

**DBS-20-02**

**「Coordinated Appointment Scheduling with  
Multiple Providers and Patient-and-Physician  
Matching Cost in Specialty Care」**

**Shenghai Zhou**

Jiangxi University of Finance and Economics

**Debiao Li**

Fuzhou University

**Yong Yin**

Graduate School of Business; Doshisha University

(Summary)

To achieve effective care, it is critical to match patients with capable physicians in specialty care. Motivated by the rising popularity of patient-and-physician matching applications in specialty care, this study optimizes the matching and appointment scheduling problems simultaneously in a stochastic environment, in which a decision-maker determines the patient-and-physician pair assignment and the starting times of services. We develop a stochastic optimization model to minimize the matching

DBS Discussion Paper Series supported by the OMRON Foundation.DBS

and operational costs (i.e., patients' waiting time costs, service providers' idle time and overtime costs). This paper is the first study that incorporates matching and appointment scheduling problems together. The benefits of combining these two problems are enormous. The experimental results show that the operational costs gap is as large as 51% between the ill-matched and the well-matched patient-and-physician scenarios. We first reformulate this problem as a two-stage optimization problem. With the analysis for the optimal solution of the second stage problem, a Benders decomposition algorithm is developed. To improve the efficiency of the proposed algorithm, we also prove a low bound of our problem and use it to construct a set of feasibility cuts. Then, we extend our method to incorporate no-shows. Our algorithm can solve problems efficiently, and it can obtain optimal solutions for medium-size problems within 2 or 3 minutes. In contrast, traditional optimal methods require nearly 2 hours. For large-size problems, our algorithm can obtain optimal solutions within 5 or 6 minutes, whereas traditional optimal methods cannot generate a result within 5 hours. Finally, numerical experiments are conducted to evaluate the performance of our proposed algorithm and to investigate the variation of the optimal solutions in different scenarios. To provide quality care as well as minimize the total cost of appointment scheduling in specialty care, we suggest that physicians should develop or train their specialties based on the local patients' disease pattern. We also disclose that the no-show has less influence on the service system when the weight of the matching cost is substantial.

**Keywords:** health care, Benders decomposition, appointment scheduling, multiple service providers, patient-and-physician matching.

# Coordinated Appointment Scheduling with Multiple Providers and Patient-and-Physician Matching Cost in Specialty Care

Shenghai Zhou<sup>a</sup>, Debiao Li<sup>b,\*</sup>, Yong Yin<sup>c</sup>

<sup>a</sup>School of Business Administration, Jiangxi University of Finance and Economics, Nanchang, China

<sup>b</sup>School of Economics and Management, Fuzhou University, Fuzhou, China

<sup>c</sup>Graduate School of Business, Doshisha University, Kyoto, Japan

---

## Abstract

To achieve effective care, it is critical to match patients with capable physicians in specialty care. Motivated by the rising popularity of patient-and-physician matching applications in specialty care, this study optimizes the matching and appointment scheduling problems simultaneously in a stochastic environment, in which a decision-maker determines the patient-and-physician pair assignment and the starting times of services. We develop a stochastic optimization model to minimize the matching and operational costs (i.e., patients' waiting time costs, service providers' idle time and overtime costs). This paper is the first study that incorporates matching and appointment scheduling problems together. The benefits of combining these two problems are enormous. The experimental results show that the operational costs gap is as large as 51% between the ill-matched and the well-matched patient-and-physician scenarios. We first reformulate this problem as a two-stage optimization problem. With the analysis for the optimal solution of the second stage problem, a Benders decomposition algorithm is developed. To improve the efficiency of the proposed algorithm, we also prove a low bound of our problem and use it to construct a set of feasibility cuts. Then, we extend our method to incorporate no-shows. Our algorithm can solve problems efficiently, and it can obtain optimal solutions for medium-size problems within 2 or 3 minutes. In contrast, traditional optimal methods require nearly 2 hours. For large-size problems, our algorithm can obtain optimal solutions within 5 or 6 minutes, whereas traditional optimal methods cannot generate a result within 5 hours. Finally, numerical experiments are conducted to evaluate the performance of our proposed algorithm and to investigate the variation of the optimal solutions in different scenarios. To provide quality care as well as minimize the total cost of appointment scheduling in specialty care, we suggest that physicians should develop or train their specialties based on the local patients' disease pattern. We also disclose that the no-show has less influence on the service system when the weight of the matching cost is substantial.

Keywords: health care, Benders decomposition, appointment scheduling, multiple service providers, patient-and-physician matching

---

\*Corresponding Email: debiaoli@fzu.edu.cn

Email addresses: zhoushenghai@jxufe.edu.cn (Shenghai Zhou), debiaoli@fzu.edu.cn (Debiao Li), yyin@mail.doshisha.ac.jp (Yong Yin)

33 1. Introduction

34 Specialty care clinics, in which specialists have extensive training and education, are designed  
35 to provide specific diagnoses and treatments. As modern medical development, specialties are dis-  
36 tinguished even in the same clinics. For example, in otorhinolaryngology, specialty care physicians  
37 (SCPs) diagnose and treat a wide range of diseases around the ear, nose, and throat regions. If a  
38 patient chooses a physician who is not in the right expertise, effective care may not be achieved. In  
39 practice, patients select SCPs either through referrals from their primary care physicians (PCPs)  
40 or by themselves. In a specialty care system with referrals, e.g., the United States, referrals by  
41 PCPs are frequently failed and often led to medical errors [1, 2]. It can be worse without the help  
42 of referrals. Especially for patients, because they have limited medical education and knowledge,  
43 but have to choose SCPs by themselves. On the other side, SCPs would prefer to see more patients  
44 that fall into their areas of clinical expertise, such that their medical training and education can be  
45 realized. Therefore, it is very critical to match patients with capable physicians in specialty care,  
46 such that the quality of care can be guaranteed to some extent.

47  
48 With the rapid development of information technology in healthcare, some applications have  
49 been invented to solve this patient and physician matching problem in specialty care recently.  
50 For example, in the United States, Specialty Care Connect ([https://armadahealth.com/patient-](https://armadahealth.com/patient-physician-matching)  
51 [physician-matching](https://armadahealth.com/patient-physician-matching)) powered by Armada Health had applied analytic models to connect patients  
52 to the right SCPs. In China, an intelligent matching system, named “Rui Zhi”, has been invented  
53 by Tencent, which is one of the biggest internet-based technology corporations in China. It has been  
54 implemented in Guangzhou women and children’s medical center since May 2018. The accuracy of  
55 diagnoses was claimed to be 94%, and the matching accuracy was reported to be 96%. Recently, a  
56 patient-and-physician matching index is proposed to measure the capability between patients and  
57 SCPs through an improved multi-disease pre-diagnosing Bayesian network model, which is based  
58 on given patient’s symptoms and physician’s specialties [3]. The experimental results show that the  
59 proposed patient-and-physician matching index increases the physician matching accuracy under  
60 various settings. Integrated with data science technologies and medical knowledge, these matching  
61 applications and mechanisms are designed to ensure the effectiveness of specialty care.

62  
63 Apart from the effectiveness, timely, and accessible coordination in specialty care is also very  
64 critical. In order to improve the efficiency of specialty care operations, appointment scheduling is  
65 an essential and efficient way to maximize physicians’ time utilities and improve patients’ satis-  
66 faction by reducing waiting time. Appointment scheduling has been well studied in primary care.  
67 In specialty care, when we make a schedule, we should take the unique characteristic of specialty  
68 into consideration, i.e., the matching between patients and physicians, so as to make a more rea-  
69 sonable schedule. To the best of our knowledge, there is no appointment scheduling literature that  
70 considers the effectiveness through the patient and physician matching. To ensure the effectiveness  
71 and efficiency of specialty care, the matching mechanism should be embedded in the appointment  
72 scheduling system, such that timely access and efficiency can be achieved at the same time.

73  
74 In an appointment scheduling problem in specialty, there are usually multiple physicians (with  
75 appointment scheduling terminology, physicians often refer to service providers) provide services.  
76 Thus, our studied appointment scheduling problem in specialty care is a multiple-provider appoint-  
77 ment scheduling problem indeed. Due to the vast majority of work on matching between patients

78 and physician preferences, in this work, we regard the matching cost between patients and physi-  
79 cians as the inputs to the appointment scheduling model.

80 To sum up, in this study, we integrate the matching cost between patients and physicians into  
81 appointment scheduling problem for the specialty care, in which the service times are stochas-  
82 tic. Our objective is to jointly optimize the assignment and job allowance for each patient, such  
83 that the weighted operational and matching costs are minimized. We first formulate our studied  
84 problem as a two-stage optimization problem based on the sample average approximation (SAA)  
85 approach. After analyzing the properties for the optimal solution of the second stage problem, an  
86 improved Benders decomposition algorithm is proposed. Finally, we conduct computational experi-  
87 ments to identify the efficiency of our proposed algorithm and investigate some managerial insights.

88  
89 The contributions of this paper are summarized as follows. First, we adapt the patient-and-  
90 physician matching cost into the appointment scheduling problem. To the best of our knowledge,  
91 we are the first to jointly consider the matching problem and the appointment scheduling problem  
92 with multiple service providers in a stochastic environment. The experimental results in Section 5  
93 give the managerial benefits of combining these two problems. The operational costs gap is as large  
94 as 51% between the ill-matched and the well-matched patients and physicians scenarios. Second,  
95 a two-stage stochastic program is formulated to minimize the weighted operational and matching  
96 costs. We analyze a low bound for the optimal objective value under any potential assignment.  
97 Third, we analyze the properties for the optimal solution of the second stage problem. On this  
98 basis, we propose an improved Benders decomposition algorithm with feasibility cuts to solve this  
99 problem efficiently. The experimental results of Section 5 show that our algorithm can obtain op-  
100 timal solutions for medium-size problems within 2 or 3 minutes. In contrast, traditional optimal  
101 methods require nearly 2 hours. For large-size problems, our algorithm can obtain optimal solu-  
102 tions within 5 or 6 minutes, whereas traditional optimal methods cannot generate a result within 5  
103 hours. Finally, several sensitivity analyses are conducted under different parameter settings. Our  
104 numerical results indicate the importance of the patient-and-physician matching. To provide qual-  
105 ity care as well as minimize the total cost of appointment scheduling in specialty care, we suggest  
106 that physicians should develop or train their specialties based on the local patients' disease pattern.

107  
108 The overview of this paper is organized as follows. The literature related to patient-and-  
109 physician matching and appointment scheduling is reviewed in Section 2. We formally describe  
110 our studied problem in Section 3. The details of our proposed improved Benders decomposition  
111 approach are stated in Section 4. We extend our method to incorporate no-shows in Section 5.  
112 In Section 6, we conduct some numerical experiments to verify our proposed method and examine  
113 some potential insights. Some managerial implications are summarized in Section 7, followed by  
114 the summary and future work in Section 8.

## 115 2. Literature Review

116 In this section, we review the literature that is most relevant to our research. In particular, we  
117 focus mainly on patient-and-physician matching problems and fundamental appointment scheduling  
118 problems.

### 119 2.1. Patient-and-physician matching

120 The patient-and-physician matching problems of specialty care are different from primary care  
121 and elective surgery. In primary care, access to services is the most critical factor for patient-and-

122 physician matching [4]. Several studies also provide evidence that physicians’ interpersonal skills  
123 affect patients’ satisfaction [5] and treatment outcomes [6]. Gong et al. [7] propose a weighted aver-  
124 age model to recommend physicians by considering the economic matching degree, medical domain  
125 matching degree, recommender influence, and region reference. A time-constraint probability factor  
126 graph model learns these features from a real-world medical data set, which was optimized by a  
127 constraint-based optimization framework. Liu et al [8] examine the patients’ preferences and choice  
128 behavior in the scheduling appointment, which include the gender effect, speed, and physician of  
129 choice. Although the physician of choice is highly correlated to the quality of care, the detail of the  
130 physician of choice is not disclosed.

131  
132 With the development of modern medical science, physicians tend to have distinguished spe-  
133 cialties, although they are in the same department. The medical skills or specialties of physicians  
134 become more critical in patient-and-physician matching. Kinchen et al. [9] analyze the significant  
135 factors affecting the choice of specialists by primary care physicians through a survey. Medical  
136 skills, appointment timeliness, and quality of specialist communication are the three most impor-  
137 tant factors. Most existing literature studies the variations and their causes in referral decision  
138 making among PCPs. However, few explore the appropriateness of the referral decision, which is  
139 mainly defined by the medical skill [2].

140  
141 Pan et al. [10] propose a dynamic preference learning algorithm to recommend physicians in  
142 specialty care by considering both patients’ preferences and their heterogeneous illness conditions.  
143 Furthermore, it is assumed that general practitioners correctly evaluated patients’ illness conditions,  
144 and there was no bias to refer physicians in specialty care. However, general practitioners may lack  
145 related expertise and have some biased information about SCP. To eliminate biases or mistakes of  
146 referrals in specialty care, a pre-diagnosing model is applied to gain a more accurate diagnosis of  
147 patients’ disease(s). Given the specialty information of a physician, a patient-and-physician match-  
148 ing index is proposed to measure the quality influence during the matching [3].

149  
150 However, the matching literature mentioned above only considers the isolated matching between  
151 patients and physicians but ignores the operational aspects (i.e., timely and accessible) during  
152 the specialty care visit. Thus, we need to further integrate the matching with the appointment  
153 scheduling problem.

## 154 2.2. Appointment Scheduling

155 In the literature, one classic appointment scheduling problem refers to the intra-day scheduling  
156 problem, which focuses on making appointment decisions on a given day. For this kind of appoint-  
157 ment scheduling problem, usually, it is assumed that only one service provider provides service. The  
158 decision-maker needs to determine the start time of each appointment so as to balance the costs  
159 for both patients and service providers. Specifically, most works, including ours, study the objec-  
160 tive of minimizing the total expected (weighted) costs of patients’ waiting times, service provider’s  
161 idle times, and overtime [11, 12]. If there is no session length constraint, some studies only take  
162 patients’ waiting times and service providers’ idle times into consideration in the objective function  
163 [13, 14, 15, 16]. While other studies, including ours, consider total expected (weighted) costs of pa-  
164 tients’ waiting times and service provider’s overtime in the objective function [17, 18, 19]. Consider  
165 patients’ behavior, some work also take patients’ no-shows into consideration [17, 20, 21, 22]. Our

166 work also consider this important patients' behavior.

167

168 The classic appointment scheduling problems with a single service provider are often considered  
169 in a stochastic environment, such as random service times [11, 23] and random no-shows [21, 22]. In  
170 order to handle those uncertainties, most works focus on developing some stochastic optimization  
171 models [11, 24, 25, 21, 22]. On this basis, the developed stochastic models are often approximated  
172 as corresponding (mixed-integer) linear programs through Sample Average Approximation (SAA)  
173 approach [26, 27, 28]. For those (mixed-integer) linear programs, when the sample size is small  
174 (e.g.,  $\leq 500$ ), they often can be solved directly [27, 28]. Nevertheless, when the sample size is large,  
175 it is difficult to achieve a high accuracy level solution within a reasonable computational time. In  
176 this case, some efficient algorithms are developed, such as the Benders decomposition algorithm  
177 [29, 30, 26], and simulation-based sequential algorithms [31]. In this work, we first formulate our  
178 studied problem as a stochastic program. And then, we also exploit the SAA approach to handle  
179 the stochastic service times. On this basis, we develop an improved Benders decomposition algo-  
180 rithm to solve the problem.

181

182 However, the vast majority of the literature focuses on service systems that only involve one  
183 service process with one service provider [11, 24, 25, 17, 23]. In addition to systems with only one  
184 service provider, systems with many service providers also prevail in practice. However, studies  
185 on appointment scheduling problems with multiple service providers are limited. To the best of  
186 our knowledge, the only appointment scheduling works with multiple service provider systems are  
187 those by [32, 33, 34, 35, 36]. Sickinger et al. [34] consider two CT-scan machines in a radiology  
188 department, in which two machines are regarded as two identical parallel providers with identical  
189 deterministic service times. Similarly, Zacharias et al. [36] also assume the service providers are  
190 identical, and the service time is also identical deterministic. In contrast, Alvarez et al. [32] con-  
191 sider stochastic service times in a two-stage service system in which two identical parallel service  
192 providers in the first stage. Their purpose is to minimize the total weighted costs of patients' wait-  
193 ing and service providers' idling. Zheng et al. [33] consider no-shows in their model; however, their  
194 model cannot be easily adapted to the overtime case. Soltani et al. [35] also consider stochastic  
195 service times, and patients no-shows for an appointment scheduling problem with multiple service  
196 providers. Different from Zheng et al [33], they consider both patients' waiting time, and service  
197 providers' idle time and overtime in the objective function.

198

199 However, all existing appointment scheduling problems do not consider the quality-related  
200 matching factors between patients and physicians. As we mentioned in the introduction section, the  
201 patient-and-physician matching is critical to improving the effectiveness in specialty care. Therefore,  
202 we integrate the multiple-provider appointment scheduling problem with the patient-and-physician  
203 matching problem in this paper.

### 204 3. Problem Formulation

205 In this paper, we consider generic specialty care with  $k$  service providers (in the rest of this  
206 paper, we use the terminology "service provider" to represent the specialty care physician). There  
207 are totally  $n$  patients needed to be scheduled within a session length  $T$ . Before making a schedule,  
208 the decision-maker has the following information: (1) The matching costs ( $m_{i,j}, i = 1, 2, \dots, n, j =$   
209  $1, 2, \dots, k$ ), which denote the cost for patient  $i$  match physician  $j$ , are assumed to be known in

210 advance. In this work, we take the matching cost as input and it can be derived from the patient-  
211 and-physician matching index proposed by [3]. Li et al. [3] measure the matching indexes between  
212 patients and physicians based on the symptom-specialty relationship through a trained Bayesian  
213 network pre-diagnosing model. We transform the matching indexes into matching costs. Generally  
214 speaking, a higher matching index indicates a better-matched patient and physician. If a patient is  
215 assigned to a capable physician with a high matching index, it most likely leads to a better  
health-216 care outcome, which has a higher potential to reduce healthcare costs. Thus, when we  
transform 217 matching indexes into matching costs, we may assume the higher the matching index,  
the lower 218 the matching cost. For simplification. in this work, we assume the matching costs  
are known. 219 (2) Through some preliminary classification, the service time of patient  $i$  at service  
provider  $j$ , 220  $d_{i,j}$  ( $i = 1, 2, \dots, n, j = 1, 2, \dots, k$ ) is an independent, not necessary identically  
distributed random 221 variable, which is known to the decision-maker.

222  
223 With the above information, the decision-maker of the specialty care needs to determine (1) as-  
224 signment problem and (2) appointment scheduling problem. Specifically, the assignment problem  
225 means how to assign those  $n$  patients to  $k$  service providers. We use  $x_{i,j}$  ( $i = 1, 2, \dots, n, j = 1, 2, \dots, k$ )  
226 to denote the decision variable for assignment problem. If patient  $i$  ( $i = 1, 2, \dots, n$ ) is assigned to  
227 service provider  $j$  ( $j = 1, 2, \dots, k$ ),  $x_{i,j} = 1$ , otherwise  $x_{i,j} = 0$ . Given the assignment, the appoint-  
228 ment scheduling problem means, for each service provider, how to determine the start times (or  
229 equivalently the job allowances) for assigned patients. We use  $s_{i,j}$  ( $i = 1, 2, \dots, n, j = 1, 2, \dots, k$ )  
230 to denote the decision variable (i.e., job allowance) for appointment scheduling problem. And we  
231 have  $s_{i,j} = 0$  if  $x_{i,j} = 0$ , which means if patient  $i$  is not assigned to service provider  $j$ , then we  
232 do not leave any job allowance for that patient  $i$  at service provider  $j$ . We assume that once the  
233 assignment is fixed, the patients assigned to each service provider would go through their service ac-  
234 cording to their index order, i.e., the service provider would handle patient  $i$  before patient  $i'$  if  $i < i'$ .

235  
236 With any given assignment, for each service provider, the corresponding appointment scheduling  
237 degenerates into a classic appointment scheduling. For convenience and clarity, in the rest of this  
238 paper, we refer to a patient in the appointment system as a job and use the terms job and patient  
239 interchangeably.

240 Due to the randomness of the stochastic service times, patients' waiting or service providers'  
241 idling might come up. We use the term  $\bar{W}_{i,j}$  ( $i = 1, 2, \dots, n, j = 1, 2, \dots, k$ ) to denote the actual  
242 waiting time of patient  $i$  before it has been seen by service provider  $j$ . Then we have  $\bar{W}_{i,j} = 0$  if  
243  $x_{i,j} = 0$ . In order to derive the waiting times logically, we introduce the virtual waiting time  $W_{i,j}$   
244 ( $i = 1, 2, \dots, n, j = 1, 2, \dots, k$ ) for patient  $i$  at service provider  $j$ . The virtual waiting time  $W_{i,j}$   
245 indicates the waiting time of patient  $i$  before she/he is seen by service provider  $j$ , regardless of the  
246 actual assignment of patient  $i$ . Given any assignment  $x_{i,j} = 1$ , patient  $i$  may actually need to wait  
247 for service provider  $j$  to be served but must not wait for other service provider  $h \neq j$ . With the  
248 definition of virtual waiting time  $W_{i,j}$ , the actual waiting time  $\bar{W}_{i,j}$  can be achieved by multiplying  
249 the virtual waiting time  $W_{i,j}$  by the assignment  $x_{i,j}$ , i.e.,  $\bar{W}_{i,j} = x_{i,j}W_{i,j}$ .

250  
251  
252 Next, we present how to derive virtual waiting times  $W_{i,j}$  recursively. In the classic single provider  
253 appointment scheduling problem, the actual waiting time is determined through the waiting time,  
254 the service time, and the job allowance of its preceding appointment recursively. However, for our  
255 problem, the actual service time is expressed as  $x_{i,j} \cdot d_{i,j}$ . Thus, in order to derive the virtual waiting



256 times  $W_{i,j}$ , we can modify the service time of our problem and formulate the virtual waiting time  
 257  $W_{i,j}$  as follows:

$$\begin{aligned} W_{i,j} &= \max\{0, x_{i-1,j}d_{i-1,j} + W_{i-1,j} - s_{i-1,j}\} \quad i = 2, \dots, n, j = 1, \dots, k \\ W_{1,j} &= 0 \quad j = 1, \dots, k \end{aligned} \quad (3.1)$$

258 Similarly, we define the virtual idle time  $I_{i,j}$  ( $i = 1, 2, \dots, n, j = 1, 2, \dots, k$ ), which denotes the  
 259 idleness of service provider  $j$  after serving patient  $i$ . Note that for any service provider  $j$ , the  
 260 summation  $\sum_{i=1}^n I_{i,j}$  equals to the actual idleness for service provider  $j$ . Thus, the notation of virtual  
 261 idle time  $I_{i,j}$  is enough to depict our performance indicator and we do not need to define the actual  
 262 idle time. In the rest of this paper, we will omit the term “virtual” for idle time  $I_{i,j}$ . Then we also  
 263 derive the idle time  $I_{i,j}$  recursively as follows:

$$I_{i-1,j} = \max\{0, -x_{i-1,j}d_{i-1,j} - W_{i-1,j} + s_{i-1,j}\} \quad i = 2, \dots, n+1, j = 1, 2, \dots, k \quad (3.2)$$

264 Note that we restrict all services that should be finished within a session length  $T$  for all service  
 265 providers. As a result, the service system may incur some overtime for some service providers. We  
 266 define  $O_j$  ( $j = 1, \dots, k$ ) to represent the overtime for service provider  $j$ , then it is derived as follows:

$$O_j = \max\{0, x_{n,j}d_{n,j} + W_{n,j} - s_{n,j}\} \quad j = 1, 2, \dots, k \quad (3.3)$$

267 For the service system, again,  $m_{i,j}$  states the cost for patient  $i$  being assigned to service provider  
 268  $j$ . Let  $c_i^W$  denote the unit waiting time cost for patient  $i$ . And let  $c_j^I$  and  $c_j^O$  denote the unit idle  
 269 time, and overtime cost for service provider  $j$ , respectively. We define a weight  $\lambda$  to balance the  
 270 matching cost and the weighted operational costs. The objective is to optimize the assignment and  
 271 appointment scheduling simultaneously, such that the expected weighted operational (i.e., waiting  
 272 costs, idling costs, and overtime costs) and matching cost is minimized, as shown in equation (3.4):

$$\sum_{j=1}^k \mathbb{E} \left[ \sum_{i=1}^n (c_i^W x_{i,j} W_{i,j} + c_j^I I_{i,j}) + c_j^O O_j \right] + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \quad (3.4)$$

273 In the objective function (3.4), the weight  $\lambda$  is used to balance different dimensions of matching  
 274 cost and the weighted operational costs. In practice, we may first test different values of  $\lambda$  to  
 275 justify the corresponding matching and operational costs, and then select an appropriate value of  
 276  $\lambda$  to implement. Thus, in our numerical analyses section, we take the values of  $\lambda$  from 0 to 100.  
 277 When  $\lambda = 0$ , it indicates that only the operational cost is considered. Besides, in the numerical  
 278 analyses, we demonstrate the effect of the matching cost by increasing  $\lambda$ , which is the same as the  
 279 standardization of the cost coefficients.

280 Next, we construct constraints for our problem. For decision variables  $x_{i,j}$  and  $s_{i,j}$ , we define the  
 281 following constraints:

$$\begin{aligned} \sum_{i=1}^n x_{i,j} &\geq 1 \quad j = 1, 2, \dots, k \\ \sum_{j=1}^k x_{i,j} &= 1 \quad i = 1, 2, \dots, n \\ \sum_{i=1}^n s_{i,j} &= T \quad j = 1, 2, \dots, k \end{aligned} \quad (3.5)$$

282 The first equation in (3.5) makes sure there should be at least one patient assigned to each  
 283 service provider, because there is no necessary to keep one service provider idle in a session length.  
 284 The second equation in (3.5) ensures each patient should be assigned to one service provider. The  
 285 third equation in (3.5) indicates all appointments should be scheduled within session length  $T$ . In  
 286 addition, we need to make sure  $s_{i,j} = 0$  if patient  $i$  is not assigned to service provider  $j$ . Thus, we  
 287 define the following inequalities to achieve this purpose:

$$s_{i,j} \leq Mx_{i,j} \quad i = 1, 2, \dots, n, j = 1, 2, \dots, k \quad (3.6)$$

288 where  $M$  is a big number.

289 Through recursive equations (3.1), (3.2), and (3.3), we derive the following equalities:

$$\begin{aligned} W_{i,j} - I_{i-1,j} &= x_{i-1,j}d_{i-1} + W_{i-1,j} - s_{i-1,j} \quad i = 2, 3, \dots, n, j = 1, 2, \dots, k \\ O_j - I_{n,j} &= x_{n,j}d_n + W_{n,j} - s_{n,j} \quad j = 1, 2, \dots, k \end{aligned} \quad (3.7)$$

290 With performance indicators derived in equations (3.1),(3.2) and (3.3), the objective is to jointly  
 291 optimize the assignment and appointment scheduling, such that the expected weighted operational  
 292 (i.e., waiting costs, idling costs and overtime costs) and matching cost is minimized. Thus, our  
 293 studied problem can be formulated as the following stochastic model (M0):

$$\begin{aligned} \text{(M0)} \quad & \min_{\mathbf{x}, \mathbf{s}} \sum_{j=1}^k \mathbb{E} \left[ \sum_{i=1}^n (c_i^W x_{i,j} W_{i,j} + c_j^I I_{i,j}) + c_j^O O_j \right] + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \\ & \text{s.t. (3.5), (3.6), (3.7)} \\ & x_{i,j} \in \{0, 1\}, s_{i,j} \geq 0 \end{aligned} \quad (3.8)$$

## 294 4. Proposed Method

295 To solve this stochastic problem, we first exploit the SAA method to handle the stochastic  
 296 service times. On this basis, an improved Benders decomposition algorithm with feasibility cuts is  
 297 developed to solve the approximated problem efficiently.

### 298 4.1. SAA-based Formulation

299 As indicated in the literature, the SAA method is an efficient scenario-based method for solv-  
 300 ing stochastic programming problems, and has been widely used to solve appointment scheduling  
 301 problems [23, 37, 26]. Specifically, with given service time distribution  $\mathbf{d}$ , we randomly generate  
 302  $H$  i.i.d. realizations. Then, the stochastic program (3.8) can be approximated by the following  
 303 deterministic program (DP):

$$\begin{aligned}
\text{(DP)} \quad & \min_{\mathbf{x}, \mathbf{s}, \mathbf{W}, \mathbf{I}, \mathbf{O}} \frac{1}{H} \sum_{h=1}^H \sum_{j=1}^k \left[ \sum_{i=1}^n (c_i^W x_{i,j} W_{i,j}^h + c_j^I I_{i,j}^h) + c_j^O O_j^h \right] + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \\
& s.t. \quad \sum_{i=1}^n x_{i,j} \geq 1 \quad j = 1, 2, \dots, k \\
& \quad \sum_{j=1}^k x_{i,j} = 1 \quad i = 1, 2, \dots, n \\
& \quad \sum_{i=1}^n s_{i,j} = T \quad j = 1, 2, \dots, k \\
& \quad s_{i,j} \leq M x_{i,j} \quad i = 1, \dots, n, j = 1, \dots, k \\
& \quad W_{i,j}^h - I_{i-1,j}^h = x_{i-1,j} d_{i-1}^h + W_{i-1,j}^h - s_{i-1,j} \quad i = 2, \dots, n, j = 1, \dots, k, h = 1, \dots, H \\
& \quad O_j^h - I_{n,j}^h = x_{n,j} d_n^h + W_{n,j}^h - s_{n,j} \quad j = 1, \dots, k, h = 1, \dots, H \\
& \quad W_{1,j}^h = 0 \quad j = 1, \dots, k, h = 1, \dots, H \\
& \quad x_{i,j} \in \{0, 1\}, s_{i,j} \geq 0
\end{aligned} \tag{4.1}$$

304 In the above deterministic program based on SAA,  $d_{i,j}^h$  denotes the realization of service time  
305  $(d_{i,j})$  under realization  $h$ , the variables  $W_{i,j}^h$ ,  $I_{i,j}^h$  and  $O_j^h$  denote the corresponding waiting time, idle  
306 time and overtime under realization  $h$ , respectively. Note that except for the original decision vari-  
307 ables  $\mathbf{x}$  and  $\mathbf{s}$ , we let the performance indicators  $\mathbf{W}, \mathbf{I}$  and  $\mathbf{O}$  as new decision variables, to linearize  
308 those performance indicators.

309 For the above deterministic program, it is a non-linear mixed-integer linear program. We can  
310 even introduce a big-M method to reformulate it as a mixed-integer linear program. However, when  
311 the problem size is large, it is difficult to achieve an optimal solution in a reasonable computational  
312 time. According to our preliminary test, when  $H = 1000$ ,  $n = 40$ ,  $k = 4$ , respectively, the optimal  
313 solution of the corresponding MILP cannot be achieved within 5 hours. Thus, we propose an  
314 improved Benders decomposition to solve above DP efficiently.  
315

#### 316 4.2. Improved Benders Decomposition

317 The Benders decomposition method is suitable for some large scale problems with special struc-  
318 ture. Our problem has the structure that the subproblem can be solved to optimality without actu-  
319 ally solving the corresponding problem, which is suitable for the Benders decomposition. Before  
320 we introduce the Benders decomposition method, we first analyze a lower bound for the objective  
321 function, which helps to generate some feasibility cuts for the proposed algorithm. We define the  
322 individual cost  $C_{i,j}$  for the original problem (M0) as follows:

$$C_{i,j} = \begin{cases} \mathbb{E}[c_{i+1}^W x_{i+1,j} W_{i+1,j} + c_j^I I_{i,j}] & i = 1, 2, \dots, n-1 \\ \mathbb{E}[c_j^O O_j + c_j^I I_{i,j}] & i = n \end{cases} \quad j = 1, \dots, k \tag{4.2}$$

323 Lemma 1. For any given assignment  $\mathbf{x}$ , the individual cost  $C_{i,j}$  is bounded as follows:

$$C_{i,j} \geq \begin{cases} x_{i+1,j} g_{i,j} & i = 1, 2, \dots, n-1 \\ g_{i,j} & i = n \end{cases} \quad j = 1, \dots, k \quad (4.3)$$

324 where

$$g_{i,j} = \begin{cases} \min_{s_{i,j}} \mathbb{E}[c_{i+1}^W [d_{i,j} - s_{i,j}]^+ + c_j^I [d_{i,j} - s_{i,j}]^-] & i = 1, 2, \dots, n-1 \\ \min_{s_{n,j}} \mathbb{E}[c_j^O [d_{i,j} - s_{i,j}]^+ + c_j^I [d_{i,j} - s_{i,j}]^-] & i = n \end{cases} \quad j = 1, \dots, k \quad (4.4)$$

325 where  $[a]^+ = \max\{a, 0\}$  and  $[a]^- = \max\{-a, 0\}$ .

326 Proof: For  $i = 1, 2, \dots, n-1, j = 1, \dots, k$ , we always have  $C_{i,j} \geq x_{i+1,j} \mathbb{E}[c_{i+1}^W W_{i+1,j} + c_j^I I_{i,j}]$ . We now  
327 bound the right hand  $\mathbb{E}[c_{i+1}^W W_{i+1,j} + c_j^I I_{i,j}]$ . Based on the definition, we

$$\begin{aligned} \mathbb{E}[c_{i+1}^W W_{i+1,j} + c_j^I I_{i,j}] &= \mathbb{E}\left[ c_{i+1}^W [W_{i,j} + d_{i,j} - s_{i,j}]^+ + c_j^I [W_{i,j} + d_{i,j} - s_{i,j}]^- \right] \\ &= \mathbb{E}_{W_{i,j}} \left\{ \mathbb{E}_{d_{i,j}} \left[ c_{i+1}^W [W_{i,j} + d_{i,j} - s_{i,j}]^+ + c_j^I [W_{i,j} + d_{i,j} - s_{i,j}]^- \right] \Big| W_{i,j} \right\} \end{aligned}$$

328 Suppose  $s_{i,j}^*$  is the optimal solution to achieve  $g_{i,j}$  such that  $s_{i,j}^* = \arg \min_{s_{i,j}} \mathbb{E}_{d_{i,j}} [c_{i+1}^W [d_{i,j} - s_{i,j}]^+ + c_j^I [d_{i,j} - s_{i,j}]^-]$ .

329 It follows that

$$g_{i,j} \leq \mathbb{E}_{d_{i,j}} [c_{i+1}^W [d_{i,j} - \tilde{s}_{i,j}]^+ + c_j^I [d_{i,j} - \tilde{s}_{i,j}]^-], \forall \tilde{s}_{i,j}$$

330

331 Therefore, for any realization of  $W_{i,j}$ , let  $\tilde{s}_{i,j} = s_{i,j} - W_{i,j}$ , we have

$$g_{i,j} \leq \mathbb{E}_{d_{i,j}} \left[ c_{i+1}^W [W_{i,j} + d_{i,j} - s_{i,j}]^+ + c_j^I [W_{i,j} + d_{i,j} - s_{i,j}]^- \Big| W_{i,j} \right]$$

332

333 By taking expectation for random variable  $W_{i,j}$  for above equation, the inequality still holds,  
334 which implies  $\mathbb{E}[c_{i+1}^W W_{i+1,j} + c_j^I I_{i,j}]$  is bounded by  $g_{i,j}$  for any  $s_{i,j}$ .

335 Similarly, we can also bound  $C_{i,j}$  when  $i = n$ . This completes the proof.  $\square$

336 Note that in Lemma 1,  $g_{i,j}$  corresponds to the optimal cost of a general Newsvendor problem  
337 [38]. Given the cumulative distribution function  $F_{i,j}$  for  $d_{i,j}$ , the optimal solution  $s_{i,j}^*$  can be achieved

338 through  $F_{i,j}(s_{i,j}^*) = \frac{c_{i+1}^W}{c_{i+1}^W + c_j^I}$  for  $i = 1, \dots, n-1, j = 1, \dots, k$ . Thus, the optimal cost  $g_{i,j}$  can be calculated  
339 easily.

340 With Lemma 1, for any given assignment  $\mathbf{x}$ , the operational cost is bounded as  $\sum_{j=1}^k [\sum_{i=1}^{n-1} x_{i+1,j} g_{i,j} +$   
341  $g_{n,j}]$ . Furthermore, Lemma 1 helps to bound some variables in the master problem in the Benders  
342 decomposition algorithm. We will introduce it later.

343

344 We observe that for any solution  $(\mathbf{x}, \mathbf{s})$ , the operational costs are decomposable by scenario  $h$   
345 and service provider  $j$ . And they can be determined through recursive equations (3.1), (3.2) and  
346 (3.3). Due to this kind of special structure, we further reformulate problem (4.1) as the following  
347 two-stage optimization problem:

$$\begin{aligned}
& \min_{\mathbf{x}, \mathbf{s}} \sum_{j=1}^k Q_j(\mathbf{x}, \mathbf{s}) + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \\
& \text{s.t. } \sum_{i=1}^n x_{i,j} \geq 1 \quad j = 1, 2, \dots, k \\
& \quad \sum_{j=1}^k x_{i,j} = 1 \quad i = 1, 2, \dots, n \\
& \quad \sum_{i=1}^n s_{i,j} = T \quad j = 1, 2, \dots, k \\
& \quad s_{i,j} \leq M x_{i,j} \quad i = 1, \dots, n, j = 1, \dots, k \\
& \quad x_{i,j} \in \{0, 1\}, s_{i,j} \geq 0
\end{aligned} \tag{4.5}$$

348 where

$$\begin{aligned}
Q_j(\mathbf{x}, \mathbf{s}) &= \min_{\mathbf{w}, \mathbf{l}, \mathbf{o}} \frac{1}{H} \sum_{h=1}^H \left[ \sum_{i=1}^n (c_i^w x_{i,j} W_{i,j}^h + c_j^l I_{i,j}^h) + c_j^o O_j^h \right] \\
& \text{s.t. } W_{i,j}^h - I_{i-1,j}^h = x_{i-1,j} d_{i-1,j}^h + W_{i-1,j}^h - s_{i-1,j} \quad i = 2, \dots, n, h = 1, \dots, H \\
& \quad O_j^h - I_{n,j}^h = x_{n,j} d_{n,j}^h + W_{n,j}^h - s_{n,j} \quad h = 1, \dots, H \\
& \quad W_{1,j}^h = 0 \quad h = 1, \dots, H
\end{aligned} \tag{4.6}$$

349 The above problems (4.5) and (4.6) called master problem (MP) and subproblem (SP), respec-  
350 tively. As we mentioned, with the optimal solution obtained from problem (4.5), the optimal cost  
351 and solution of problem (4.6) can be achieved through recursive equations (3.1), (3.2) and (3.3)  
352 without actually solving the optimization problem. Thus, the remaining problem is how to use the  
353 solution of the second stage problem to verify the optimality of the master problem.

354  
355 Next, we analyze the optimal solution of the dual problem of  $Q_j(\mathbf{x}, \mathbf{s})$ , which helps to find the  
356 optimal solution for the first problem (4.5). For each scenario  $h$ , let  $\alpha_{i,j}^h (i = 1, 2, \dots, n, h = 1, 2, \dots, H)$   
357 be the dual decision variable for the second problem (4.6), the dual form of the operational cost for  
358 service provider  $j$  under scenario  $h$  can be derived as follows:

$$\begin{aligned}
& \max_{\alpha} \sum_{j=1}^k \sum_{i=1}^n (x_{i,j} d_{i,j}^h - s_{i,j}) \alpha_{i,j}^h \\
& \text{s.t. } \alpha_{i-1,j}^h - \alpha_{i,j}^h \leq c_i^w x_{i,j} \quad i = 2, \dots, n, h = 1, \dots, H \\
& \quad -\alpha_{i,j}^h \leq c_j^l \quad i = 2, \dots, n, h = 1, \dots, H \\
& \quad \alpha_{n,j}^h \leq c_j^o \quad h = 1, \dots, H
\end{aligned} \tag{4.7}$$

359 By the strong dual theorem, we can obtain the optimal solution of the dual problem (4.7)  
360 without actually solve it:

$$\alpha_{i,j}^h = \begin{cases} \alpha_{i+1,j}^h + c_{i+1}^W x_{i+1,j} & \text{if } W_{i+1,j}^h > 0 \\ -c_j^I & \text{otherwise} \end{cases} \quad i = 1, \dots, n-1, h = 1, \dots, H \quad (4.8)$$

$$\alpha_{n,j}^h = \begin{cases} c_j^O & \text{if } O_j^h > 0 \\ -c_j^I & \text{otherwise} \end{cases} \quad h = 1, \dots, H \quad (4.9)$$

361 To solve the two-stage program (4.5) efficiently, we propose an improved Benders decomposition  
 362 algorithm based on the optimal solution of (4.8) and (4.9). Following the idea of general Benders  
 363 decomposition, the procedures of the Benders decomposition algorithm for our studied two-stage  
 364 program are as follows: (1) Formulate the master problem MP as a relaxation form by replacing  
 365 the optimal value  $Q_j(\mathbf{x}, \mathbf{s})$  with a new decision variable  $\theta_j$  ( $\theta_j \geq 0$ ). And then solve the new master  
 366 problem and find an optimal solution  $\mathbf{s}$  and sends it to the SP. (2) Evaluate whether the optimal  
 367 solution obtained in the new master problem violates the optimality. If it does, then add optimality  
 368 cuts generated by incorporating the optimal solution of the SP to the MP and go back to proce-  
 369 dure (1); otherwise, the solution is globally optimal. (3) Repeat the above two procedures until an  
 370 optimal solution is found.

371  
 372 In this work, we also improve the standard Benders decomposition by adding the feasibility cuts  
 373  $\theta_j \geq \sum_{i=1}^{n-1} x_{i+1,j} g_{i,j} + g_{n,j}$   $j = 1, \dots, n$ , which derives from Lemma 1. The set of feasibility cuts restrict  
 374 the feasible region, which helps to solve the problem more efficiently. As for the optimality cuts  
 375  $\{L(\mathbf{x}, \mathbf{s}) \geq 0\}$  that come from the optimal solution of the dual of the SP, due to the special structure of  
 376 SP, the optimal solution can be achieved easily without actually solving the optimization problem.

377 The pseudocode of the proposed improved Benders decomposition is presented in Algorithm 1.

378 Algorithm 1

379 Step 1: Input: Service time realization  $\mathbf{d}$ , parameters  $c^W, c^I, c^O$ ,  $\lambda$ , and  $\mathbf{m}$ . Set the set of  
 380 optimality cuts  $\{L(\mathbf{x}, \mathbf{s}) \geq 0\} = \emptyset$ .

Step 2: Solve the master problem

$$\begin{aligned} (MP) \quad & \min_{\mathbf{x}, \mathbf{s}, \theta} \sum_{j=1}^k \theta_j + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \\ & s.t. \sum_{j=1}^k x_{i,j} = 1 \quad i = 1, \dots, n \\ & \sum_{i=1}^k x_{i,j} \geq 1 \quad j = 1, \dots, n \\ & \sum_{i=1}^n s_{i,j} = T \quad j = 1, 2, \dots, k \\ & s_{i,j} \leq M x_{i,j} \quad i = 1, \dots, n, j = 1, \dots, k \\ & \theta_j \geq \sum_{i=1}^{n-1} x_{i+1,j} g_{i,j} + g_{n,j} \quad j = 1, \dots, n \\ & x_{i,j} \in \{0, 1\}, s_{i,j} \geq 0 \\ & L(\mathbf{x}, \mathbf{s}) \geq 0 \end{aligned}$$

381 and record an optimal solution  $(\mathbf{x}^*, \mathbf{s}^*, \theta^*)$ .

382 Step 3: Given  $(\mathbf{x}^*, \mathbf{s}^*)$  obtained from above MP, calculate the corresponding performance indi-  
 383 cators  $\mathbf{W}, \mathbf{I}$  and  $\mathbf{O}$  by recursive equations (3.1), (3.2) and (3.3). And record  $\sum_{j=1}^k Q_j^*(\mathbf{x}^*, \mathbf{s}^*)$ . On  
 384 this basis, determine the optimal solution  $(\alpha)$  to problem 4.7 for each scenario through recursive  
 385 equations (4.8) and (4.9).

386 Step 4: If  $\sum_{j=1}^k \theta_j^* \geq \sum_{j=1}^k Q_j^*(\mathbf{x}^*, \mathbf{s}^*)$ , then  
 387 stop and return  $(\mathbf{x}^*, \mathbf{s}^*, \theta^*)$  as an optimal solution.

388 Else

389 add the set of cuts  $\theta_j \geq \frac{1}{H} \sum_{h=1}^H [\sum_{i=1}^n d_{i,j}^h \alpha_{i,j}^h x_{i,j} - \sum_{i=1}^n \alpha_{i,j}^h s_{i,j}]$  ( $j = 1, \dots, k$ ), to the set of  
 390 optimality cuts  $\{L(\mathbf{x}, \mathbf{s}) \geq 0\}$ . And go to Step 2.

391 End if.

## 392 5. Incorporating No-shows

393 In this section, the proposed method is extended to solve the matching and appointment schedul-  
 394 ing problem by considering no-shows. Let  $p_i$  denote the show up probability for patient  $i$ , which is  
 395 known to the decision-maker. Let  $z_i$  indicate whether patient  $i$  shows up for her appointment (i.e.,  
 396  $z_i = 1$  with probability  $p_i$ ) or not (i.e.,  $z_i = 0$  with probability  $1 - p_i$ ). And we also assume that the  
 397 no-shows are independent for patients.

398 The key idea is to treat the no-show patient as a “ghost” patients with 0 service time. Let  $\tilde{d}_{i,j}$   
 399 denote the service time in the presence of no-shows. We can calculate the actual service time in  
 400 the system through  $\tilde{d}_{i,j} = z_i d_{i,j}$ , where  $d_{i,j}$  is the service time without no-shows studied previously.  
 401 By abusing notations, let  $W_{i,j}, I_{i,j}$ , and  $O_j$  denote the corresponding virtual waiting time, idle time,  
 402 and over time, respectively. Then, we have

$$\begin{aligned} W_{i,j} &= \max\{0, x_{i-1,j} \tilde{d}_{i-1,j} + W_{i-1,j} - s_{i-1,j}\} \quad i = 2, \dots, n, j = 1, \dots, k \\ I_{i-1,j} &= \max\{0, -x_{i-1,j} \tilde{d}_{i-1,j} - W_{i-1,j} + s_{i-1,j}\} \quad i = 2, \dots, n+1, j = 1, 2, \dots, k \\ O_j &= \max\{0, x_{n,j} \tilde{d}_{n,j} + W_{n,j} - s_{n,j}\} \quad j = 1, 2, \dots, k \\ W_{1,j} &= 0 \quad j = 1, \dots, k \end{aligned} \quad (5.1)$$

403 In the presence of no-shows, if one patient is a no-show, we can regard the actual waiting time  
 404 for him/her as zero. Thus, in the objective function, we only need to count the waiting time of  
 405 those patients who actually show up. On this basis, the total cost under the no-show case is

$$\sum_{j=1}^k \mathbb{E} \left[ \sum_{i=1}^n (c_i^W x_{i,j} z_i W_{i,j} + c_j^I I_{i,j}) + c_j^O O_j \right] + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \quad (5.2)$$

406 Note that  $z_i$  is independent of  $(z_1, \dots, z_{i-1})$ , it must be independent of  $W_{i,j}$ . Thus,  $\mathbb{E} [c_i^W x_{i,j} z_i W_{i,j}] =$   
 407  $p_i \mathbb{E} [c_i^W x_{i,j} W_{i,j}]$ . Let  $\tilde{c}_i^W = p_i c_i^W$ , then the above equation (5.2) is equivalent to

$$\sum_{j=1}^k \mathbb{E} \left[ \sum_{i=1}^n (\tilde{c}_i^W x_{i,j} W_{i,j} + c_j^I I_{i,j}) + c_j^O O_j \right] + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \quad (5.3)$$

408 Note that equation (5.3) has the same form as equation (3.4), except the notations  $\tilde{c}_i^W$  and  $\tilde{d}_{i,j}$ .  
 409 Therefore, the proposed method in section 4 can be applied for the case with no-shows.

410

411 6. Numerical Analyses

412 In this section, we conduct numerical experiments to evaluate the performance of our proposed  
 413 algorithm, and study the influence of different parameters on the optimal assignment and schedule.  
 414 Specifically, we intend to compare the computational time between our proposed method and the  
 415 benchmark. The benchmark is solving the corresponding MILP through CPLEX directly, which  
 416 we will introduce later. Moreover, we investigate the optimal solutions under different scenarios  
 417  $(\lambda, (n, k))$  and the effect of no-shows. We assume an i.i.d. normal distribution for the service time,  
 418 i.e.,  $d_{i,j} \sim N(\mu, \sigma^2), i = 1, \dots, n, j = 1, \dots, k$ , which has been widely used in the appointment scheduling  
 419 literature [11, 24]. Following the literature[21, 20, 36], we set the unit waiting time, idling time and  
 420 overtime costs as  $c_i^W = 0.2, c_j^I = 1, c_j^O = 1.5$ .

421 6.1. Performance of Improved Decomposition Algorithm

422 In this subsection, we study the performance of our proposed decomposition algorithm. We  
 423 first reformulate the original deterministic problem (*DP*) as a deterministic mixed-integer linear  
 424 program (*MILP*). And then we solve the MILP directly with CPLEX and use the corresponding  
 425 results as the benchmark. Through the big-M transformation, the deterministic model (*DP*) can  
 426 be reformulated as the following mixed-integer linear program:

$$\begin{aligned}
 \text{(DMILP)} \quad & \min_{\mathbf{x}, \mathbf{s}, \mathbf{W}, \bar{\mathbf{W}}, \mathbf{I}, \mathbf{O}} \frac{1}{H} \sum_{h=1}^H \sum_{j=1}^k \left[ \sum_{i=1}^n (c_i^W \bar{W}_{i,j}^h + c_j^I I_{i,j}^h) + c_j^O O_j^h \right] + \lambda \sum_{i=1}^n \sum_{j=1}^k m_{i,j} x_{i,j} \\
 \text{s.t.} \quad & \sum_{i=1}^n x_{i,j} \geq 1 \quad j = 1, 2, \dots, k \\
 & \sum_{j=1}^k x_{i,j} = 1 \quad i = 1, 2, \dots, n \\
 & \sum_{i=1}^n s_{i,j} = T \quad j = 1, 2, \dots, k \\
 & s_{i,j} \leq M x_{i,j} \quad i = 1, \dots, n, j = 1, \dots, k \\
 & W_{i,j}^h - I_{i-1,j} = x_{i-1,j} d_{i-1,j}^h + W_{i-1,j}^h - s_{i-1,j} \quad i = 2, \dots, n, j = 1, \dots, k, h = 1, \dots, H \\
 & O_j^h - I_{n,j}^h = x_{n,j} d_{n,j}^h + W_{n,j}^h - s_{n,j} \quad j = 1, \dots, k, h = 1, \dots, H \\
 & \bar{W}_{i,j}^h \geq W_{i,j}^h + (x_{i,j} - 1)M \quad i = 2, \dots, n, j = 1, \dots, k, h = 1, \dots, H \\
 & \bar{W}_{i,j}^h \leq W_{i,j}^h \quad i = 2, \dots, n, j = 1, \dots, k, h = 1, \dots, H \\
 & W_{1,j}^h = 0 \quad j = 1, \dots, k, h = 1, \dots, H \\
 & x_{i,j} \in \{0, 1\}, s_{i,j} \geq 0
 \end{aligned} \tag{6.1}$$

427 Throughout this section, we randomly generate  $H = 1000$  i.i.d. realizations based on given  
 428 service time distribution. And then we solve the MILP directly with CPLEX. Finally, we illustrate  
 429 the superiority of our method by comparing the running time with the benchmark.

430 The parameters are presented as follows.



- 431 • The service times of all jobs at each server  $d_{i,j}$  follow a normal distribution  $N(20, 16)$ .
- 432 • The matching cost  $m_{i,j}$  of patient  $i$  for physician  $j$  is generated by a uniform distribution  
433  $U[0, 1]$ .
- 434 • The number of patients and physicians appear as pairs  $(n, k)$ . We consider four pairs, i.e.,  
435  $(20, 2)$ ,  $(30, 3)$ ,  $(40, 4)$  and  $(50, 5)$ .
- 436 • The session length is set at  $T = 1.5 \cdot \mu \cdot n/k$ ;

437 For each pair  $(n, k)$ , we generate 5 problem instances, for each problem instance, the  $\lambda$  is randomly  
438 generated from interval  $[10, 1000]$ . As a result, there are a total of 20 problem instances in our  
439 computational experiments. All instances are solved by calling CPLEX 12.6 on Matlab R2016a that  
440 runs on a PC with an Intel i5-4590 CPU and 12 GB memory. In our computational experiments, we  
441 set the limit of computational times to 5 hours (i.e., 18,000 seconds). For our developed algorithm,  
442 we set the absolute tolerance at 0.01. Because the optimal objective value from our method is almost  
443 the same (the absolute tolerance is within 0.1) with that from the benchmark (those instances can  
444 be solved within 5 hours), we do not display the optimal value in this work. We compare the  
445 average, minimum and maximum computational times, which are summarized in Table 1.

Table 1: Comparison on the computational time between our method and benchmark (in second)

patients-physician pairs	Our method			Benchmark		
	Min. time	Avg. time	Max. time	Min. time	Avg. time	Max. time
(20,2)	67	76	90	854	947	1,042
(30,3)	107	149	196	4,670	7,127	15,374
(40,4)	181	213	269	15,929	18,000	18,000
(50,5)	278	320	369	18,000	18,000	18,000
overall	158	189	231	9,863	11,018	13,104

446 As shown in Table 1, the average, minimum and maximum computational times of the proposed  
447 method are significantly shorter than the benchmark. Computational time increases with the  
448 increasing of  $n$  and  $k$  in pair  $(n, k)$  for both our method and the benchmark. As  $(n, k)$  increases, the  
449 computational time of our method increases slowly, while that of the benchmark increases rapidly.  
450 When  $(n, k)$  reaches  $(40, 4)$  and  $(50, 5)$ , the benchmark cannot obtain the optimal solution before  
451 reaching the time limit (i.e., Ave.time =18,000), while our method can solve all problem instances  
452 optimality in a reasonable computational time. These facts indicate that our method is indeed  
453 much more efficient than the benchmark. Moreover, there is no significant difference when  $\lambda \geq 100$   
454 according to the experimental results. Therefore, we test  $\lambda$  ranged from 0 to 100 in the following  
455 sections.

#### 456 6.2. Analysis of $\lambda$ in different scenarios

457 In this subsection, we analyze the effect of  $\lambda$  in different scenarios on the optimal solution and  
458 the objective value. Specifically, we fix the pair of patients and physicians as  $(20, 2)$  and test differ-  
459 ent values of matching cost  $\mathbf{m}$  and  $\lambda$ . During the numerical experiments, we construct four different  
460 scenarios for the matching between patients and physicians, i.e., the number of patients who match

461 the first physician best is 2, 4, 6, 8, respectively, which we use 1 : 9, 2 : 8, 3 : 7 and 4 : 6 to represent  
 462 those four scenarios. We refer to scenario 1 : 9 as the most imbalanced patients-physician scenario.  
 463 As for the matching cost  $\mathbf{m}$ , we randomly generate from  $U[0.1, 0.3]$  if one patient match with the  
 464 best-matched physician; otherwise, we randomly generate from  $U[0.8, 1]$ . Other parameters setting  
 465 are the same as subsection 6.1.

466  
 467 We first study the variant of the total cost in different scenarios as  $\lambda$  increases. As shown in  
 468 Figure 1, all matching patterns appear an increasing trend as  $\lambda$  increases. Furthermore, the more  
 469 imbalanced of patients-physician matching, the higher total the cost it has. This because the im-  
 470 balanced matching pattern would result in a higher matching cost to the optimal solution, and  
 471 further amplify the total cost through  $\lambda$ .

472

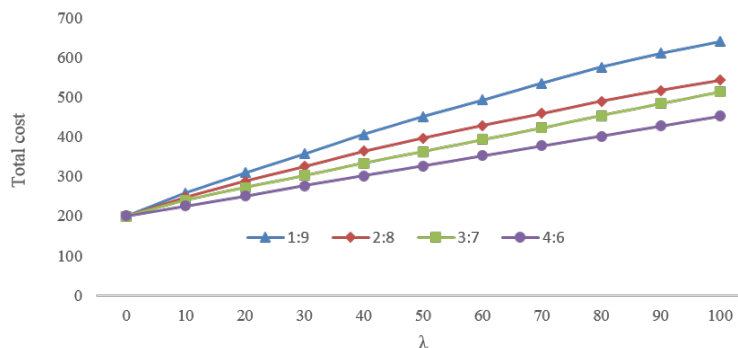


Figure 1: Comparison of the optimal total cost

473 Then, we break down the total cost and further study the variant of its corresponding opera-  
 474 tional and unit matching costs. The operational cost refers to the total expected weighted costs of  
 475 patients' waiting times and service providers' idle times and overtimes under the optimal solution,  
 476 i.e.,  $\sum_{j=1}^k \mathbb{E} \left[ \sum_{i=1}^n (c_i^W x_{i,j}^* W_{i,j} + c_j^I I_{i,j}) + c_j^O O_j \right]$ . The unit matching cost refers to the total matching  
 477 cost without weighted by  $\lambda$  under the optimal solution, i.e.,  $\sum_{j=1}^k \sum_{i=1}^n m_{i,j} x_{i,j}^*$ . As shown in Figure 2,  
 478 both the operational cost and unit matching cost of the scenario 4 : 6 remains stable as  $\lambda$  changes.  
 479 However, for other matching scenarios, the operational costs increase as  $\lambda$  increases, while the unit  
 480 matching costs exhibit an opposite trend. The reason is that when  $\lambda$  is large, the optimal assign-  
 481 ment would result in an unbalance operational pattern to avoid a considerable matching cost, thus  
 482 leads to higher operational costs. But for the more balanced matching scenario (e.g., scenario 4 : 6),  
 483 the weight  $\lambda$  has no significant impact on the assignment. Besides, both the operational cost and  
 484 unit matching cost witness a more significant fluctuation when the matching scenario is unbalanced  
 485 (e.g., scenario 1 : 9). However, for the operational cost, it seems that the more stable matching  
 486 scenario, the less operational cost it has. The difference of operational costs between scenarios 1 : 9  
 487 and 4 : 6 goes up to 51%, while the difference of unit matching costs between them goes up to 54%.

488

489 Moreover, we study the pattern of patients assigned to the first physician as  $\lambda$  changes for dif-  
 490 ferent matching scenarios. As shown in Figure 3, we can see that the number of patients assigned

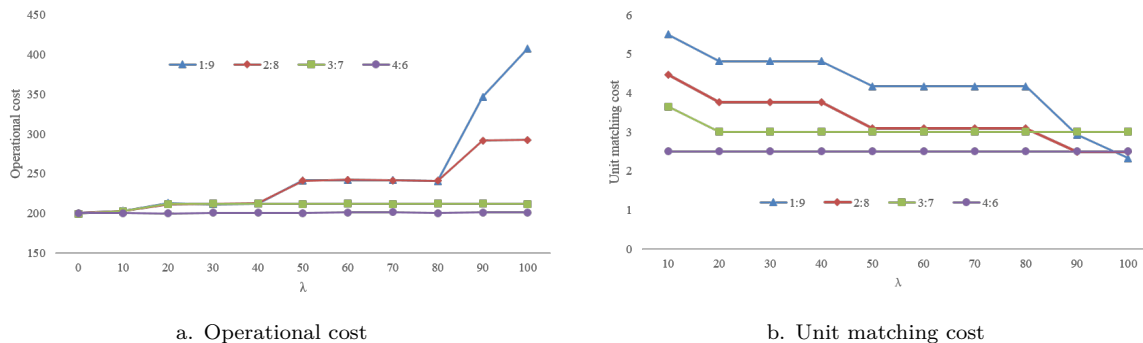


Figure 2: Comparison of the operation cost and the unit matching cost

491 to the first physician exhibits a non-increasing trend as  $\lambda$  increases for all four matching scenarios.  
 492 This result indicates that when the weight of matching is large, the optimal solution would assign  
 493 more patients to the most matched physician. Otherwise, the optimal solution would balance the  
 494 operational workload for each physician. Furthermore, when the scenario is more balanced (i.e.,  
 495 scenario 4 : 6),  $\lambda$  is not a significant factor to the patient assignment. Another interesting ob-  
 496 servation is that there exist some overlaps for four scenarios and scenario 1 : 9 overlaps all other  
 497 scenarios. It may due to scenario 1 : 9 is the most imbalanced case. When  $\lambda$  increases, patients  
 498 tend to be assigned to the first physician, such that a higher matching cost is avoided.  
 499

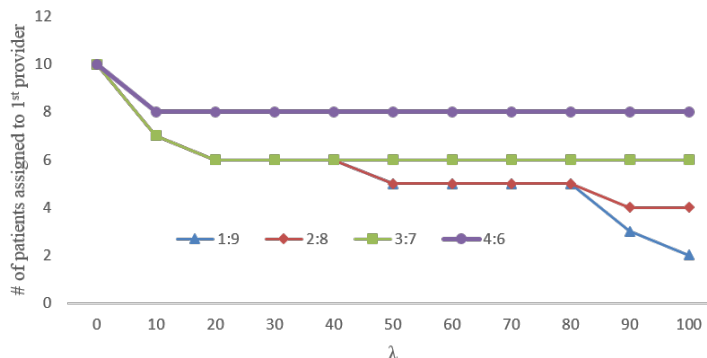


Figure 3: Comparison of the number of patients assigned to the first physician

### 500 6.3. Analyses of $\lambda$ with no-shows

501 We further analyze the effect of no-shows on the optimal solution and the objective value. We  
 502 assume an i.i.d. show up indicator, i.e.,  $p_i = p$ , and set the show up probability  $p$  takes values  
 503 from 0.6 to 0.9 with an increment of 0.1. By considering the no-shows, we reset the session length  
 504 as  $T = 1.5\mu p n/k$ , such that it can adapt according to patients' no-show behaviors. From Figure 2a  
 505 , we observe a more substantial fluctuation of the operational cost when the matching scenario is

506 unbalanced. Thus, we set  $(n, k)$  at  $(20, 2)$  and choose scenarios 1 : 9 to analyze objective value and  
 507 performance indicators. Moreover, there is a small variation of the operational cost when  $\lambda \leq 50$ .  
 508 Therefore, to obtain an apparent comparison result for performance indicators, we test two different  
 509 values of  $\lambda$  (i.e.,  $\lambda = 50, 100$ ), and examine the impact of  $\lambda$  on performance indicators. The other  
 510 parameters (e.g.,  $m, d, c_i^W, c_j^I, c_j^O$ ) are the same with scenario 1 : 9 in section 6.2.  
 511

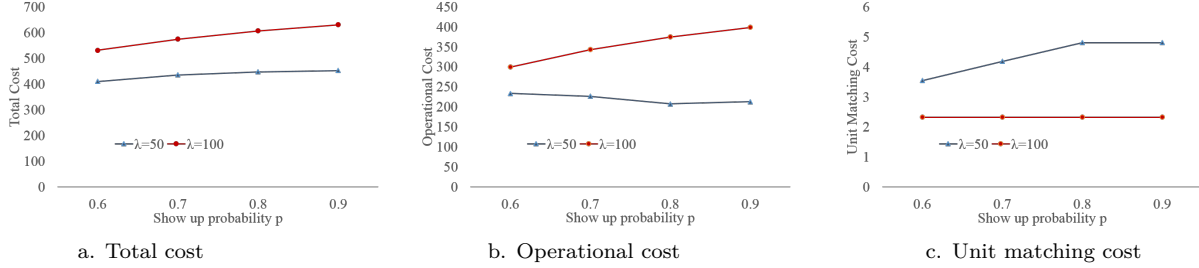


Figure 4: Comparison of the optimal cost with no-shows

512 The optimal total, operational, and unit matching costs are presented in Figure 4. As the same  
 513 as we discussed in section 6.2,  $\lambda$  has the a positive impact on the total and operational costs and  
 514 negative impact on the unit matching cost. It is intuitive that the total cost increases as the show  
 515 up probability  $p$  increase, as shown in Figure 4a. We also observe the same effect of the show up  
 516 probability on the unit matching cost when  $\lambda = 50$ . However, when  $\lambda$  is substantial (e.g.,  $\lambda = 100$ ),  
 517 the show up probability has no influence on the unit matching cost, as shown in Figure 4c. It is  
 518 more surprising to observe that the operational cost has an increasing trend when  $\lambda = 100$  while an  
 519 decreasing trend when  $\lambda = 50$ . Therefore, we decompose the operational cost and further analyze  
 520 the total waiting times, idle times, and overtimes, as shown in Figure 5.

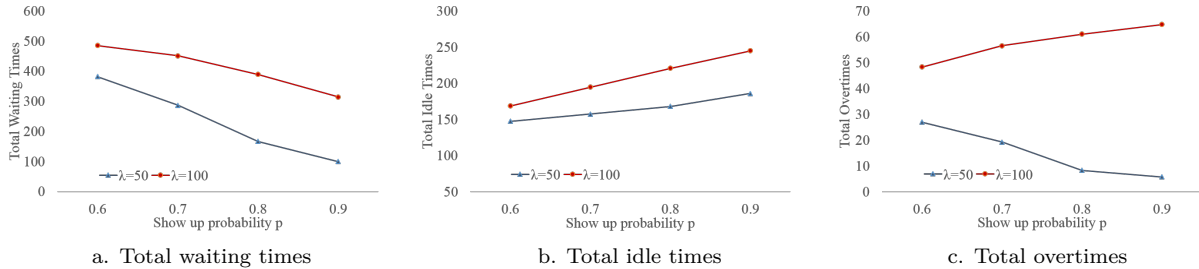


Figure 5: Comparison of operational indicators

521 From Figures 5a and 5b, we can see that the total waiting times under both values of  $\lambda$  decrease  
 522 as the show up probability  $p$  increases, while the total idle times witness an opposite trend. This  
 523 phenomenon can be explained as follows. In each configuration combination of  $(\lambda, p)$ , the average  
 524 job allowance for each patient can be approximated as  $1.5\mu = 30$  roughly (Note that we set the  
 525 session length as  $T = 1.5 \cdot \mu \cdot p \cdot n/k$ ). However, the average service time is  $\mu = 20$ . It would be  
 526 more idleness than waiting for the service system. Thus, as the show up probability increases, the  
 527 number of show up patients increases, which leads to the decreasing of waiting times and increasing  
 528 of idle times. This kind of trend is amplified as  $\lambda$  increase. Furthermore, from Figure 5c, we can

529 observe that the total overtimes decrease when  $\lambda = 50$  and increase when  $\lambda = 100$  as the show up  
530 probability increases. The server system has a higher probability to suffer an overload when  $\lambda = 100$   
531 than  $\lambda = 50$ . Since the average job allowance is considerably longer than the average service time,  
532 the more show up patients, the more idleness, thus less overtimes when  $\lambda = 50$ . However, when 533  $\lambda =$   
100, the optimal assignment has a high chance to overload for one server. Therefore, more 534 overtime  
may occur as the number of show up patients increases. As a result, when  $\lambda = 100$ , the 535 total overtime  
increases as the show up probability  $p$  increases.

## 536 7. Managerial Implications

537 In this section, we summarize some managerial insights from the numerical results of this study.  
538 First, a more balanced supply and demand result in a lower total cost. From Figure 1, we can  
539 observe that for any fixed  $\lambda$ , the more imbalanced scenario would result in a higher total cost.  
540 This result indicates that if we want to achieve a lower total cost for the system, a more balanced  
541 scenario, i.e., scenario 4 : 6, is more desirable. Furthermore, as  $\lambda$  increases, the total cost increase  
542 for all scenarios. We also can observe that the more unbalanced scenario, i.e., scenario 1 : 9, the  
543 faster-increasing speed of total cost. This result implies that when matching becomes more impor-  
544 tant, the service provider should pay more effort to balance physicians' specialties and patients'  
545 diseases.

546  
547 Second, the weight of the matching cost has a positive impact on the operational cost. From  
548 Figure 2, we can see that there is no significant difference for the operational costs among different  
549 scenarios when  $\lambda$  is small (e.g.,  $\lambda \leq 10$ ). As  $\lambda$  increases, the marginal operational cost increases  
550 while the marginal unit matching cost decreases. The operational cost increases at a swift speed  
551 as matching becomes more critical. It is because the overtime cost of matched physicians and  
552 waiting time of patients increases significantly while other mismatched physicians are idle for the  
553 unbalanced scenario, i.e., scenario 1 : 9.

554  
555 Third, we would suggest the service provider to develop or train the specialty set based on the  
556 disease pattern of local patients, such that the workload among physicians can be more balanced.  
557 From Figure 3 we can see that when the weight  $\lambda$  is large (e.g.,  $\lambda = 100$ ), the workload among  
558 different service providers is imbalanced, which means there may be a waste of resources to some  
559 degree. Therefore, the specialty set of physicians is critical to balance the physician workload.

560  
561 Fourth, when the weight of the matching cost is not substantial, we would suggest penalizing  
562 the no-show patient, such that the operational cost is minimized. However, when the weight of the  
563 matching cost is considerable, which means the matching is more critical to the health care quality,  
564 the patient has a higher motivation to see the matched physician. Therefore, the no-show has less  
565 influence on the service system.

## 566 8. Conclusion and future work

567 In this paper, we jointly optimize a matching problem and appointment problem with multiple  
568 service providers, in which the decision-maker determines how to assign patients to physicians and  
569 when to start serving patients for each service provider. We assume the service times are stochastic

570 and all patients would arrive at the service system punctually at their scheduled times. The ob-  
571 jective is to minimize total weighted matching costs and operational costs (patients' waiting time  
572 costs, service providers' idle time and overtime costs). To solve the problem, we first reformulate  
573 the studied problem as a two stage optimization problem based on the SAA approach. And then  
574 the properties for the optimal solution of the second stage problem are analyzed. On this basis,  
575 an improved benders decomposition algorithm is proposed to solve this problem efficiently. We  
576 also extend our method to incorporate no-shows. Finally, we conduct computational experiments  
577 to evaluate the efficiency of our proposed algorithm and investigate the variation of the optimal  
578 solutions yielded in different scenarios.

579  
580 Our studies mainly show that: (1) the integrated matching problem and appointment scheduling  
581 can be formulated as a two-stage optimization problem; (2) the improved Benders decomposition  
582 algorithm is efficient to solve our studied problem in a reasonable time; (3) the optimal solution  
583 would assign patients to the most matched physicians if the matching cost dominate operational  
584 cost, otherwise, the optimal solution would balance the workload of physicians as much as possible;  
585 and (4) the no-show has less influence on the service system when the weight of the matching cost  
586 is substantial.

587  
588 Our work can be extended in several aspects. We assume a fixed sequence for patients at each  
589 service provider, due to different type of patients, it may be valuable to optimize the sequence for  
590 patients at each service provider. Furthermore, unpunctuality is also inevitable in practice, which  
591 may lead to patients arriving out of order. In the future, we may also handle these stochastic  
592 factors.

#### 593 Acknowledgment

594 The first author is supported by the the Science and Technology Foundation of Jiangxi Educa-  
595 tional Committee (#GJJ190287). The second author is funded by NSFC #71801051, and the third  
596 author is funded by Omron research fund.

#### 597 References

- 598 [1] M. Weiner, G. El Hoyek, L. Wang, P. R. Dexter, A. D. Zerr, A. J. Perkins, F. James, R. Juneja,  
599 A web-based generalist–specialist system to improve scheduling of outpatient specialty consul-  
600 tations in an academic center, *Journal of General Internal Medicine* 24 (6) (2009) 710–715.
- 601 [2] A. Mehrotra, C. B. Forrest, C. Y. Lin, Dropping the baton: specialty referrals in the united  
602 states, *The Milbank quarterly* 89 (1) (2011) 39–68.
- 603 [3] D. Li, S. Chen, X. Chen, C.-A. Chou, Learning and optimizing for the patient-  
604 physician matching problem in specialty care, *SSRN* (2019) Available at SSRN:  
605 <https://ssrn.com/abstract=3450184>.
- 606 [4] E. D. Güneş, H. Yaman, B. Çekyay, V. Verter, Matching patient and physician preferences in  
607 designing a primary care facility network, *Journal of the Operational Research Society* 65 (4)  
608 (2014) 483–496.

- 609 [5] R. Crow, H. Gage, S. Hampson, J. Hart, A. Kimber, L. Storey, H. Thomas, The measure-  
610 ment of satisfaction with healthcare: implications for practice from a systematic review of the  
611 literature., *Health Technology Assessment* 6 (32) (2002) 1–244.
- 612 [6] R. Beck, R. Daughtridge, P. Sloane, Physician–patient communication in the primary care  
613 office: a systematic review., *The Journal of the American Board of Family Medicine* 15 (2002)  
614 25–38.
- 615 [7] J. Gong, H. Cheng, L. Wang, Individual doctor recommendation in large networks by con-  
616 strained optimization, *International Journal of Web Services Research* 12 (2015) 16–28.
- 617 [8] N. Liu, S. R. Finkelstein, M. E. Kruk, D. Rosenthal, When waiting to see a doctor is less  
618 irritating: Understanding patient preferences and choice behavior in appointment scheduling,  
619 *Management Science* 64 (5) (2018) 1975–1996.
- 620 [9] K. Kinchen, L. Cooper, D. Levine, N.-Y. Wang, N. Powe, Referral of patients to specialists:  
621 Factors affecting choice of specialist by primary care physicians, *Annals of family medicine* 2  
622 (2004) 245–52.
- 623 [10] X. Pan, J. Song, F. Zhang, Dynamic recommendation of physician assortment with patient  
624 preference learning, *IEEE Transactions on Automation Science and Engineering* 16 (1) (2019)  
625 115–126.
- 626 [11] B. Denton, D. Gupta, A sequential bounding approach for optimal appointment scheduling,  
627 *IIE Transactions* 35 (11) (2003) 1003–1016.
- 628 [12] S. De Vuyst, H. Bruneel, D. Fiems, Computationally efficient evaluation of appointment sched-  
629 ules in health care, *European Journal of Operational Research* 237 (3) (2014) 1142–1154.
- 630 [13] E. N. Weiss, Models for determining estimated start times and case orderings in hospital  
631 operating rooms, *IIE transactions* 22 (2) (1990) 143–150.
- 632 [14] H.-S. Lau, A. H.-L. Lau, A fast procedure for computing the total system cost of an appoint-  
633 ment schedule for medical and kindred facilities, *IIE Transactions* 32 (9) (2000) 833–839.
- 634 [15] H.-Y. Mak, Y. Rong, J. Zhang, Sequencing appointments for service systems using inventory  
635 approximations, *Manufacturing & Service Operations Management* 16 (2) (2014) 251–262.
- 636 [16] A. Kuiper, B. Kemper, M. Mandjes, A computational approach to optimized appointment  
637 scheduling, *Queueing Systems* 79 (1) (2015) 5–36.
- 638 [17] R. Hassin, S. Mendel, Scheduling arrivals to queues: A single-server model with no-shows,  
639 *Management Science* 54 (3) (2008) 565–572.
- 640 [18] Q. Kong, C.-Y. Lee, C.-P. Teo, Z. Zheng, Scheduling arrivals to a stochastic service delivery  
641 system using copositive cones, *Operations Research* 61 (3) (2013) 711–726.
- 642 [19] H.-Y. Mak, Y. Rong, J. Zhang, Appointment scheduling with limited distributional informa-  
643 tion, *Management Science* 61 (2) (2014) 316–334.
- 644 [20] T. Cayirli, K. K. Yang, S. A. Quek, A universal appointment rule in the presence of no-shows  
645 and walk-ins, *Production and Operations Management* 21 (4) (2012) 682–697.

- 646 [21] L. W. Robinson, R. R. Chen, A comparison of traditional and open-access policies for appointment  
647 scheduling, *Manufacturing & Service Operations Management* 12 (2) (2010) 330–346.
- 648 [22] C. Zacharias, M. Pinedo, Appointment scheduling with no-shows and overbooking, *Production  
649 and Operations Management* 23 (5) (2014) 788–801.
- 650 [23] M. A. Begen, R. Levi, M. Queyranne, A sampling-based approach to appointment scheduling,  
651 *Operations Research* 60 (3) (2012) 675–681.
- 652 [24] L. W. Robinson, R. R. Chen, Scheduling doctors’ appointments: optimal and empirically-based  
653 heuristic policies, *IIE Transactions* 35 (3) (2003) 295–307.
- 654 [25] G. C. Kaandorp, G. Koole, Optimal outpatient appointment scheduling, *Health Care Manage-  
655 ment Science* 10 (3) (2007) 217–229.
- 656 [26] B. Jiang, J. Tang, C. Yan, A stochastic programming model for outpatient appointment  
657 scheduling considering unpunctuality, *Omega* 82 (2019) 70–82.
- 658 [27] Y. Chen, Y.-H. Kuo, P. Fan, H. Balasubramanian, Appointment overbooking with different  
659 time slot structures, *Computers & Industrial Engineering* 124 (2018) 237–248.
- 660 [28] Y.-H. Kuo, H. Balasubramanian, Y. Chen, Medical appointment overbooking and optimal  
661 scheduling: tradeoffs between schedule efficiency and accessibility to service, *Flexible Services  
662 and Manufacturing Journal* (2019) 1–30.
- 663 [29] R. R. Chen, L. W. Robinson, Sequencing and scheduling appointments with potential call-in  
664 patients, *Production and Operations Management* 23 (9) (2014) 1522–1538.
- 665 [30] G. Xiao, M. Dong, J. Li, L. Sun, Scheduling routine and call-in clinical appointments with  
666 revisits, *International Journal of Production Research* 55 (6) (2017) 1767–1779.
- 667 [31] I. Bendavid, Y. N. Marmor, B. Shnits, Developing an optimal appointment scheduling for  
668 systems with rigid standby time under pre-determined quality of service, *Flexible Services and  
669 Manufacturing Journal* 30 (1-2) (2018) 54–77.
- 670 [32] H.-J. Alvarez-Oh, H. Balasubramanian, E. Koker, A. Muriel, Stochastic appointment schedul-  
671 ing in a team primary care practice with two flexible nurses and two dedicated providers,  
672 *Service Science* 10 (3) (2018) 241–260.
- 673 [33] B. Zheng, S. W. Yoon, M. T. Khasawneh, et al., An overbooking scheduling model for outpa-  
674 tient appointments in a multi-provider clinic, *Operations Research for Health Care* 6 (2015)  
675 1–10.
- 676 [34] S. Sickinger, R. Kolisch, The performance of a generalized bailey–welch rule for outpatient  
677 appointment scheduling under inpatient and emergency demand, *Health care management  
678 science* 12 (4) (2009) 408.
- 679 [35] M. Soltani, M. Samorani, B. Kolfal, Appointment scheduling with multiple providers and  
680 stochastic service times, *European Journal of Operational Research* 277 (2) (2019) 667–683.
- 681 [36] C. Zacharias, M. Pinedo, Managing customer arrivals in service systems with multiple identical  
682 servers, *Manufacturing & Service Operations Management* 19 (4) (2017) 639–656.



- 683 [37] C. Mancilla, R. Storer, A sample average approximation approach to stochastic appointment  
684 sequencing and scheduling, *IIE Transactions* 44 (8) (2012) 655–670.
- 685 [38] E. N. Weiss, Models for determining estimated start times and case orderings in hospital  
686 operating rooms, *IIE Transactions* 22 (2) (1990) 143–150.