered here operate on a fairly large scale, with our specific doctor seeing about 66 patients in each shift (am or pm).

Assuming that our goal is to understand the arrival process over a single day and possibly to make improvements in this process, the next question is the *level of variability in the appointment-generated arrivals*. Are the arrivals highly regular or not? The analysis is devoted to the case in

which the arrivals exhibit significant variability. An initial rough classification of the variability is the *dispersion* or variance-to-mean ratio V/M of the daily totals.

The remaining classification is aimed at exposing the primary sources of the variability observed in the arrivals. Careful analysis is then devoted to identifying and quantifying the important sources of that variability. Here, it is natural to start with the schedule.

Step 2. The Schedule Given that the actual arrivals are irregular, we ask if the scheduled arrivals are also irregular, additionally exhibiting a significant level of variability. For our doctor in the clinic, we found that the schedule is indeed quite irregular, exhibiting significant variability and that too can be roughly quantified based on the dispersion of the daily totals. In fact, we concluded that the primary source of variability in the arrivals is the variability in the schedule. This is supported by the fact that both the variance and the dispersion of the scheduled daily totals are greater than for the actual arrivals.

Whether the schedule is regular or not, we want to identify the fundamental deterministic framework, if possible. In general, a first step in analyzing the schedule is to infer this framework. An orderly framework might be communicated by system managers, but it is important to consider data showing what actually happened. From our examination of the schedule for the 22 am shifts, we were able to identify a stationary framework involving small batches of arrivals at ten-minute intervals during the time interval [8:50, 12:20].

We then ask what are the major deviations from this framework. In the present analysis, we found that the batch sizes in each time slot are variable, but the largest deviation from that framework was due to extra arrivals scheduled outside the main time interval.

In general, it is evidently important to determine whether the service system is a *high-demand or low-demand system*. Is the variability due to an uncertain ability to fill the schedule in the presence of low demand or because of an uncertain response to pressures to meet high demand? Or do we see a combination of these? We concluded that our doctor in the clinic consistently operates as a high-demand system, with a significant response to high demand. In particular, 12 of the 22 days are overloaded (OL), and the remaining 10 are at capacity (AC).

We then come to the distribution of the scheduled arrivals in the main interval. We concluded in §2.5 and §2.6 that the scheduled arrivals in the 22 daily time slots in the main interval can be regarded as i.i.d. random variables with the distribution in (3), which has mean 2.76. We found relatively low variability in the scheduled arrivals within this main interval.

Step 3. Adherence to the Schedule We next shift attention to the adherence to the schedule. Here, we focused on three ways that the arrivals might not adhere to the schedule: (i) no-shows, (ii) extra unscheduled arrivals and (iii) deviations in actual arrival times from the scheduled times. Since our clinic data included cancellations, no-shows were easily identifiable as scheduled arrivals that never occurred. Given that all arrivals were included in our clinic data and that our definition of the schedule was based on its value at the end of the previous day, we defined unscheduled arrivals as arrivals that were scheduled and that arrived on the current day.

It is well known that no-shows and unscheduled arrivals can be quite frequent in appointment-generated arrival processes. However, in the clinic studied here, there were relatively low percentages of no-shows and unscheduled arrivals. In particular, the average percentage of no-shows for our doctor in the clinic was about 8.5%. This level was fairly constant over the day but was somewhat higher during the first intervals of the am shift. The average number of unscheduled arrivals in the clinic was only about 2 per day, which was 3% of the daily total. About half of those occurred outside the main interval, again indicating effort to respond to high demand.

Significant deviations in the actual arrival time from the scheduled arrival times were observed, with values of approximately 60 minutes, but most were due to early arrivals. Only about 6% of the arrivals were late by more than 15 minutes. Overall, we conclude that the adherence to the schedule was very good relative to that in other appointment systems.

7. Conclusions

The Principal Source of Variability Is the Schedule. In this paper, we have examined an appointment-generated arrival process for one doctor in an endocrinology clinic. As a consequence of the appointment system, the arrival process tends to be much less variable than a Poisson process, but is also not nearly a regular deterministic arrival process. The dispersion (variance-to-mean ratio) is about 0.3. As others have observed before, some variability is due to no-shows, extra unscheduled arrivals and deviations of the actual arrival times from the scheduled appointment times, but §3.3 shows that the dominant source of variability in the arrival process is the schedule itself. In particular, surprisingly, the inequality in (6) shows that the dispersion of the daily schedule is actually greater than the dispersion of the daily arrivals itself.

New Stochastic Arrival Process Models. Our data analysis has culminated in both a detailed stochastic model in §2.8 and §3 and a simplified stochastic model in §5.3 that can be used to simulate the arrival process of patients to see the doctor in the clinic. The fitting process should be

useful for analyzing the other doctors in this clinic as well as for other applications, and simulation experiments can be used to evaluate operational procedures in the clinic.

What Is Generalizable? (i) Variations of the specific arrival process stochastic models developed here may be useful for analyzing other outpatient clinics, but what we think is widely generalizable is *the data-analysis process, rather than the model*. Consistent with earlier work, we advocate carefully examining no-shows, extra unscheduled arrivals and punctuality. However, before doing those steps, we recommend looking at randomness in the schedule. It may even be important to view the schedule as a stochastic process. We do not have data on the original demand in the current analysis, but we would also advocate collecting information on requests for appointments, including ones that were not scheduled or that were moved to alternate days and times. Additionally, we recommend determining how the schedule relates to the original demand.

(ii) The specific arrival process models may also be useful more widely. Especially promising is the parsimonious model with Gaussian daily totals and, given those daily totals, i.i.d. arrival times within the day with a non-uniform probability density that takes account of the earliness and lateness of the patients. It is reasonable to anticipate that the earliness or lateness will alter the arrival rate during the day, as we have discovered.

(iii) Even more broadly, it is important to recognize that appointment-generated arrival processes are likely to be neither solely deterministic and evenly spaced nor solely Poisson; rather, many systems will have variability in between those two extremes, just as we have seen.

What Is the Practical Relevance? In this paper, we have not performed a complete performance analysis of the endocrinology outpatient clinic, so we have not yet improved the performance of that clinic. However, based on the long history of modeling and analysis of outpatient clinics briefly surveyed in §1.1, modeling and analysis can improve system performance. Thus, we did this work with the conviction that improved arrival process models can produce improved performance.

We see two principal ways that the stochastic model of the appointment system can be used to improve the performance of the clinic, and similar stochastic models can also be used to improve performance in other appointment system applications. First, the model provides a basis for analyzing the performance of the clinic with the given arrival process by conducting standard performance (queueing) analyses after incorporating an additional detailed analysis of the patient processing and flow after arrival, which we do not consider here. Second, the model can be used to consider alternative scheduling strategies to achieve various objectives, such as reducing the variability of

the schedule and thus reducing the variability in the doctor workloads or ensuring that patients with urgent needs have limited delays in getting an appointment.

Classification of Appointment-Generated Arrival Processes. In addition to gaining a better understanding of the appointment-generated arrival process in the endocrinology clinic, we have learned how to think about appointment-generated arrival processes more generally. A useful first step when considering appointment systems and appointment-generated arrival processes is to classify the system as done in Table 6 for our analysis of the clinic. For any new appointment system to be considered, we recommend seeking this information. After evaluating both the schedule and the adherence to the schedule by comparing them to what is desired, one could consider ways to improve both the schedule and the adherence.

Acknowledgments

Support was received from NSF grant CMMI 1265070. The authors thank Dr. Sang-Man Jin in the Division of Endocrinology and Metabolism at Samsung Medical Center for providing advice and Mohamad Soltani for a helpful comment.

References

- Bailey, N. T. J. 1952. A study of queues and appointment systems in hospital outpatient departments, with special reference to waiting times. *Journal of the Royal Statistical Society A* **14** 185–199.
- Cayirli, T., E. Veral. 2003. Outpatient scheduling in health care: a review of literature. *Production and Operations Management* **12**(4) 519–549.
- Chakraborty, S., K. Muthuraman, M. Lawley. 2010. Sequential clinical scheduling with patient no-shows and general service-time distributions. *IIE Transactions* **42**(5) 354–366.
- Chand, S., H. Moskowitz, J. B. Norris, S. Shade, D. R. Willis. 2009. Improving the patient flow at an outpatient clinic: study of sources of variability and improvement factors. *Health Care Management Science* **12** 325–340.
- Feldman, J., N. Liu, H. Topaloglu, S. Ziya. 2014. Appointment scheduling under patient preference and no-show behavior. *Operations Research* **62**(4) 794–811.
- Fetter, R. B., J. D. Thompson. 1965. The simulation of hospital systems. *Operations Research* **13**(5) 689–711.
- Green, L. V., S. Savin, M. Murray. 2007. Providing timely access to care: what is the right patient panel size? *Joint Commission Journal on Quality and Patient Safety* **33**(4) 211–218.
- Guo, M., M. Wagner, C. West. 2004. Outpatient scheduling - a simulation approach. R. G. Ingalls, M. D. Rossetti, J. S. Smith, B. A. Petyers, eds., *Proceedings of the 2004 Winter Simulation Conference*. 1981–1987.
- Gupta, D., B. Denton. 2008. Appointment scheduling in health care: Challenges and opportunities. *IIE transactions* **40**(9) 800–819.

- Hall, R. W., ed. 2006. *Patient Flow: Reducing Delay in Healthcare*. Springer, New York.
- Hall, R. W., ed. 2012. *Handbook of Healthcare System Scheduling*. Springer, New York.
- Harper, P. R., H. M. Gamlin. 2003. Reduced outpatient waiting times with improved appointment scheduling: a simulation modelling approach. *OR Spectrum* **25**(3) 207–222.
- Hassin, R., S. Mendel. 2008. Scheduling arrivals to queues: A single-server model with no-shows. *Management Science* **54**(3) 565–572.
- Jacobson, S. H., S. N. Hall, J. R. Swisher. 2006. Discrete-event simulation of health care systems. R. W. Hall, ed., *Patient Flow: Reducing Delay in Healthcare Delivery*, chap. 8. Springer, 211–252.
- Jouini, O., S. Benjaafar. 2012. Queueing systems with appointment-driven arrivals, non-punctual customers, and no-shows. *Working paper*.
- Jun, J. B., S. H. Jacobson, J. R. Swisher. 1999. Application of discrete-event simulation in health care clinics: a survey. *The Journal of the Operational Research Society* **50**(2) 109–123.
- Kaandorp, G. C., G. Koole. 2007. Optimal outpatient appointment scheduling. *Health Care Management Science* **10**(3) 217–229.
- Kim, S.-H., P. Vel, W. Whitt, W. C. Cha. 2015a. Analysis of arrival data from an endocrinology clinic. Columbia University, <http://www.columbia.edu/~ww2040/allpapers.html>.
- Kim, S.-H., P. Vel, W. Whitt, W. C. Cha. 2015b. Poisson and non-Poisson properties in appointment-generated arrival processes: the case of an endocrinology clinic. *Operations Research Letters* **43** 247–253.
- Liu, N., S. Ziya. 2014. Panel size and overbooking decisions for appointment-based services under patient no-shows. *Production and Operations Management* **23**(12) 2209–2223.
- Liu, N., S. Ziya, V. G. Kulkarni. 2010. Dynamic scheduling of outpatient appointments under patient no-shows and cancellations. *Manufacturing and Service Operations Management* **12**(2) 347–364.
- Luo, J., V. G. Kulkarni, S. Ziya. 2012. Appointment scheduling under patient no-shows and service interruptions. *Manufacturing and Service Operations Management* **14**(4) 670–684.
- Moore, C. G., P. Wilson-Witherspoon, J. C. Probst. 2001. Time and money: effects of no-shows at a family practice residency clinic. *Family Medicine - Kansas City* **33**(7) 522–527.
- Neal, R. D., D. A. Lawlor, V. Allgar, M. Colledge, S. Ali, A. Hassey, C. Portz, A. Wilson. 2001. Missed appointments in general practice: retrospective data analysis from four practices. *British journal of general practice* **51**(471) 830–832.
- Ross, S. M. 1996. *Stochastic Processes*. 2nd ed. Wiley, New York.
- Swartzman, Gordon. 1970. The patient arrival process in hospitals: statistical analysis. *Health services research* **5**(4) 320.
- Swisher, J. R., J. B. Jun, S. H. Jacobson, O. Balci. 2001. Modeling and analyzing a physician clinic environment using discrete-event (visual) simulation. *Computers and Operations Research* **28** 105–125.

Table 7 The extra unscheduled arrivals, i.e., the same-day arrivals $B_{u,j}$ scheduled for slot j on each of the 22 days.

time slot	22 days in July-October 2013																Avg	Var	Var/Avg						
7:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
8:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00					
8:10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
8:20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
8:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
8:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
8:50	0	0	0	0	0	0	0	0	0	1	0	0	1	0	0	1	0	0.14	0.12	0.90					
9:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.09	0.09	0.95					
9:10	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	1	0.14	0.12	0.90					
9:20	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00					
9:30	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00					
9:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0.09	0.09	0.95					
9:50	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0.09	0.09	0.95					
10:00	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0.05	0.05	1.00					
10:10	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0.05	0.05	1.00					
10:20	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0.14	0.12	0.90					
10:30	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0.09	0.09	0.95					
10:40	0	1	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0.14	0.12	0.90					
10:50	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0.09	0.09	0.95					
11:00	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0.09	0.09	0.95					
11:10	0	0	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0.09	0.09	0.95					
11:20	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0.05	0.05	1.00					
11:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
11:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
11:50	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0.09	0.09	0.95					
12:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00							
12:10	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0.14	0.12	0.90					
12:20	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00					
12:30	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0.14	0.12	0.90					
12:40	0	0	0	0	0	0	0	0	0	1	2	0	0	0	0	0	0	0.14	0.22	1.60					
12:50	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0.09	0.09	0.95					
13:00	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0.09	0.09	0.95					
Daily Total	2	3	0	1	1	1	1	3	2	1	3	8	3	2	2	2	3	4	0	2	3	1	2.18	2.82	1.29
[8:50, 12:20] Total	2	3	0	1	1	1	1	2	1	2	5	2	1	2	2	3	3	0	1	2	1		1.68	1.27	0.76

Wang, R., O. Jouini, S. Benjaafar. 2010. Queueing systems with appointment-driven arrivals, non-punctual customers and no-shows. Working paper, Singapore Management University.

Wang, R., O. Jouini, S. Benjaafar. 2014. Service systems with finite and heterogeneous customer arrivals. *Manufacturing and Service Operations Management* **16**(3) 365–380.

Welch, J. D., N. T. J. Bailey. 1952. Appointment systems in hospital outpatient departments. *The Lancet* **259**(6718) 1105–1108.

Whitt, W. 1982. Approximating a point process by a renewal process: two basic methods. *Operations Research* **30** 125–147.

Whitt, W. 1983. The queueing network analyzer. *Bell System Technical Journal* **62**(9) 2779–2815.

Zacharias, C., M. Pinedo. 2014. Appointment scheduling with no-shows and overbooking. *Production and Operations Management* **23**(5) 788–801.

Zonderland, M. E., R. J. Boucherie. 2012. Queueing networks in healthcare systems. R. W. Hall, ed., *Handbook of Healthcare System Scheduling*, chap. 9. Springer, 201–244.

Appendix

Table 8 The unscheduled arrivals that actually arrived ($B_{a|u,j}$) for slot j on each of the 22 days.

Slot	Different Days																				Avg	Var	Var/Avg		
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0					
7:50	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
8:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
8:10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
8:20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
8:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
8:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
8:50	0	0	0	0	0	0	0	0	0	0	1	0	0	1	0	0	1	0	0	0	0.14	0.12	0.90		
9:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.09	0.09	0.95		
9:10	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	1	0	0	0	0	0.14	0.12	0.90		
9:20	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00		
9:30	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00		
9:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0.09	0.09	0.95		
9:50	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0.09	0.09	0.95		
10:00	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00		
10:10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
10:20	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	1	0	0.14	0.12	0.90		
10:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0.05	0.05	1.00		
10:40	0	1	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0.14	0.12	0.90		
10:50	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0.09	0.09	0.95		
11:00	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0.09	0.09	0.95		
11:10	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00		
11:20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0.05	0.05	1.00		
11:30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
11:40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
11:50	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0.09	0.09	0.95		
12:00	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.00				
12:10	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	1	0	0.14	0.12	0.90		
12:20	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.05	0.05	1.00		
12:30	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0.09	0.09	0.95		
12:40	0	0	0	0	0	0	0	0	0	1	2	0	0	0	0	0	0	0	0	0	0.14	0.22	1.60		
12:50	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	0	0.09	0.09	0.95		
13:00	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0.09	0.09	0.95		
Daily Total	2	3	0	1	1	1	1	2	2	0	3	7	3	1	2	2	3	4	0	2	2	1	1.95	2.43	1.24
[8:50, 12:20] Total	2	3	0	1	1	1	1	2	0	2	4	2	0	2	2	3	3	0	1	2	1	1.55	1.21	0.78	

Table 9 The difference between the number of patients scheduled to arrive for slot j ($B_{s,j}$) and the number of patients who actually arrived for slot j ($B_{a,j}$) on each of the 22 days. The summary statistics are based on absolute values.

Slot	Different Days																				Avg	Var	Var/Avg		
	0	-1	0	0	0	0	0	0	0	0	0	0	0	0	0	-3	0	0	0	0					
7:50	0	-1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.18	0.44	2.43		
8:00	0	-2	-3	-1	-7	-2	-2	-4	1	-2	-4	-4	-2	-5	-1	-3	-2	-4	0	-1	-2	0	2.36	3.10	1.31
8:10	-3	0	0	0	-2	-4	-1	1	-3	0	-4	0	-2	-3	-2	0	0	-2	0	0	1.23	2.09	1.70		
8:20	-2	-1	-3	-3	-1	-2	-2	-5	-3	-1	0	-3	0	-1	-3	-4	0	-4	0	-2	-2	1.91	2.18	1.14	
8:30	0	-2	-3	-1	-2	-2	-1	-2	0	-3	-2	-2	-2	-3	-3	-2	-1	-4	-5	-7	2.32	2.51	1.08		
8:40	-3	-5	-4	-5	-3	-3	-3	-3	-1	-2	0	-2	0	-3	-1	-1	-1	-3	0	-2	-2	2.18	2.16	0.99	
8:50	2	2	4	3	3	2	1	3	2	3	2	3	2	1	0	0	1	0	1	0	2.00	2.00	1.00		
9:00	-1	-2	1	0	1	-1	-2	-1	-3	-2	0	0	-1	-3	0	1	-1	0	0	0	1.09	0.94	0.87		
9:10	0	1	1	0	-3	0	1	0	-1	0	-3	1	-2	3	1	-1	-2	1	0	-2	1	1.18	1.01	0.86	
9:20	0	-1	-1	-2	-1	-4	-1	-2	-3	-2	-1	1	1	-1	1	-1	2	-2	0	1	1.41	0.82	0.59		
9:30	-2	1	0	-1	1	2	1	1	1	-3	-4	1	-1	0	1	-3	0	2	1	1	1	1.32	1.08	0.82	
9:40	2	2	-1	0	-2	-2	0	2	0	2	-3	-3	1	1	-1	0	-4	-2	-1	-2	1.50	1.12	0.75		
9:50	1	2	-2	2	0	-1	-2	0	-1	2	-1	-1	1	2	0	2	1	0	-3	1	3	1.32	0.80	0.61	
10:00	-4	0	1	0	-1	3	-4	1	1	2	0	1	-1	-2	0	1	3	2	0	0	2	1.36	1.58	1.16	
10:10	0	-1	2	3	1	0	1	2	2	-1	1	3	-1	-2	-2	2	0	1	1	0	2	1.32	0.80	0.61	
10:20	1	1	1	2	0	2	1	1	2	0	-2	3	2	1	2	-1	2	0	0	2	1	1.23	0.76	0.62	
10:30	0	-3	-3	0	1	1	0	-1	1	0	-1	1	2	-3	0	1	1	-2	1	1	-1	1.09	0.94	0.87	
10:40	2	0	0	1	2	-2	-1	1	0	0	0	0	-1	0	0	0	2	0	-4	0	-1	0.77	1.14	1.47	
10:50	0	-3	2	2	-2	0	0	-1	1	0	2	1	0	0	3	2	-1	1	1	-1	2	1.18	0.92	0.78	
11:00	0	1	1	1	-1	-2	3	1	-1	4	3	1	1	2	2	0	2	1	0	3	-2	1.50	1.12	0.75	
11:10	-2	0	2	-1	3	2	2	1	1	0	1	0	1	1	-1	2	1	1	1	1	0	1.23	0.76	0.62	
11:20	1	3	-1	3	3	3	1	2	2	1	0	1	-3	1	-2	2	3	0	2	1	2	1.73	0.97	0.56	
11:30	0	0	2	1	2	1	1	-1	1	3	0	3	2	2	-2	2	-1	1	1	1	1	1.36	0.72	0.53	
11:40	3	1	2	2	0	0	3	3	1	0	1	-2	1	-1	2	1	0	0	3	1	1	1.32	1.08	0.82	
11:50	-1	1	3	2	3	0	0	-2	-1	-1	1	1	2	2	-1	2	1	2	2	-1	1	1.41	0.63	0.45	
12:00	1	3	1	0	2	3	4	0	0	2	2	2	2	2	0	-2	3	0	2	1	1	1.59	1.40	0.88	
12:10	3	2	1	-1	2	3	2	2	0	2	2	1	3	2	1	1	-1	2	-1	1	2	1.59	0.73	0.46	
12:20	2	2	3	2	2	1	3	1	2	2	3	2	3	1	1	2	0	1	-1	2	2	1.82	0.63	0.35	
12:30	2	1	0	-1	0	2	3	2	0	2	2	2	3	3	1	0	4	3	1	2	2	1.73	1.26	0.73	
12:40	0	0	0	0	2	2	3	3	0	3	2	1	0	2	3	3	2	3	0	0	2	1.41	1.68	1.19	
12:50	0	0	0	0	0	0	1	4	0	0	0	0	2	2	0	1	0	4	0	-1	4	0.86	2.03	2.35	
13:00	1	0	0	0	1	0	0	0	-1	0	0	0	0	0	0	0	0	0	0	0	0	0.14	0.12	0.90	
Daily Total	39	44	48	40	52	52	48	49	45	38	48	46	40	42	47	43	48	43	44	37	45	44	44.64	17.67	0.40
[8:50, 12:20] Total	28	32	35	29	36	35	34	29	27	31	36	30	30	27	32	26	32	30	23	29	31	25	30.32	12.61	0.42
All slot avg	2.0	2.0	2.2	1.9	2.0	2.0	2.4	2.3	1.9	2.2	2.1	1.9	2.3	2.2	2.2	2.2	2.1	2.2	2.1	2.2	2.4	1.39	1.42	1.02	
All slot var	1.5	1.9	2.2	1.9	1.8	1.5	1.8	1.3	1.5	1.7	1.5	1.5	1.6	1.5	1.3	1.7	1.8	1.6	1.6	2.2	1.8	1.6	(across all days)		
All slot var/avg	0.7	1.0	1.0</																						