

Clinic Scheduling Models with Overbooking for Patients with Heterogeneous No-show Probabilities*

Bo Zeng, Ji Lin and Mark Lawley

email: {bzeng, lin35, malawley}@purdue.edu

Weldon School of Biomedical Engineering, Purdue University, West Lafayette, IN 47907

Ayten Turkcan,

e-Enterprise Center, Purdue University, West Lafayette, IN 47907

email: aturkcan@purdue.edu

February 26, 2008

Abstract

Clinical overbooking is intended to reduce the negative impact of patient no-shows on clinic operations and performance. In this paper, we study the clinical scheduling problem with overbooking for heterogeneous patients, i.e. patients who have different no-show probabilities. We consider the objective of maximizing expected profit, which includes revenue from patients and costs associated with patient waiting times and physician overtime. We show that the objective function with homogeneous patients, i.e. patients with the same no-show probability, is multimodal. We also show that this property does not hold when patients are heterogeneous. We identify properties of an optimal schedule with heterogeneous patients and propose a local search algorithm to find local optimal schedules. Then, we extend our results to sequential scheduling and propose two sequential scheduling procedures. Finally, we perform a set of numerical experiments and provide managerial insights for health care practitioners.

Key Words clinical scheduling, overbooking, patient no-show, multimodularity

1 Introduction

The majority of patient care in the U.S. (80 – 90% [2, 4]) is provided by outpatient clinics. Clinic operations are typically driven by appointment schedules, and appointment scheduling is often cited by clinic managers as a major opportunity for improvement. Cayirli and Veral [3] provide a comprehensive review of research in outpatient appointment scheduling. They state that most analytical research does not consider factors such as patient no-shows, walk-ins and emergency. Nevertheless, these factors have significant adverse effect on operational efficiency, total revenue and patient satisfaction.

Among these factors, patient no-show is of particular concern because it wastes the available capacity of valuable resources (physician, staff, equipment) and limits clinic access to the patient population. Cayirli and Veral [3] mention that no-show rates are 5-30%. However, Rust et al. [17] report that for some health care settings, such as public pediatric clinics, no-show rates can reach 80%. To reduce the negative impact of patient no-show, clinic schedulers often overbook. However, naive overbooking can lead to longer patient waiting times, clinic overtime, and deteriorating outcomes for patients who leave without being seen [8, 18]. Thus, modeling and analysis must be used to develop a scheduling methodology that properly balances these competing objectives.

*Supported by NSF Grant CMMI 0729463

38 An ideal overbooking model depends on four characteristics. The first is a valid patient no-show
39 description that captures the real pattern of patient behavior. The second is the underlying service
40 model that reflects the operational dynamics of the clinic. The third is an objective function that reflects
41 the performance concern of clinic managers. And the last is an efficient algorithm that can generate
42 schedules of desired quality in a timely fashion. We give a brief review of existing clinic overbooking
43 models categorized according to these four characteristics. We also include some studies that do not
44 explicitly consider overbooking but that can be used to obtain overbooked schedules.

45 No-show probabilities can be correlated to factors such as reservation lead time and patient demo-
46 graphics, see Garuda et al. [5] for details. Even though no-show usually differs by patient, almost all
47 overbooking studies assume that patients are homogeneous, i.e. all patients have the same no-show
48 probability. Liu and Liu [12], Laganga and Lawrence [10, 11], Kim and Giachetti [8] and Kaandorp and
49 Koole [7] consider a single no-show rate for all patients in their models. Muthuraman and Lawley [15]
50 provide an exception by explicitly modeling different no-show probabilities.

51 With clinic dynamics, most researchers develop single server models (a schedule is created for a single
52 physician). Kim and Giachetti [8] and Laganga and Lawrence [10, 11] assume that the service times of
53 patients are deterministic while Kaandorp and Koole [7] and Muthuraman and Lawley [15] use queuing
54 based models with exponential service times. Liu and Liu [12] consider a model for multiple servers and
55 investigate service times with exponential and general distributions.

56 The performance criteria considered in appointment scheduling models includes revenue from pa-
57 tients, patient waiting time/cost, physician overtime/cost and physician idle time/cost. Kim and Gia-
58 chetti [8] consider expected revenue and physician overtime. Since the cost of patient waiting time is
59 not included in their model, all patients are assumed to arrive at the beginning of a clinic session, which
60 implies a single block scheduling model. Laganga and Lawrence [11] consider revenue from patients and
61 costs of patient waiting time and physician overtime, in both linear and quadratic objective functions.
62 Kaandorp and Koole [7] explicitly include the cost of physician idle time, as do Liu and Liu [12]. Muthu-
63 raman and Lawley [15] consider revenue from patients and costs from patients waiting time and clinic
64 overtime.

65 In most cases, computing the optimal schedule is computationally intractable and thus most schedul-
66 ing algorithms are heuristics or simulation based methods [7, 8, 10, 11, 12]. The research by Kaandorp
67 and Koole [7] is of special interest because they show that their model is multimodular. Multimodularity
68 is a property of functions in discrete space, similar to convexity in continuous space, which guarantees
69 that a locally optimal solution is also globally optimal. In contrast to the work just mentioned, Muthu-
70 raman and Lawley [15] consider sequential scheduling in which the schedule is constructed as patients
71 seeking appointments call clinic schedulers. Patients must be given their appointment time before the
72 call ends, and thus once a patient appointment is added to the schedule, it is typically not feasible to
73 alter the time. In this case, the set of patients to be scheduled is not initially known and deciding when
74 a schedule is complete becomes a problem. Although the authors provide a scheduling algorithm and
75 derive optimal stopping criteria, the optimal sequential schedule is not characterized. In this study, we
76 derive some properties of an optimal schedule, which can be used to design better algorithms.

77 The existing studies do not adequately address the question of how the scheduling problem for hetero-
78 geneous patients is different from that of homogeneous patients and whether modeling the heterogeneous
79 nature of patient no-show can lead to superior schedules, particularly in a sequential setting where sched-
80 ules have to be constructed as patients call-in. Even though different no-show probabilities are taken
81 into account in [15], the patients are treated equally while scheduling. The decision to accept (or reject)
82 a patient is given by only looking at the increase (or decrease) in the objective function. However,
83 scheduling patients with high no-show probabilities leads to higher variabilities in daily workload of
84 clinics. In this study, we also investigate the effect of variability in no-show rates on the quality of the
85 resulting schedules.

86 The remainder of the paper is structured as follows. Section 2 provides the notation and defines the
87 basic problem. In Section 3, we study the structure of the optimization model and prove that it is not
88 multimodular in general. In Section 4, after deriving some important properties of optimal schedules,
89 we propose a local search algorithm to obtain a good schedule and discuss its extension to sequential

90 scheduling settings. In Section 5, we present a computational study to compare the proposed algorithms
 91 with the existing methods. Section 6 concludes with some managerial insights that can improve patients
 92 scheduling and clinic performance in practice.

93 2 Problem Definition

94 We first state our assumptions and then present required notation. We assume that all patient arrivals to
 95 the clinic are scheduled (no walk-ins) and that the patient population can be partitioned into categories
 96 based on no-show probability. We further assume that the clinic day is partitioned into a set of time
 97 slots and that patient appointment times coincide with the beginning of a slot. We also assume that all
 98 arriving patients are punctual and that patients are served according to a first-come-first-serve protocol.
 99 Finally, we assume that service times are exponential and that they are independent and identically
 100 distributed across patients. Notation is as follows:

I	set of slots
i	slot set index, $i \in \{1, \dots, I \}$
t_i	length of slot i
J	set of patient types based on no-show
j	patient type, $j \in \{1, \dots, J \}$
X_i	number of patients arriving at start of slot i
Y_i	number of patients in the system at the end of slot i , overflow from slot i to slot $i + 1$
L_i	number of service completions in slot i
λ	service rate
c_i	unit overflow cost from slot i to slot $i + 1$ ($i \neq I $), $c_i \geq 0$
c_I	unit overflow cost for $i = I $, unit overtime cost, typically $c_{ I } > c_i$
r	revenue per patient, $r > 0$
p_j	probability that patient of type j arrives as scheduled, ($p_1 > \dots > p_{ J }$)
n_j	number of patients of type j
S	a schedule ($\in \mathbb{Z}^{ I \times J }$)
$\Delta_{i,j}$	unit matrix such that cell i, j is 1, all others 0
$F(S)$	objective function value of S
$R(S)$	overflow matrix of schedule S
$Q(S)$	arrival matrix of schedule S
$a \sim b$	random variables a and b are iid

101 For a given set of heterogeneous patients, we formulate the following overbooking model to obtain
 102 an optimal schedule S that maximizes the expected total profit.

$$\begin{aligned}
 \max F(S) &= r \sum_{i \in I} E[X_i] - \sum_{i \in I} c_i E[Y_i] \\
 \text{s.t.} \quad &\sum_{i \in I} S_{i,j} \leq n_j \\
 &S_{i,j} \in \mathbb{Z} \quad \forall i \in I, j \in J
 \end{aligned} \tag{1}$$

103 We note that $X_i + Y_{i-1}$ is the number of patients in the system at the beginning of slot i and that
 104 the number of patients in the system at the end of slot i is given as:

$$Y_i = \max\{Y_{i-1} + X_i - L_i, 0\}. \tag{2}$$

105 To compute probabilities for X_i and Y_i , Muthuraman and Lawley [15] introduce two matrices, an
 106 arrival matrix $[Q_{i,l}]$ such that $Q_{i,l}$ is the probability that l patients arrive at the beginning of slot i , and

107 an overflow matrix $[R_{i,k}]$ such that $R_{i,k}$ is the probability that k patients overflow from slot i to slot
 108 $i + 1$. These are computed as follows:

$$Q_{i,l}(S) = Pr(X_i = l) = \sum_{\pi \in \Omega} \prod_{j \in J} \frac{S_{i,j}!}{\pi_j!(S_{i,j} - \pi_j)!} p_j^{\pi_j} (1 - p_j)^{S_{i,j} - \pi_j},$$

109 where $\pi = \{\pi_1, \dots, \pi_{|J|}\}$ with $\pi_j \in \mathbb{Z}_+$ for $j \in J$, $\sum_{j \in J} \pi_j = l$ and Ω is the set of all such vectors.

$$R_{i,m}(S) = \begin{cases} \sum_l \sum_k (1 - F_{L_i}(l + k)) Q_{i,l} R_{i-1,k} & \text{if } m = 0 \\ \sum_l \sum_k f_{L_i}(l + k - m) Q_{i,l} R_{i-1,k} & \text{if } m \geq 1 \end{cases} \quad (3)$$

110

$$f_{L_i}(k) = e^{-\lambda t_i} \frac{(\lambda t_i)^k}{k!}$$

$$F_{L_i}(k) = \sum_{\tilde{k}=0}^{k-1} f_{L_i}(\tilde{k}).$$

111 Given these equations, we can compute $E[X_i] = \sum_l l Q_{i,l}$ and $E[Y_i] = \sum_k k R_{i,k}$.

112 Typically, optimization problems such as (1) arising from appointment service systems are very
 113 difficult to solve since the objective functions are nonlinear and decision variables are discrete. However,
 114 it has recently been shown that if the objective function is multimodular over \mathbb{Z}^n , a property similar to
 115 convexity in \mathbb{R}^n , and constraints are simple upper or lower bound constraints, a well-defined local search
 116 algorithm can be used to obtain (global) optimal solutions, see Hajek [6], Altman et al. [1] and Koole
 117 and van der Sluis [9]. Based on these results, Kaandorp and Koole [7] prove that their scheduling model
 118 for homogeneous patients is multimodular and implement a local search algorithm to obtain an optimal
 119 schedule. As a natural extension, it is important to see whether our overbooking model is multimodular.
 120 If so, we can use the results to obtain an optimal scheduling method, and if not we are justified in seeking
 121 heuristic approaches. Section 3 addresses this problem.

122 3 Structure of the Overbooking Scheduling Model

123 In this section, we investigate the multimodularity of the scheduling model given in (1). As an aid to
 124 the reader, we make the following informal note about multimodularity before providing the definition.
 125 Let f be function on \mathbb{Z}^m . When we join the integer points of f by lines, we obtain a new function g on
 126 \mathbb{R}^m . g is convex if and if f is multimodular. This implies that a local optimum is also a global optimum.

127 More formally, let \vec{e}_i be the i^{th} standard unit vector of \mathbb{R}^m . Then, we define a set of vectors
 128 $\Gamma = \{\vec{v}_0, \dots, \vec{v}_m\} \in \mathbb{Z}^m$ such that $\vec{v}_0 = -\vec{e}_1$, $\vec{v}_i = \vec{e}_i - \vec{e}_{i+1}$, for $i = 1, \dots, m - 1$ and $\vec{v}_m = \vec{e}_m$.

129 **Definition 1.** A function $f : \mathbb{Z}^m \rightarrow \mathbb{R}$ is multimodular if for all $\vec{x} \in \mathbb{Z}^m$, $\vec{u}, \vec{v} \in \Gamma$, $\vec{u} \neq \vec{v}$,

$$f(\vec{x} + \vec{u}) + f(\vec{x} + \vec{v}) \geq f(\vec{x}) + f(\vec{x} + \vec{u} + \vec{v}). \quad (4)$$

130 Because of the connection between multimodular and convex functions, Koole and van der Sluis [9]
 131 propose a local search algorithm that searches all the neighbors of a particular point x in the form
 132 $\vec{x} + \sum_{\vec{v} \in U} \vec{v}$ where U is a subset of Γ . They show that this local search will lead to an optimal solution
 133 of f . Later, Kaandorp and Koole [7] use this concept to obtain an optimal schedule for their scheduling
 134 model. In this section, we use the following equivalent form

$$f(\vec{x} + \vec{u}) - f(\vec{x}) \geq f(\vec{x} + \vec{v} + \vec{u}) - f(\vec{x} + \vec{v}) \quad (5)$$

135 to verify the multimodularity of a function. We note that (5) can be interpreted as the improvement
 136 from perturbing \vec{x} by \vec{u} is greater or equal to that from perturbing $\vec{x} + \vec{v}$ by \vec{u} . Because \vec{u} and \vec{v} are

137 closely related to the unit vectors \vec{e}_i for some i , we first derive some result on the improvement of f
 138 obtained from perturbing \vec{x} by \vec{e}_i . This result will be frequently used in this section to help us simplify
 139 the proof of multimodularity.

140 **Proposition 1.** *For a given schedule S^0 , we have*

$$\frac{F(S^0 + \Delta_{i^*, j_1}) - F(S^0)}{F(S^0 + \Delta_{i^*, j_2}) - F(S^0)} = \frac{p_{j_1}}{p_{j_2}} \quad (6)$$

141 for all $i^* \in I$ and $j_1, j_2 \in J$.

142 *Proof.* Assume that W is a patient with no-show probability p_{j_1} being added to slot i^* in schedule S^0
 143 such that the schedule is updated by $S^1 = S^0 + \Delta_{i^*, j_1}$. We use X_i^0 and Y_i^0 to denote the number of
 144 arrivals in slot i and the size of overflow from slot i , respectively, for schedule S^0 . Then, we define X_i^1
 145 and Y_i^1 for S^1 similarly. Also, we introduce $P_i(i^*)$ to be the conditional probability that the arrival of
 146 patient W increases the overflow from slot i to $i + 1$ by 1.

147 Let \mathcal{W} denote the arrival of patient W . Then, from (1), on the condition of \mathcal{W} , we have

$$\begin{aligned} F(S^1) - F(S^0) &= rp_{j_1} + \sum_{i \in I} c_i E[Y_i^1 - Y_i^0] \\ &= rp_{j_1} - (1 - p_{j_1}) \sum_{i \in I} 0 - p_{j_1} \sum_{i \in I} c_i E[Y_i^1 - Y_i^0 | \mathcal{W}] \\ &= p_{j_1} (r - \sum_{i \in I} c_i P_i(i^*)). \end{aligned} \quad (7)$$

148 It can be easily seen that if $P_i(i^*)$ is independent of p_{j_1} for all i , we have $F(S^0 + \Delta_{i^*, j_2}) - F(S^0) =$
 149 $p_{j_2} (r - \sum_{i \in I} c_i P_i(i^*))$. Then, the conclusion follows.

150 In fact, we observe that the only situation where the arrival of patient W leads to one more patient
 151 overflowing from slot i is the case where for each slot k such that $i^* \leq k \leq i$, the number of patients
 152 served is less than or equal to the number patients in slot k in schedule S^0 . So, we have

$$P_i(i^*) = \begin{cases} \prod_{k=i^*}^i Pr(L_k \leq X_k^0 + Y_{k-1}^0) & \text{if } i^* \leq i \\ 0 & \text{otherwise.} \end{cases} \quad (8)$$

153 Clearly, both $P_i(i^*)$ and $r - \sum_{i \in I} c_i P_i(i^*)$ are independent of the no-show probability of patient W . The
 154 desired results follows. \square \square

155 Because the result of Proposition 1 is about the value difference from unit changes, we call the result
 156 of (7) the *local perturbation* of a given schedule. Next, we show that the concept of local perturbation can
 157 help us simplify the proof of multimodularity significantly as compared to that of Kaandorp and Koole
 158 [7]. Because (1) is a maximization problem, we use \mathfrak{F} to denote $-F$ and verify that \mathfrak{F} is multimodular.
 159 When $|J| = 1$, we use S_i instead of $S_{i,j}$ and use p to denote the patient show-up probability in our
 160 derivation.

161 **Theorem 2.** *When $|J| = 1$, \mathfrak{F} is a multimodular function over $\mathbb{Z}^{|I|}$.*

162 *Proof.* Because for the simple cases where \vec{u} or \vec{v} is $-\vec{e}_1$ or $\vec{e}_{|I|}$, (5) can easily be proven using the
 163 argument similar to the following, we focus on the general case where neither \vec{u} and \vec{v} are standard unit
 164 vectors.

165 Without loss of generality, we let $\vec{u} = \vec{e}_k - \vec{e}_{k+1}$ and $\vec{v} = \vec{e}_l - \vec{e}_{l+1}$ such that $1 \leq k < l \leq |I| - 1$. To
 166 make use of our local perturbation, for any particular schedule S , we define two base schedules S^L and
 167 S^R for LHS and RHS of (5) as

$$S^L = S - \vec{e}_{k+1} \quad (9)$$

168 and

$$S^R = S - \vec{e}_{k+1} + \vec{e}_l - \vec{e}_{l+1} \quad (10)$$

169 with $S_{k+1} \geq 1$ and $S_{l+1} \geq 1$.

170 By using S^L and S^R , both LHS and RHS can be interpreted as the difference of two local pertur-
 171 bations. Correspondingly, we use $X_i^L, X_i^R, Y_i^L, Y_i^R$ to denote the number of arrivals in slot i and the
 172 number of overflows from slot i in S^L and S^R , respectively. We also use $P_i^L(i_0)$ and $P_i^R(i_0)$ to denote
 173 the overflow effect from adding one more patient to slot i_0 in schedule S^L and S^R on the condition of
 174 this patient's arrival.

175 From (8), we have (11) for LHS of (5)

$$\begin{aligned}
 \mathfrak{F}(S + \vec{u}) - \mathfrak{F}(S) &= (\mathfrak{F}(S + \vec{u}) - \mathfrak{F}(S^L)) - (\mathfrak{F}(S) - \mathfrak{F}(S^L)) \\
 &= p \sum_{i=k}^{|I|} c_i P_i^L(k) - p \sum_{i=k+1}^{|I|} c_i P_i^L(k+1) \\
 &= p \{c_k P_k^L(k) + Pr(L_k \leq X_k^L + Y_{k-1}^L) \sum_{i=k+1}^{|I|} c_i \prod_{h=k+1}^i Pr(L_h \leq x_h^L + Y_{h-1}^L) - \\
 &\quad \sum_{i=k+1}^{|I|} c_i \prod_{h=k+1}^i Pr(L_h \leq x_h^L + Y_{h-1}^L)\} \tag{11} \\
 &= p \{c_k P_k^L(k) + (Pr(L_k \leq X_k^L + Y_{k-1}^L) - 1) \sum_{i=k+1}^{|I|} c_i P_i^L(k+1)\}.
 \end{aligned}$$

176 Note that, in the first equality of (11), the first sum evaluates overflow characteristics when the patient
 177 is added to slot k and the second sum evaluates overflow characteristics when the patient is added to
 178 slot $k+1$. Similarly, we have (12) for RHS of (5).

$$\begin{aligned}
 \mathfrak{F}(S + \vec{u} + \vec{v}) - \mathfrak{F}(S + \vec{v}) &= (\mathfrak{F}(S + \vec{u}) - \mathfrak{F}(S^R)) - (\mathfrak{F}(S) - \mathfrak{F}(S^R)) \\
 &= p \sum_{i=k}^{|I|} c_i P_i^R(k) - p \sum_{i=k+1}^{|I|} c_i P_i^R(k+1) \\
 &= p \{c_k P_k^R(k) + (Pr(L_k \leq X_k^R + Y_{k-1}^R) - 1) \sum_{i=k+1}^{|I|} c_i P_i^R(k+1)\}. \tag{12}
 \end{aligned}$$

179 Because S^L and S^R have same number of patients per slot up to and including slot $l-1 \geq k$, we
 180 have $P_k^L(k) = P_k^R(k) = Pr(L_k \leq X_k^L + Y_{k-1}^L) = Pr(L_k \leq X_k^R + Y_{k-1}^R)$. Also, because $Pr(L_k \leq$
 181 $X_k^L + Y_{k-1}^L) - 1 \leq 0$, it is sufficient to show that

$$\sum_{i=k+1}^{|I|} c_i P_i^L(k+1) \leq \sum_{i=k+1}^{|I|} c_i P_i^R(k+1).$$

182 Observe that S^R can be obtained from S^L by reassigning one patient who is in slot $l+1$ to slot l . Let
 183 W be a such patient. Then, we can compare (11) and (12) conditioned on the arrival of W , \mathcal{W} . Clearly,
 184 if $\neg \mathcal{W}$, we have (11) and (12) are same. If \mathcal{W} , we have $X_i^R \sim X_i^L + 1$, $X_{l+1}^R \sim X_l^L - 1$, $Y_i^R \sim Y_i^L$ for
 185 $i = k, \dots, l-1$ and $X_i^R \sim X_i^L$ for $i \neq l, l+1$. Furthermore, on the condition of \mathcal{W} , we also have

$$Pr(L_l \leq X_l^R + Y_{l-1}^R) = Pr(L_l \leq X_l^L + Y_{l-1}^L + 1) = Pr(L_l \leq X_l^L + Y_{l-1}^L) + Pr(L_l = X_l^L + Y_{l-1}^L + 1).$$

186 Next, we compare the dynamics of queuing model in S^L and S^R in the case where $L_l \leq X_l^L + Y_{l-1}^L$.
 187 Because $X_l^R \sim X_l^L + 1$, $Y_{l-1}^R \sim Y_{l-1}^L$ and $L_l \leq X_l^L + Y_{l-1}^L$, we have $Y_l^R \sim Y_l^L + 1$. Also, because
 188 $X_{l+1}^R \sim X_{l+1}^L - 1$, we have $X_{l+1}^R + Y_{l+1}^R \sim X_{l+1}^L + Y_{l+1}^L$ and therefore $Y_{l+1}^R \sim Y_{l+1}^L$. From the fact that
 189 $X_j^R \sim X_j^L$ for $j \geq l+1$, we have $X_j^R + Y_{j-1}^R \sim X_j^L + Y_{j-1}^L$ for $j = l+1, \dots, |I|$. Clearly, from slot $l+1$
 190 to slot $|I|$, the queuing dynamics in schedule S^R and S^L are identical, as shown in Figure 1.

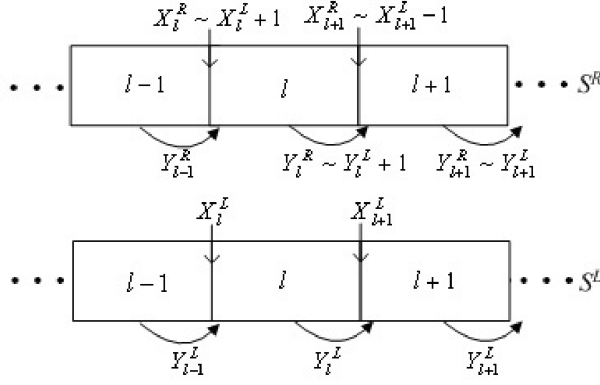


Figure 1: Dynamics of Slot Queuing System in S^R and S^L Conditioned on \mathcal{W} and $L_l \leq X_l^L + Y_{l-1}^L$

191 Let $H_i^{R*} = \prod_{j=l+1}^i Pr(L_j \leq X_j^R + Y_{j-1}^R)$ given $L_l \leq X_l^L + Y_{l-1}^L$ and $H_i^{R'} = \prod_{j=l+1}^i Pr(L_j \leq$
 192 $X_j^R + Y_{j-1}^R)$ given $L_l = X_l^L + Y_{l-1}^L + 1$ for $i \geq l+1$. H_i^{R*} represents the probability that \mathcal{W} causes
 193 additional overflow from slot i and $H_i^{R'}$ represents the probability that \mathcal{W} causes no additional overflow
 194 flow from slot i . Further, $H_i^{R*} Pr(L_l \leq X_l^L + Y_{l-1}^L) = \prod_{j=l}^i Pr(L_j \leq X_j^L + Y_{j-1}^L)$ for $i \geq l+1$. As a
 195 consequence, we obtain

$$\begin{aligned}
 \sum_{i=k+1}^{|I|} c_i P_i^R(k+1) &= \sum_{i=k+1}^{l-1} c_i P_i^L(k+1) + c_l P_{l-1}^L(k+1) (Pr(L_l \leq X_l^L + Y_{l-1}^L) + Pr(L_l = X_l^L + Y_{l-1}^L + 1)) \\
 &\quad + P_{l-1}^L(k+1) \left(\sum_{i=l+1}^{|I|} c_i H_i^{R*} Pr(L_l \leq X_l^L + Y_{l-1}^L) + \sum_{i=l+1}^{|I|} c_i H_i^{R'} Pr(L_l = X_l^L + Y_{l-1}^L + 1) \right) \\
 &\geq \sum_{i=k+1}^{l-1} c_i P_i^L(k+1) + c_l P_{l-1}^L(k+1) Pr(L_l \leq X_l^L + Y_{l-1}^L) + \\
 &\quad P_{l-1}^L(k+1) Pr(L_l \leq X_l^L + Y_{l-1}^L) \sum_{i=l+1}^{|I|} c_i H_i^{R*} \\
 &= \sum_{i=k+1}^{|I|} c_i P_i^L(k+1).
 \end{aligned} \tag{13}$$

196 The inequality follows from the fact that $c_i \geq 0$ for $i \in I$ and probabilities are non-zero. The last equality
 197 follows from the definition of $P_i^L(k)$ in (8). \square \square

198 Since (1) is multimodular for $|J| = 1$, we can apply the local search method by Koole and van der
 199 Sluis [9] to obtain an optimal schedule. However, for the general case where $|J| \geq 2$, we have the following
 200 result. Note that we express a schedule S in the form of a vector, $[S_{1,1}, \dots, S_{1,|J|}, \dots, S_{|I|,1}, \dots, S_{|I|,|J|}] \in$
 201 $\mathbb{Z}^{|I||J|}$.

202 **Theorem 3.** *The function \mathfrak{F} is not multimodular over $\mathbb{Z}^{|I||J|}$ for $|J| \geq 2$.*

203 *Proof.* It is sufficient to show that for some $\vec{u}, \vec{v} \in \Gamma$, (5) does not hold. Let $\vec{u} = \vec{e}_l - \vec{e}_{l+1}$ and
 204 $\vec{v} = \vec{e}_{l+k|J|} - \vec{e}_{l+k|J|+1}$ for some l, k such that $1 \leq l, l+k|J|+1 \leq |I||J|$. Denote $j_0 = \lfloor \frac{l}{|J|} \rfloor + 1$. Then,
 205 we observe that the operation corresponding to \vec{u} (\vec{v} , respectively) is to move one patient of p_{j_0} (p_{j_0+k} ,
 206 respectively) from slot $i_0 + 1 = l + 1 - (j_0 - 1)|J|$ and to slot i_0 .

207 Similar to our proof for Theorem 2, for a particular schedule S , we define two base schedule S^L and
 208 S^R for LHS and RHS of (5) as

$$S^L = S - \vec{e}_{l+1} \quad (14)$$

209 and

$$S^R = S - \vec{e}_{l+1} + \vec{e}_{l+k|J|} - \vec{e}_{l+k|J|+1}. \quad (15)$$

210 We also use X_i^L, X_i^R, Y_i^L , and Y_i^R to denote arrival and overflow in S^L and S^R respectively. Comparing
 211 S^L and S^R , we observe that $S_{i_0, j_0+k}^R = S_{i_0, j_0+k}^L + 1$, $S_{i_0+1, j_0+k}^R = S_{i_0+1, j_0+k}^L - 1$ and $S_{i, j}^R = S_{i, j}^L$ for all
 212 other (i, j) . We can easily see that S^R the scheduling resulting when we reassign one patient of type
 213 $j_0 + k$ from slot $i_0 + 1$ to slot i_0 in schedule S^L . Let W be a such patient.

214 Again, we use $P_i^L(i_0)$ and $P_i^R(i_0)$ to denote the overflow effect from adding one more patient to slot
 215 i_0 in schedule S^L and S^R conditioned on \mathcal{W} . If $\neg\mathcal{W}$, $S^R = S^L$. So, we need only consider the case where
 216 W shows up.

217 Similar to the proof of Theorem 2, for LHS of (5), we have

$$\begin{aligned} \text{LHS} = \mathfrak{F}(S + \vec{u}) - \mathfrak{F}(S) = & c_{i_0} Pr(L_{i_0} \leq X_{i_0}^L + Y_{i_0-1}^L) \\ & + (Pr(L_{i_0} \leq X_{i_0}^L + Y_{i_0-1}^L) - 1) \left(\sum_{i=i_0+1}^{|I|} c_i \prod_{k=i_0+1}^i Pr(L_k \leq X_k^L + Y_{k-1}^L) \right) \end{aligned} \quad (16)$$

218 For RHS of (5), we have

$$\begin{aligned} \text{RHS} = \mathfrak{F}(S + \vec{u} + \vec{v}) - \mathfrak{F}(S + \vec{v}) = & c_{i_0} Pr(L_{i_0} \leq X_{i_0}^R + Y_{i_0-1}^R) \\ & + (Pr(L_{i_0} \leq X_{i_0}^R + Y_{i_0-1}^R) - 1) \left(\sum_{i=i_0+1}^{|I|} c_i \prod_{k=i_0+1}^i Pr(L_k \leq X_k^R + Y_{k-1}^R) \right) \end{aligned} \quad (17)$$

219 Next, we compare the value of (16) and (17) conditioned on the physician's performance in slot i_0 ,
 220 i.e. L_{i_0} . Note that under \mathcal{W} , we have $X_{i_0}^R + Y_{i_0}^R \sim X_{i_0}^L + Y_{i_0}^L + 1$ and $X_{i_0+1}^R \sim X_{i_0+1}^L - 1$.

221 **Case (i)** $L_{i_0} = X_{i_0}^L + Y_{i_0-1}^L + 1$

222 Because $Pr(L_{i_0} \leq X_{i_0}^R + Y_{i_0-1}^R) = 1$, (17) is equal to c_{i_0} . However, (16) is equal to
 223 $-\sum_{i=i_0+1}^{|I|} c_i \prod_{k=i_0+1}^i Pr(L_k \leq X_k^L + Y_{k-1}^L)$. So, LHS - RHS < 0 because $c_{|I|} > c_{i_0+1} \geq 0$.

224 **Case (ii)** $L_{i_0} \leq X_{i_0}^L + Y_{i_0-1}^L - 1$

225 For this case, we have LHS = RHS = c_{i_0} .

226 Because the probability of both cases is nonzero, we conclude that (5) does not hold when $|J| \geq 2$. $\square \square$

227 Since the multimodular property does not hold for the general case, we do not expect to obtain
 228 an optimal schedule without implementing an exhaustive search. These observations motivate us to
 229 develop a local search algorithm that is efficient and can be used to obtain schedules with good quality.
 230 We present our study on the solution methodology in Section 4.

231 4 Local Search Algorithm and Sequential Heuristics for Clinical 232 Scheduling

233 Because our scheduling model for heterogeneous patients is not multimodular as shown in Section 3, we
 234 first propose a local search algorithm to find good schedules in Section 4.1. The main assumption of
 235 the proposed local search algorithm is that the set of patients is known in advance. In many situations,
 236 clinics do not know the set of patients that should be scheduled and appointment schedules are generated
 237 sequentially along with the patient call-in process. So, in Section 4.2, we extend our study to sequential
 238 scheduling and propose two sequential scheduling procedures.

4.1 Local Search Algorithm

First, we derive an important property of optimal schedules, that can be used to generate initial schedules. Then, we define the neighborhood of a given schedule and propose dominance rules to reduce the search space in a local search algorithm. Finally, we give the basic steps of the proposed algorithm.

Theorem 4 shows that an optimal schedule prefers patients with lower no-show probabilities.

Theorem 4. *Let S^* be an optimal schedule of (1). Let $j, j_0 \in J$ with $p_j > p_{j_0}$ and $n_j, n_{j_0} > 0$. If $\sum_{i \in I} S_{i, j_0}^* \geq 1$, then $\sum_{i \in I} S_{i, j}^* = n_j$.*

Proof. Let $j, j_0 \in J$ with $p_j > p_{j_0}$ and $n_j, n_{j_0} > 0$. Let S^* be an optimal schedule such that for some $i_0 \in I$, $S_{i_0, j_0}^* > 0$, and suppose $\sum_{i \in I} S_{i, j}^* < n_j$.

Let $\tilde{S} = S^* - \Delta_{i_0, j_0}$. Then, $F(\tilde{S}) \leq F(S^*)$. From the proof of Proposition 1, we have $p_{j_0} r \geq p_{j_0} (\sum_{i \in I, i \geq i_0} c_i P_i(i_0))$ where $P_i(i_0)$ is the probability of overflow from slot i incurred by adding one patient in slot i_0 to \tilde{S} on the condition of this patient's arrival. Consider the schedule $\hat{S} = \tilde{S} + \Delta_{i_0, j}$. Since $p_{j_0} r \geq p_{j_0} (\sum_{i \in I, i \geq i_0} c_i P_i(i_0))$, we have $r \geq (\sum_{i \in I, i \geq i_0} c_i P_i(i_0))$, and thus $p_j r \geq p_j (\sum_{i \in I, i \geq i_0} c_i P_i(i_0))$. From this, we get $p_j (r - (\sum_{i \in I, i \geq i_0} c_i P_i(i_0))) > p_{j_0} (r - \sum_{i \in I, i \geq i_0} c_i P_i(i_0))$. Thus, $F(\hat{S}) > F(S^*)$, a contradiction. \square

Theorem 4 shows that patients with lower no-show probabilities contribute more than those with higher no-show probabilities. We propose a local search algorithm that schedules patients according to their no-show probabilities. Before explaining the local search algorithm in detail, we define the neighborhood of a given schedule.

Definition 2. *We say schedule S^1 is a neighbor of schedule S^0 if it satisfies the following:*

- (1) $S^1 = S^0 \pm \Delta_{i_0, j_0}$ for some $i_0 \in I$ and $j_0 \in J$, i.e. S^1 is obtained by adding/removing one patient of type j_0 ;
- (2) $S^1 = S^0 - \Delta_{i_0, j_0} + \Delta_{i_1, j_0}$ where $i_0 \neq i_1$, i.e. S^1 is obtained by reassign one patient of type j_0 from slot i_0 to slot i_1 ;
- (3) $S^1 = S^0 - \Delta_{i_0, j_0} + \Delta_{i_1, j_0} - \Delta_{i_1, j_1} + \Delta_{i_0, j_1}$ where $i_0 \neq i_1$ and $j_0 \neq j_1$, i.e. S^1 is obtained by switching the slots for two patients of different no-show probabilities.

Note that the size of the neighborhood is $O(\max\{|I||J|, n\}^2)$ where n is the number of patients in the schedule. However, in the course of local search, the size of the effective neighborhood can further be reduced as follows:

Proposition 5. *Let S^0 be a given schedule that is feasible to (1).*

- (i) *Assume that $S^k = S^0 - \Delta_{i_0, j_k} + \Delta_{i_1, j_k}$ for $k = 1, 2$ and $p_{j_2} > p_{j_1}$. If $F(S^1) > F(S^0)$, then $F(S^2) > F(S^1) > F(S^0)$;*
- (ii) *Assume that $S^k = S^0 - \Delta_{i_0, j_0} + \Delta_{i_1, j_0} - \Delta_{i_1, j_k} + \Delta_{i_0, j_k}$ for $k = 1, 2$, and $i_0 < i_1$. If $F(S^1) > F(S^0)$ and $p_{j_2} > p_{j_1} > p_{j_0}$, then $F(S^2) > F(S^1) > F(S^0)$; if $F(S^1) > F(S^0)$ and $p_{j_2} < p_{j_1} < p_{j_0}$, then $F(S^2) > F(S^1) > F(S^0)$.*

Proof. It is very straightforward to show (i) using the results of Proposition 1 and Theorem 4. Here, we focus on the more difficult proof of (ii).

Let x and y be two patients of types j_0 and j_1 , respectively. Let $S|A\overline{B}$ denote schedule S on the condition that all patients in A arrive as scheduled and all patients in B are no-shows. The expected profit of schedules S^1 and S^2 are calculated by conditioning the no-show scenarios of patients.

$$F(S^0) = (1 - p_{j_0})(1 - p_{j_1})F(S^0|\overline{xy}) + p_{j_0}(1 - p_{j_1})F(S^0|x\overline{y}) + (1 - p_{j_0})p_{j_1}F(S^0|\overline{xy}) + p_{j_0}p_{j_1}F(S^0|xy)$$

$$F(S^1) = (1 - p_{j_0})(1 - p_{j_1})F(S^1|\overline{xy}) + (1 - p_{j_1})p_{j_0}F(S^1|x\overline{y}) + (1 - p_{j_0})p_{j_1}F(S^1|\overline{xy}) + p_{j_0}p_{j_1}F(S^1|xy)$$

Note that $(S^0|\overline{xy}) = (S^1|\overline{xy})$ and $(S^1|xy) = (S^0|xy)$. Let $S^* = S^0 - \Delta_{i_0, j_0} - \Delta_{i_1, j_1}$ and $P_i^*(k)$ be the overflow probability from slot k defined in (8). It can be easily seen that $(S^0|\overline{xy}) = (S^* + \Delta_{i_1, j_1}|y)$ and $(S^0|x\overline{y}) = (S^* + \Delta_{i_0, j_0}|x)$. Similar result holds for S^1 . Then, we have the following:

$$\begin{aligned}
F(S^1) - F(S^0) &= p_{j_0}(1 - p_{j_1})F(S^* + \Delta_{i_1, j_0}|x) + (1 - p_{j_0})p_{j_1}F(S^* + \Delta_{i_0, j_1}|y) \\
&\quad - \{p_{j_0}(1 - p_{j_1})F(S^* + \Delta_{i_0, j_0}|x) + (1 - p_{i_0})p_{j_1}F(S^* + \Delta_{i_1, j_1}|y)\} \\
&= p_{j_0}(1 - p_{j_1})\{F(S^* + \Delta_{i_1, j_0}|x) - F(S^* + \Delta_{i_0, j_0}|x)\} \\
&\quad + (1 - p_{j_0})p_{j_1}\{F(S^* + \Delta_{i_0, j_1}|y) - F(S^* + \Delta_{i_1, j_1}|y)\} \\
&= p_{j_0}(1 - p_{j_1})\left\{\sum_{i \in I} c_i P_i^*(i_0) - \sum_{i \in I} c_i P_i^*(i_1)\right\} \\
&\quad + (1 - p_{j_0})p_{j_1}\left\{\sum_{i \in I} c_i P_i^*(i_1) - \sum_{i \in I} c_i P_i^*(i_0)\right\} \\
&= \left\{\sum_{i \in I} c_i \{P_i^*(i_0) - P_i^*(i_1)\}\right\} \left\{p_{j_0}(1 - p_{j_1}) - (1 - p_{j_0})p_{j_1}\right\}
\end{aligned}$$

280 For the case where $p_{j_0} < p_{j_1}$, we observe that the second term in the last equality is always negative
281 because $p_{j_0} < p_{j_1}$ and $1 - p_{j_1} < 1 - p_{j_0}$. Since $F(S^1) - F(S^0) > 0$, the first term of last equality should be
282 negative. The first term is independent of no-show probabilities of x, y and $p_{j_0}(1 - p_{j_2}) - (1 - p_{j_0})p_{j_2} <$
283 $p_{j_0}(1 - p_{j_1}) - (1 - p_{j_0})p_{j_1}$. Therefore, $F(S^2) > F(S^1) > F(S^0)$.

284 Similarly, we can prove the desired results for the case where $p_{j_0} > p_{j_1}$. □ □

285 We observe that the results of Proposition 5 provide guidelines such that local movements can be
286 implemented according to patient no-show probabilities. In particular, using Theorem 4 and Proposition
287 5, we can obtain better schedules with reduced computational effort. Next, we describe our local search
288 algorithm in detail.

289 For a given schedule S , we define $\bar{j}_i = \arg \max\{p_j : S_{i,j} \geq 1\}$, $\underline{j}_i = \arg \min\{p_j : S_{i,j} \geq 1\}$ and j^* as
290 the patient type with the lowest no-show probability to be scheduled. From Definition 2, we note that
291 there are four types of neighbors of S obtained by the following local movements: add, remove, reassign
292 and switch, which are numbered by 1, ..., 4 respectively.

293 **Algorithm 1.**

294 (1) *Initialization:* $S = \emptyset$.

295 (2) *Local Search:*

296 For $l = 1$ to 4

297 **if** $l = 1$: (neighbors obtained by adding) $F_1^* = \max\{F(S + \Delta_{i, j^*}) : i \in I\}$;

298 **if** $l = 2$: (neighbors obtained by removing) $F_2^* = \max\{F(S - \Delta_{i, \bar{j}_i}) : i \in I\}$;

299 **if** $l = 3$: (neighbors obtained by reassigning) $F_3^* = \max\{F(S - \Delta_{i, \bar{j}_i} + \Delta_{k, \bar{j}_i}) : i, k \in I\}$;

300 **if** $l = 4$: (neighbors obtained by switching) $F_4^* = \max\{F_4^1, F_4^2\}$ where

$$F_4^1 = \max\{F(S - \Delta_{i, \bar{j}_i} + \Delta_{k, \bar{j}_i} - \Delta_{k, \underline{j}_k} + \Delta_{i, \underline{j}_k}) : 1 \leq i < k \leq |I|, p_{\bar{j}_i} > p_{\underline{j}_k}\};$$

301 and

$$F_4^2 = \max\{F(S - \Delta_{i, \underline{j}_i} + \Delta_{k, \underline{j}_i} - \Delta_{k, \bar{j}_k} + \Delta_{i, \bar{j}_k}) : 1 \leq i < k \leq |I|, p_{\bar{j}_i} < p_{\underline{j}_k}\}.$$

302

303 end for

304 (3) Find the best schedule, S^* , in the neighborhood of S , i.e. $S^* = \arg \max\{F_l^* : l = 1, \dots, 4\}$.

- 305 (4) If $F(S^*) \geq F(S)$, update current schedule $S = S^*$ and go back to Step (2). Otherwise, go to Step
306 (5).
- 307 (5) Return the local optimal schedule S .

308 In Algorithm 1, the number of neighbors search for a schedule need to be searched is $O(\max\{|I|^2, n\})$,
309 which is smaller than the actual neighborhood size.

310 4.2 Sequential Scheduling Methods

311 As mentioned earlier, many clinics generate schedules in a sequential fashion. Typically, a patient calls
312 requesting an appointment. The scheduler will either add the patient to an existing schedule and give
313 an appointment time or reject the patient. Muthuraman and Lawley [15] propose a myopic scheduling
314 method, which sequentially assigns calling patients to the slot that most increases the expected profit
315 of the resulting schedule. It is called myopic since it does not take the possibility of future call-ins
316 into account when making the current assignment. Patients are added to a schedule until the expected
317 total profit starts decreasing. Although the authors consider heterogeneous patients, their method does
318 not differentiate patients according to their no-show probabilities while generating schedules. However,
319 better schedules can be generated by considering the no-show probabilities of patients. We propose two
320 sequential scheduling algorithms that do this by using the properties explained in Section 4.1.

321 Let S^0 be a fixed schedule for $n - 1$ patients and assume that we need to schedule the n^{th} patient
322 of type j . The patient will be inserted into the schedule if adding this patient increases the objective
323 function value. Corollary 6 shows that the decision of accepting (or rejecting) a patient is independent
324 of the patient's no-show probability.

325 **Corollary 6.** *For the myopic scheduling method in [15], the decision to accept (or reject) the n^{th} patient
326 of type j is independent of p_j .*

327 *Proof.* The myopic scheduling method in [15] is as follows.

328 **Step 1.** Set $S_{i,j}^0 = 0$ for all $i \in I$ and $j \in J$.

329 **Step 2.** Wait for k^{th} patient call.

330 **Step 3.** The k^{th} patient calls in and the patient is of type j_0 .

331 **Step 4.** Compute $F(S^0 + \Delta_{i,j_0})$ for $i \in I$ and set $i^* = \arg \max\{F(S^0 + \Delta_{i,j_0}) : i \in I\}$.

332 **Step 5.** If $F(S^0 + \Delta_{i^*,j_0}) > F(S^0)$, assign the k^{th} patient to slot i^* and update $S^0 = S^0 + \Delta_{i^*,j_0}$ and
333 go to Step 2. Otherwise, stop.

334 Assume that the n^{th} patient is assigned to slot i_n . Let $S = S^0 + \Delta_{i_n,j}$. From Proposition 1, we have
335 $F(S) - F(S^0) = p_j(r - \sum_{i \in I} c_i P_i(i_n))$ where $P_i(i_n)$ is independent of p_j and can be computed without
336 the n^{th} patient. Then, let i_n^* denote the slot index that yields the minimal $\sum_{i \in I} c_i P_i(i_n)$. It is clear that
337 if $\sum_{i \in I} c_i P_i(i_n^*) \leq r$, we can assign patient n to slot i_n^* to increase the expected total profit. Otherwise,
338 patient n will be rejected. Therefore, the decision on the n^{th} patient is independent of p_j . \square \square

339 Based on the results in Proposition 1 and Corollary 6, it is anticipated that overbooking many
340 patients with high no-show probabilities cannot provide the most desirable results. One major drawback
341 of the myopic scheduling method Muthuraman and Lawley [15] is that the objective function depends
342 on the call-in sequence. Different call-in sequences generate schedules which have high variability in the
343 objective function. In order to show the effect of call-in sequences, we consider two sequences. In the
344 first sequence, there are sufficiently many patients of type j_1 before any patient of type j_2 ($p_{j_1} > p_{j_2}$).
345 In the second sequence, there are sufficiently many patients of type j_2 before any patient of type j_1 . We
346 apply the myopic scheduling algorithm [15] to generate schedules S_1 and S_2 , respectively. Clearly, S_1
347 has patients of type j_1 and S_2 has patients of type j_2 . Figure 2 shows $\frac{F(S_1)}{F(S_2)}$ as a function of $p_{j_1} - p_{j_2}$ for

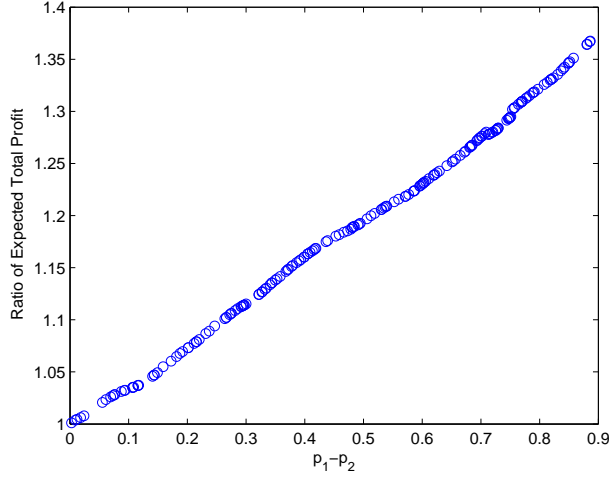


Figure 2: Ratio of Expected Total Profit $F(S_1)/F(S_2)$ vs. $p_{j_1} - p_{j_2}$

348 200 pairs of randomly generated call-in sequences. The difference between $F(S_1)$ and $F(S_2)$ increases
 349 as $p_2 - p_1$ increases.

350 One may think that it is rare to have all patients at the beginning of the sequence with high no-show
 351 probabilities. However, it is commonly observed that patients who make reservations at earlier times
 352 tend to have higher no-show probabilities and they often use up the physician's capacity before patients
 353 with low no-show probabilities can be scheduled [13, 14, 16]. Therefore, a sequential scheduling method,
 354 which accepts all patients regardless of their no-show probabilities, may not generate good schedules in
 355 terms of expected profit.

356 In sequential scheduling, there is an existing schedule and the patients who will call-in should be
 357 scheduled or rejected. Assume that S^0 is the existing schedule and there are \bar{n}_j patients for all $j \in J$
 358 that should be scheduled or rejected. The following is the revised model of (1) in the sequential scheduling
 359 setting, which adds new patients to an existing schedule:

$$\begin{aligned}
 \max \quad & G(S) = r \sum_{i \in I} \sum_{j \in J} S_{i,j} p_j - \sum_{i \in I} c_i \sum_k k R_{i,k} \\
 \text{s.t.} \quad & \sum_{i \in I} S_{i,j} - \sum_{i \in I} S_{i,j}^0 \leq \bar{n}_j \\
 & S_{i,j} - S_{i,j}^0 \geq 0 \\
 & S_{i,j} \in \mathbb{Z} \forall i \in I, j \in J
 \end{aligned} \tag{18}$$

360 We give the following result as a corollary of Theorem 4 which can be easily proven using similar argument
 361 to that of Theorem 4.

362 **Corollary 7.** Assume that S^* is an optimal schedule to (18). For $j, j_0 \in J$ with $p_j > p_{j_0}$, either
 363 $\bar{n}_j = \sum_{i \in I} (S_{i,j}^* - S_{i,j}^0)$ or $\sum_{i \in I} (S_{i,j_0}^* - S_{i,j_0}^0) = 0$, i.e. $S_{i,j_0}^* = S_{i,j_0}^0$ for all $i \in I$.

364 Corollary 7 shows that it is always better to include patients of low no-show probabilities into an
 365 existing schedule before capacity limit is reached. Also, by Proposition 1 and Corollary 7, our local
 366 search algorithm still works for any given existing schedule and patient set. From Corollary 7 and
 367 the observation in Figure 2, it is easy to see that limiting the number of patients with high no-show
 368 probabilities in the schedule will be an effective way to improve its performance. Following this line, we
 369 propose two sequential scheduling methods: the restricted myopic scheduling and the forecasting-based
 370 scheduling methods.

371 The basic idea of the restricted myopic scheduling method is using upper bounds to restrict the
372 number of patients with high no-show probabilities in the schedule. We set upper bounds, \overline{B}_j , on the
373 number of patients of type j for $j \neq 1$ in the schedule such that $\overline{B}_1 \geq \overline{B}_2 \geq \dots \geq \overline{B}_{|J|}$. Let b_j be the
374 number of patients of type j in the current schedule. The basic steps of the restricted myopic scheduling
375 method are as follows:

376 **Restricted Myopic Sequential Scheduling Method 1.**

377 **Step 1.** Set $b_j = 0$ and $S_{i,j}^0 = 0$ for all $i \in I$ and $j \in J$.

378 **Step 2.** Wait for k^{th} patient call.

379 **Step 3.** The k^{th} patient call occurs and is of type j_0 .

380 **Step 4.** If $b_{j_0} + 1 > \overline{B}_{j_0}$, do not accept patient k and go to Step 2.

381 **Step 5.** Perform the traditional myopic scheduling algorithm to compute the best slot i_0 for patient k .

382 **Step 6.** If $F(S^0 + \Delta_{i_0, j_0}) < F(S^0)$, i.e. adding patient k decreases the expected total profit, stop.
383 Otherwise, update $b_{j_0} = b_{j_0} + 1$ and $S^0 = S^0 + \Delta_{i_0, j_0}$ and go to Step 2.

384 The restricted myopic scheduling method is very simple to implement. However, this method does
385 not consider potential call-ins in the future. We propose another sequential scheduling algorithm, which
386 considers forecasted patient requests from current time to the appointment day. Specifically, for each call-
387 in patient, we generate a schedule from the existing schedule considering the forecasted future patients.
388 We believe that the information about anticipated patients contributes to the scheduling algorithm in
389 two ways: (i) this information can limit the number of patients with high no-show probabilities in the
390 final schedule and (ii) slot allocation decisions anticipate possible future patient call-ins.

391 In this study, we simply use average historical data to predict future patient demand. Assume that
392 we are generating forecasted patient demand for a day that is T minutes ahead from current time.
393 We can use the average of q pieces of historical patient demand data that happened T or less minutes
394 before virtual appointment days. We denote the forecasted patient demand by \bar{n}_j for $j \in J$. To avoid
395 overestimating future patient arrivals, we may discount our forecasting by α with $0 \leq \alpha \leq 1$. When $\alpha \bar{n}_j$
396 is not an integer, we can round it to the nearest integer.

397 **Forecasting-based Sequential Scheduling Method 1.**

398 **Step 1.** Set $S_{i,j}^0 = 0$ for all $i \in I$ and $j \in J$.

399 **Step 2.** Wait for k^{th} patient call.

400 **Step 3.** The k^{th} patient calls in and the patient is of type j_0 .

401 **Step 4.** Predict future patient demand and obtain \bar{n}_j for $j \in J$.

402 **Step 5.** Perform the scheduling algorithm that considers requests $\alpha \bar{n}_1, \dots, \alpha \bar{n}_{j_0-1}, \alpha \bar{n}_{j_0} + 1, \alpha \bar{n}_{j_0+1}, \dots, \alpha \bar{n}_{|J|}$
403 to generate a schedule S^* from S^0 for (18).

404 **Step 6.** If $\exists i \in I$ such that $S_{i, j_0}^* - S_{i, j_0}^0 \geq 1$, assign the k^{th} patient to slot i and update $S^0 = S^0 + \Delta_{i, j_0}$.

405 **Step 7.** If $S_{i,j}^* - S_{i,j}^0 = 0$ for $i \in I$ and $j \in J$, stop. Otherwise, go to Step 2.

406 When $\alpha = 0$, the forecasting-based scheduling method reduces to the myopic scheduling method in
407 [15].

408 Comparing these two sequential scheduling methods, the restricted myopic method is conservative
409 because it mostly considers available patient information while the forecasting based method is aggressive
410 since it heavily uses predicted information on potential patient calls. Note that the successful application
411 of both of them requires that the physician has enough patient demand which is always the case in
412 practice. In Section 5, we perform a computational study to compare the proposed algorithms with the
413 traditional myopic scheduling algorithm in [15].

5 Computational Study

We perform a computational study to test the performance of proposed algorithms. We consider three experimental settings. In the first setting, we show the effect of considering heterogeneous patients instead of homogeneous patients. In the second setting, we compare the proposed sequential scheduling algorithms with the traditional myopic scheduling algorithm in [15]. In the last setting, we analyze the effect of overflow cost on expected profit.

Throughout our experiments, we assume that a clinic session is 4 hours and partitioned into 8 equal length slots. The service rate λ , which is equal to 2, is constant during the session. Unless explicitly mentioned, $r = 100$, $c_i = 40$ for $i \neq |I|$ and $c_{|I|} = 200$.

5.1 Homogeneous versus Heterogeneous Patients

A major contribution of this study is that the variability of no-show rates is taken into consideration while designing the scheduling algorithms. Algorithm 1 is used to schedule heterogeneous patients. The heuristic algorithm proposed by Kaandorp and Koole [7] is used to schedule homogeneous patients. We consider three types of patients. p_2 is set to 0.5, and p_1 and p_3 are randomly generated such that $(p_1 + p_3)/2 = 0.5$. We assume equal number of patients in each group ($n_1 = n_2 = n_3 = n$). We consider different values for n ($n = 1, \dots, 12$) to analyze the effect of variance on expected total profit for different population sizes. The variance of no-show rates is derived from the variance of p_1 , p_2 and p_3 .

Figure 3 shows the results of 400 randomly generated problems. We observe that the expected profit obtained from the heterogeneous scheduling model dominates the one obtained from the homogeneous scheduling model for all population sizes. The impact of variance of no-show rates on expected profit becomes more significant as the number of patients increases. Figure 4 highlights the difference for 6 and 12 patients. The consideration of heterogeneous patients leads to greater improvements when variance is greater. Especially, when $n = 12$, the improvement on total expected profit could reach up to 20%. Algorithm 1 schedules more patients with low no-show probabilities. However, the total number of patients scheduled is less. As a consequence, the variance in expected profit is less than the one obtained by homogeneous model.

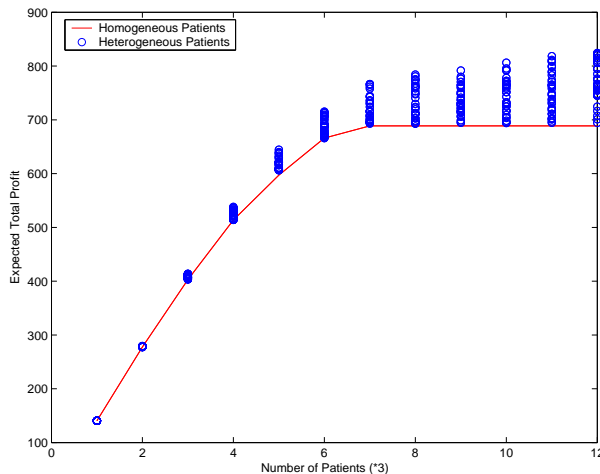


Figure 3: Expected Total Profit vs. Number of Patients

5.2 Sequential Scheduling

We compare the proposed sequential scheduling methods with the traditional myopic scheduling method by Muthuraman and Lawley [15]. We set $|J| = 2$, $p_1 = 0.8$ and $p_2 = 0.2$. We first randomly generate 100

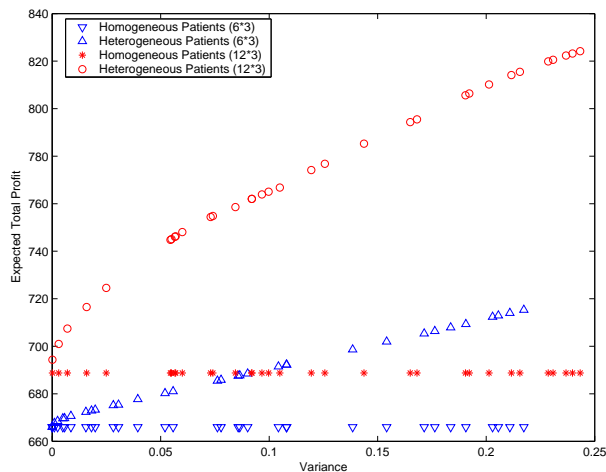


Figure 4: Impact of Variance for 6×3 and 12×3 Patients

443 call-in sequences that span over 30 days. We assume that the call-in rate increases as the call-in time
 444 gets closer to the appointment time. At the beginning, the call-in rate for patients of type j_1 is smaller
 445 than the rate for patients of type j_2 . As time goes on, it increases and finally becomes larger than that
 446 for patients of type j_2 . Let $\tilde{\lambda}_1$ and $\tilde{\lambda}_2$ be the arrival rates of patients of types j_1 and j_2 , respectively.
 447 Specifically, once a call-in of type j_1 is generated, we update $\tilde{\lambda}_1 = \gamma \tilde{\lambda}_1$ where γ is a randomly generated
 448 positive number that is larger than 1. Similarly, we keep updating $\tilde{\lambda}_2$ by a randomly generated number
 449 that is less than 1. We generally set parameters in a way such that the number of expected call-ins of
 450 both type 1 and 2 throughout 480×30 mins are more than the number of expected services, $\lambda \times 8 = 16$.
 451 In our experiments, we set the initial values $\tilde{\lambda}_1 = \frac{1}{600}$ and for $\tilde{\lambda}_2 = \frac{1}{300}$. Random numbers are generated
 452 from $(1, 1 \pm 0.05]$ respectively.

453 First, we consider the restricted myopic scheduling algorithm. Figure 5 shows the expected total profit
 454 for restricted myopic scheduling algorithm (for both $\bar{B}_2 = 4$ represented by *, and $\bar{B}_2 = 8$ represented
 455 by \times) and the traditional myopic method for 100 randomly generated sequences. The results from the
 456 restricted myopic method always dominates those from the traditional myopic scheduling method of [15].
 457 The results from the proposed method are very stable (with less variance), while the traditional myopic
 458 method gives results with high variance. As expected, the proposed method obtains better results when
 459 the upper bound on the number of patients of type j_2 is lower.

460 To predict potential patient call-ins for a call-in sequence, we randomly choose 4 other sequences to
 461 predict numbers of call-ins and set $\alpha = 0.5$ to discount risk. Figure 6 shows the expected total profit
 462 for the forecasting-based sequential scheduling method (represented by *) and the traditional myopic
 463 scheduling method for 100 randomly generated call-in sequences. The results show the advantage of
 464 using forecasted patient demand to generate schedules sequentially. In cases where many potential
 465 patients with low no-show probabilities are expected, it would be wiser to reserve the capacity for these
 466 patients rather than adding patients with high no-show probabilities into the schedule. Obviously, this
 467 argument justifies the open-access scheduling model in which the capacity is kept until the day before
 468 the appointment day or the appointment day since their no-show probabilities are low in those days.

469 In fact, both our theoretical analysis in Sections 2-4 and our computational study on the performance
 470 of the restricted myopic scheduling method and the forecasting-based scheduling method show that
 471 patients with low no-show probabilities should not be restricted by open access model. Actually, our
 472 models show that allowing patients with low no-show probabilities to make appointments ahead is
 473 effective.

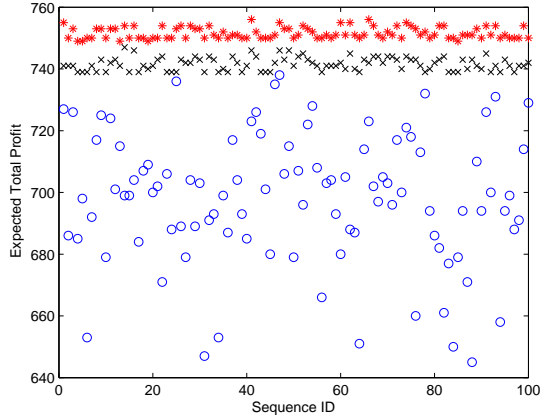


Figure 5: Restricted Myopic Scheduling Method vs Myopic Scheduling Method

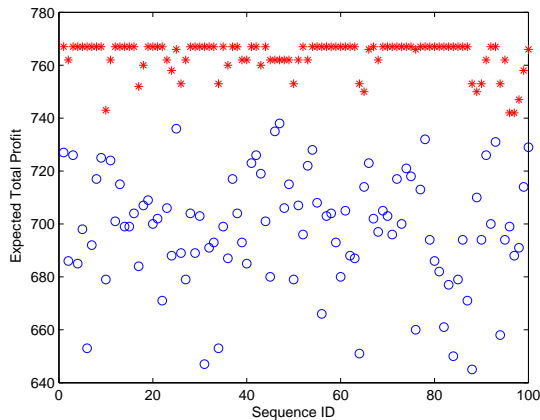


Figure 6: Forecasting-base Method vs Traditional Myopic Method

474 **5.3 Cost of Patient Waiting Time**

475 The expected revenue and overtime cost can typically be estimated by health care practitioners. However,
 476 the value of patient waiting time is determined subjectively. Since the cost associated with patient waiting
 477 time is not an actual cost paid by the clinics, it is important to investigate the effect of different cost
 478 values on scheduling and expected profit. As discussed in Section 2, patient waiting time is directly
 479 related to the size of overflow between slots. So, we control the value of patient waiting time through
 480 changing the value of c_i for $i \neq |I|$. Similarly, the cost of physician overtime can be controlled through
 481 changing the value of c_I . In our experiment, we assume $c_i = c_j$ if $i \neq j$ for $i, j \in I \setminus \{|I|\}$.

482 We consider two schedules generated by the traditional myopic scheduling method from two calling
 483 sequences described in Section 4.2 because of their simple patient structures. For these two sequences,
 484 we set $p_1 = 0.8$ and $p_2 = 0.2$. Figure 7 shows the effect of c_i on expected total profit for both sequences.
 485 As overflow cost (c_i) increases, expected total profit decreases for both schedules. However, the expected
 486 profit of the schedule for the second call-in sequence decreases faster than that of the schedule for the first
 487 sequence. When c_i is small (which means that the physician time is more valuable than patient waiting
 488 time), the performances of two schedules are close to each other. In such cases, differentiating patients
 489 by their no-show probabilities is not very beneficial. When c_i is increasing, the difference between the
 490 two schedules becomes larger. If the patients' waiting time is considered as a significant part of the

491 performance measure, the number of patients with high no-show probabilities should be restricted.

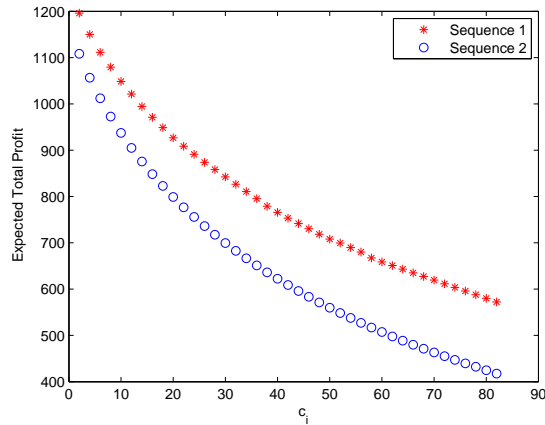


Figure 7: Expected Total Profit vs c_i

492 Laganga and Lawrence [11] mention that the *net overbooking utility*, which is the expected net return
493 generated by overbooking, is larger in the case where patients have high no-show probabilities than in
494 the case where patients have low no-show probabilities. They further mention that this phenomenon is
495 more significant when the cost of patient waiting time and physician overtime are high. According to our
496 computational results in Figure 7 and Figure 8, the expected total profit of schedules generated using
497 overbooking decrease when cost of patient waiting time and physician overtime increases. Furthermore,
498 the speed of decrease is faster in the case where patients have higher no-show probabilities. These results
499 indicate that overbooking can compensate the loss from patient no-shows to some extent. However,
500 reducing patients' no-show rates should have higher priority than applying overbooking.

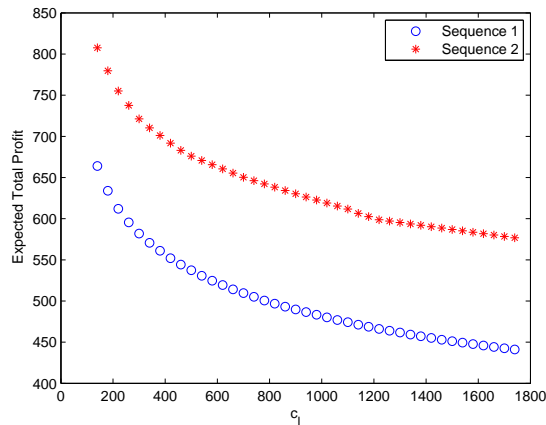


Figure 8: Expected Total Profit vs c_I

501 6 Concluding Remarks and Managerial Insights

502 Various scheduling models with overbooking have been proposed to help health care providers alleviate
503 the negative effects of patient no-shows. However, to the best of our knowledge, all existing studies
504 either assume that patients are homogeneous in terms of their no-show probabilities, or do not consider

505 the impact of different no-show rates on general performance measures. In this paper, we systematically
506 study a clinical scheduling model with overbooking for a set of heterogeneous patients, i.e. their no-show
507 probabilities are different. We prove that, unlike the overbooking model for homogeneous patients, the
508 model for heterogeneous patients is not multimodular. It is very difficult to obtain an optimal schedule
509 since the local optimal solution is not guaranteed to be global optimal. We develop a guided local search
510 algorithm based on the properties of an optimal schedule. We observe that homogeneous overbooking
511 models using the mean value of show-up probabilities are not enough to build high quality schedules. The
512 variance of no-show probabilities have a significant impact on the performance of overbooked schedules.
513 Further, we show the disadvantages of the traditional myopic sequential scheduling method and propose
514 two improved sequential scheduling algorithms that give better schedules.

515 Next, we provide some managerial insights based on our theoretical derivations and computational
516 results. These insights can help health care practitioners better manage clinic scheduling when patients'
517 no-show probabilities are different but can be estimated.

- 518 1. Clustering patients according to their no-show probabilities and using our clinical scheduling meth-
519 ods for heterogeneous patients will help to build schedules with better performances.
- 520 2. Patients with low no-show probabilities are always preferable in schedule generation. This result
521 justifies the open-access scheduling approach, because no-show probabilities increase as the interval
522 between the call-in time and appointment time increases. However, appointments for patients with
523 low no-show probabilities can be made earlier.
- 524 3. Overbooking is beneficial for open-access scheduling systems, because it reduces fluctuations in
525 clinic workload and helps to control demand over time.
- 526 4. The traditional myopic scheduling method proposed in [15] performs well when there is enough
527 patient with low no-show probabilities at the beginning of the call-in sequence. Its performance
528 can be improved significantly by restricting the number of patients with high no-show probabilities
529 in the schedule or using the information of potential patient call-ins.
- 530 5. If costs of patient waiting time and physician overtime are high, few patients with high no-show
531 probabilities should be scheduled.
- 532 6. To reduce overtime cost, patients with low no-show probabilities should be assigned into early slots
533 and patients with high no-show probabilities should be assigned to later slots.

534 Future research directions include extending our research to multiple physician (server) systems,
535 since several physicians collaborate and share the same set of patients. Another direction is developing
536 scheduling methods considering cancelations and unpunctual arrivals. Finally, as [11, 15] point out,
537 overbooking models can also be used in other appointment-based service systems such as law offices,
538 counseling centers and photo studios.

539 References

- 540 [1] E. Altman, B. Gaujal, and A. Hordijk. Multimodularity, convexity and optimization properties.
541 *Mathematics of Operations Research*, 25:324–347, 2000.
- 542 [2] T. Bodenheimer and K. Grumbach. *Understanding Health Policy: A Clinical Approach*. Lange
543 Medical Books / McGraw-Hill, Medical Publishing Division, New York, third edition, 2002.
- 544 [3] T. Cayirli and E. Veral. Outpatient scheduling in health care: a review of literature. *Production*
545 *and Operations Management*, 12:519–549, 2003.
- 546 [4] Centers for Medicare and Office of the Actuary Medicaid Services. National health care expenditures
547 projections: 2005-2015.

- 548 [5] S. Garuda, R. Javalgi, and V. Talluri. Tackling no-show behavior: A market-driven approach.
549 *Health Marketing Quarterly*, 15:25–44, 1998.
- 550 [6] B. Hajek. Extremal splitting of point processes. *Mathematics of Operations Research*, 10:543–556,
551 1985.
- 552 [7] G. Kaandorp and G. Koole. Optimal outpatient appointment scheduling. *Health Care Management
553 Science*, 10:217–229, 2007.
- 554 [8] S. Kim and R. Giachetti. A stochastic mathematical appointment overbooking model for healthcare
555 providers to improve profits. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems
556 and Humans*, 36:1211–1219, 2006.
- 557 [9] G. Koole and E. van der Sluis. Optimal shift scheduling with a global service level constraint. *IIE
558 transactions*, 35:1049–1055, 2003.
- 559 [10] L. Laganga and S. Lawrence. Clinic overbooking to improve patient access and provider productivity.
560 *Decision Sciences*, 38:251–276, 2007.
- 561 [11] L. Laganga and S. Lawrence. An appointment overbooking model to improve client access and
562 provider productivity. Technical report, University of Colorado at Boulder, 2007.
- 563 [12] L. Liu and S. Liu. Dynamic and static job allocation for multi-server systems. *IIE Transactions*,
564 30:845–854, 1998.
- 565 [13] M. Murray and DM. Berwick. Advanced access: reducing waiting and delays in primary care.
566 *JAMA: The Journal of the American Medical Association*, 289:1035–1040, 2003.
- 567 [14] M. Murray and C. Tantau. Same-day appointments: exploding the access paradigm. *Family Practice
568 Management*, 7:45–50, 2000.
- 569 [15] K. Muthuraman and M. Lawley. A stochastic overbooking model for outpatient clinical scheduling
570 with no-shows. *to appear in IIE Transactions*, 2008.
- 571 [16] G. Randolph, M. Murray, J. Swanson, and P. Margolis. Behind schedule: improving access to care
572 for children one practice at a time. *Pediatrics*, 113:230–237, 2004.
- 573 [17] C. Rust, N. Clark, W. Clark, D. Jones, and W. Wilcox. Patient appointment failures in pediatric
574 resident continuity clinics. *Archives of Pediatrics and Adolescent Medicine*, 149(6):693–695, 1995.
- 575 [18] W. Shonick and B. Klein. An approach to reducing the adverse effects of broken appointments in
576 primary care systems: Development of a decision rule based on estimated conditional probabilities.
577 *Medical Care*, 15:419–429, 1977.