1
2

# Capacitated preventive health infrastructure planning with accessibility-based service equity

3   Hongzhi Lin, Ph.D.[1]; Min Xu, Ph.D.[2]

4   **Abstract**

5   Hard-to-access health infrastructure is likely to lead to increased morbidity and
6   mortality. The optimal layout of health facilities is undoubtedly of great significant for
7   disease control and prevention. This study aims to propose a method to provide
8   equitable access to capacitated preventive health facilities, which captures the key
9   features of facility congestion in a competitive choice environment. The problem is
10  formulated as a bilevel non-linear integer programming model. The upper level is a bi-
11  objective programming model subject to investment budget constraint, where the
12  primary objective is to minimize the maximum probability of balking (i.e., denied to
13  access service) and the secondary objective is to minimize the maximum queueing time.
14  The lower level is a user equilibrium analogous model resulting from the user choice
15  of facility location. It determines the allocation of users to facilities by a defined
16  generalized cost. An efficient heuristic algorithm is designed according to the bilevel
17  structure where the genetic algorithm (GA) with elite strategy is developed to solve the
18  upper level problem and the method of successive averages (MSA) is adopted to solve
19  the lower level problem. An illustrative case study is employed to validate the
20  performance of the proposed methods, and a number of interesting results and
21  managerial insights are provided with sensitivity analysis.

22  ***Keywords***: *health services, queueing, facility location, bilevel programming, user*
23  *equilibrium*

24

25  **1. Introduction**

26      The health infrastructure planning is an essential part of urban planning. It usually
27  means the planning of hospitals while the planning of preventive health facilities is

[1]Associate Professor, School of Economics and Management, Southeast University, Nanjing 211189, China. ORCID: https://orcid.org/0000-0002-4766-2621. Email:linhz@seu.edu.cn

[2]Assistant Professor, Department of Industrial and Systems Engineering, The Hong Kong Polytechnic University, Hung Hom, Hong Kong (Corresponding author). Email: xumincee@gmail.com

28  usually neglected. However, preventive health service, such as screening, examination,

29  isolation, and vaccination, is necessary for urban development. It is of utmost

30  importance since it can make massive savings on health expenditure by early detection.

31  This is a painful lesson from the raging COVID-19 pandemic. In fact, before the current

32  COVID-19 pandemic, three historically important epidemics had occurred since 2000:

33  severe acute respiratory syndrome (SARS) in 2003, Middle East respiratory syndrome

34  (MERS) in 2013, and Ebola virus disease in 2014. The arising monkeypox virus has

35  already attracted our attention. The health issue is a real problem disturbing urban

36  development. If the diseases can be detected and controlled earlier, the society would

37  not suffer from huge economic damage and life losses. Therefore, the authorities around

38  the world begin to realize the importance of preventive health facilities. In fact, the

39  users usually face barriers in accessing appropriate, timely, and affordable preventive

40  health service so far. The planning of preventive health infrastructure for disease control

41  and prevention is an urgent problem that has practical implication for urban planning

42  community.

43  A noticeable disparity in the accessibility to health facilities among different zones,

44  however, is found in theory and practice. This paper tries to propose a method to design

45  a health facility network for disease prevention, with the aim of improving service

46  equity in terms of accessibility. The accessibility usually refers to a measure of the ease

47  of reaching destinations or activities distributed in space. There are various ways to

48  measure the spatial accessibility to facilities. Unlike the conventional definition of

49  accessibility, the implication of accessibility here is straightforward and intuitive which

50  is defined as the accessible demand. In fact, there are two sources of inaccessible

51  demand: one is demand lost due to insufficient coverage and the other is demand lost

52  due to congested facility (Abouee-Mehrizi et al., 2011; Berman et al., 2006). For the

53  first source, demand is elastic with respect to cost and customers are usually assigned

54  to the closest facilities to maximize system total demand (Berman and Drezner, 2006;

55  Davari et al., 2016; Marianov, 2003; Zhang et al., 2010). For the second source,

56  customers could be denied to access the service (i.e., occurrence of balking) upon their

57  arrival, due to capable space. It is seldom explored in past studies because the

58  incorporation of limited capacity is not easy. There are only two closely related

59  references to the best of our knowledge. Marianov et al. (2008) studied the capacitated

60  facility network design problem. They defined travel cost as the travel time and

61  queueing delay but ignored the balking cost. This creates to the paradoxical situation

62    where customers will choose facilities with a greater likelihood of balking, because it
63    will reduce their overall time spent in the system while ignore the inaccessible demand.
64    Motivated by this, Dan and Marcotte (2019) defined user utility considering additional
65    balking cost and formulated a model to maximize the overall accessible demands.
66    However, although system accessibility is maximized, the probability of balking
67    between different facilities could be disparate. This could result in serious service
68    inequity issue. Therefore, this study tries to propose a method to deploy capacitated
69    health facilities to alleviate accessibility-based service inequity arising from balking.

70    The service equity issue is one of the critical problems concerned by the users,
71    especially for public health services. Tao et al. (2014) and Zhang et al. (2016) proposed
72    to locate health facilities by maximizing equity in accessibility. The disparity in
73    accessibility to health facilities is noticed and optimized. They adopted a general
74    definition of facility accessibility and minimized the variance of accessibility.
75    Mousazadeh et al. (2018) suggested to design an accessible, stable, and equitable health
76    service network where the equity is incorporated by maximizing the minimum service
77    level of each residential zone. It is the well-known John Rawls's social justice approach
78    where the welfare of the worst group is maximized. Filippi et al. (2021) found that the
79    equitable treatment of users was usually neglected. They suggested a way to
80    compromise between efficiency and equity. Pourrezaie-Khaligh et al. (2022) proposed
81    a bi-objective approach for health facility location problem considering both equity and
82    accessibility. The objective is to minimize system costs, maximize accessibility, and
83    minimize inequality among all demand nodes. They employed the accessibility index
84    introduced by Wang and Tang (2013). Different from conventional way of equity
85    measured using the variance of individual accessibility, they defined equity based on a
86    minimum envy criterion. All in all, although accessibility-based service equity has
87    started few attention recently, the congestion effect and user choice behavior have not
88    been incorporated yet.

89    Service facility location problems have been widely studied because of numerous
90    real-life applications. Most literature is concerned with various versions of the problem
91    where users are simply assigned to closest facilities, while sidesteps the important issue
92    of user choice behavior, as well as the effect of congestion. In fact, the users have
93    freedom to choose facilities. In addition to the travel time, the waiting time of a user at
94    a congested facility also has a significant influence on her/his choice (Marianov et al.,
95    2005; Marianov et al., 2008). It is a congestion game problem. From the perspective of

96  user choice behavior, previous studies could fall into two categories: (i) system optimal
97  models, where users are directed by a central decision-maker to optimize system
98  performances; and (ii) user choice models, where users are free to choose a facility. The
99  congestion at a facility is beginning to be introduced both. Let us give a brief overview
100 on them separately.

101    The system optimal models accounts for the major part of facility location problems
102 where the users are assigned to the closest facilities. They are also known as all-or-
103 nothing allocation, or winner-takes-all allocation. Verter and Lapierre (2002)
104 investigated the problem of locating preventive health facilities using a system optimal
105 model. The travel time was assumed to be the only determinant of facility choice and
106 users would go to the closest facility without considering the congestion effect.
107 Although the users from the same residential node can be directed to different facilities
108 in theory, the optimization problems will have an optimal solution where all-or-nothing
109 allocation is adopted (Castillo et al., 2009). Zhang et al. (2009) further incorporated
110 congestion effect at a facility where the users are assumed to visit the facility with
111 minimum total cost including travel time and queueing time. The queueing time can
112 also be incorporated as a constraint (Davari et al., 2016). Multi-objective location
113 problems are also proposed recently where multiple performances are evaluated (Dogan
114 et al., 2020; Erdoğan et al., 2019). As the outbreak of COVID-19, Risanger et al. (2021)
115 recently proposed a system optimal model to select pharmacies for COVID-19 testing
116 to ensure accessibility.

117    The user choice models are emerging ways of facility location problems. Most
118 location models assume that all the demand originating at a particular node is served by
119 the same closest facility. This is not so in competitive situations where the users are free
120 to choose a facility. In this case, the users at each demand node may choose different
121 facilities to patronize. The more attractive the facility for users at a certain demand node,
122 the larger the percentage it captures the demand originating there. The formulation of
123 user choice behavior is the foundation of facility network design. However, the user
124 choice behaviors in facility location problems, are usually sidestepped intentionally or
125 unintentionally. Although the literature concerning facility location is vast, few studies
126 have incorporated user choice behaviors (Dan and Marcotte, 2019). Generally speaking,
127 the user choice models can be classified into two categories: one is proportional
128 allocation and the other is equilibrium allocation. The proportional allocation can be
129 further classified into Huff-based allocation which can revert to a gravity model with

130    pre-specified parameters (Gu et al., 2010; Tao et al., 2016) and logit-based allocation

131    where a multinomial logit function is used to model the probability that users choose

132    facility (Abouee-Mehrizi et al., 2011; Filippi et al., 2021; Kucukyazici et al., 2020).

133    However, the proportional allocation cannot account for congestion effect. It is well-

134    known that as a facility captures more users, it becomes more congested, resulting in

135    longer queueing times. In fact, this effect makes the service facility less attractive and,

136    consequently, user capture is reduced, leading to an eventual user equilibrium state

137    where no user can further reduce his cost by unilaterally changing his behavior.

138    Therefore, the equilibrium allocation was suggested recently which includes

139    deterministic user equilibrium when utility is deterministic and stochastic user

140    aquarium when stochastic utility is assumed (Dan and Marcotte, 2019; Zhang and

141    Atkins, 2019). However, it is few incorporated due to the computation complexity. The

142    facility location problem with equilibrium allocation is still a cutting-edge problem

143    deserved to be explored.

144        This study makes four main theoretical and practical contributions for urban planning

145    community. (i) We propose a way to improve the accessibility-based service equity for

146    capacitated preventive health infrastructure planning. The equitable accessible flow is

147    achieved by minimizing the maximum probability of balking. (ii) A bilevel decision

148    structure is adopted where the upper level is urban planners and the lower level is

149    facility users. The congestion effect is incorporated in the user utility function including

150    queueing time and probability of balking. (iii) The users competing with each other will

151    lead to user equilibrium state. An equivalent mathematical programming model is

152    proposed to predict facility demand volumes at equilibrium state. (iv) A generic

153    efficient and effective solution algorithm is proposed and validated, which is also

154    applicable for other public service facility planning problems. (v) Several interesting

155    findings and managerial insights for urban planners are provided based on

156    computational experiments.

157        The remainder of this paper is organized as follows. Section 2 describes the problem

158    and formulates it as a bi-objective bilevel programming model. Section 3 proposes a

159    heuristic algorithm to solve the bilevel problem. Section 4 presents the computational

160    results for the model with managerial insights. Finally, conclusions and future research

161    directions are provide in section 5.

162    **2. Problem modeling**

163        Let $H = (\mathbf{N}, \mathbf{L})$ be a road network with a set of nodes $\mathbf{N}$ and a set of links $\mathbf{L}$.

The nodes represent either demand concentrations, facility locations, or road intersections, and links are main transportation arteries between nodes. We assume that the demand rate requiring preventive health service at population node $i(i \in \mathbf{N})$ follows the Poisson process with an average rate $h_i$. The set of candidate locations for health facilities is $\mathbf{M}$, and $\mathbf{S}$ is the set of chosen locations where $\mathbf{S} \subset \mathbf{M}$. The shortest path travel time from demand node $i(i \in \mathbf{N})$ to facility location $j(j \in \mathbf{M})$ is denoted by $t_{ij}$. The government has a limited budget $B$ that can be used to build facilities with associated servers. We assume that servers at all of the facilities are the same, and service time is exponentially distributed providing service to $\mu$ clients per unit of time on average. We also assume that clients are homogenous, their arrivals to each facility follow poisson distributions and the queueing discipline is first-come first served (FCFS). These assumptions are reasonable for walk-in facilities, which applies to most routine health services in many countries or regions. Thus, facility $j(j \in \mathbf{M})$ here is assumed to behave as a $M/M/s_j/K_j$ queueing system, where $M$ denotes Markovian (or poisson) arrivals or departures distribution, or equivalently exponential interarrival or service time distribution, $s_j$ denotes the number of servers at facility $j$, and $K_j$ is the capable number of clients at facility $j$ due to physical constraint. Whenever there are $K_j$ clients at facility $j$, any arriving client is denied to access and leaves the system as a lost client. The value of $K_j$ is predetermined for each facility location, depending on specific conditions. This assumption is not loss of generality since it could be extended to other queueing system based on estimation from available data.

The problem is to make location and associated capacity decisions, with the aim of equitable probability of balking, subject to the budget constraint $B$. Three sets of decision variables are defined as follows:

$$y_j = \begin{cases} 1 & \text{if a facility is opened at location } j, \forall j \in \mathbf{M}, \\ 0 & \text{otherwise,} \end{cases}$$

$s_j$ = number of servers at facility location $j$, $\forall j \in \mathbf{M}$,

$x_{ij}$ = number of clients from demand node $i$ to location $j$, $\forall i \in \mathbf{N}$, $j \in \mathbf{M}$.

Therefore, for a chosen set $\mathbf{S} = \{j : j \in \mathbf{M}, y_j = 1\}$, we have

$$\sum_{j \in \mathbf{S}} x_{ij} = h_i, \quad \forall i \in \mathbf{N}. \tag{1}$$

Let $\lambda_j$ denotes the arrival rate of clients at facility $j$, $\forall j \in \mathbf{M}$, then we have

$$195 \qquad \lambda_j = \sum_{i \in \mathbf{N}} x_{ij}, \quad \forall j \in \mathbf{M}. \qquad (2)$$

196 It defines the demand at each facility as the sum of demands originating from all the

197 demand nodes. Given the arrival rate $\lambda_j$ and the number of servers $s_j$ at facility $j$,

198 the probability that there are $n$ clients in the queue is

$$199 \qquad p_{nj}(\lambda_j, s_j) = \begin{cases} \dfrac{\rho_j^n}{n!} p_{0j}, & \text{if } 0 \le n \le s_j, \\[2ex] \dfrac{\rho_j^n}{s_j! s_j^{n-s_j}} p_{0j}, & \text{if } s_j \le n \le K_j, \end{cases} \qquad (3)$$

200 where $\rho_j = \lambda_j / \mu$ is the intensity of the queueing process and the probability of no

201 client is

$$202 \qquad p_{0j} = \left[1 + \sum_{n=1}^{s_j} \frac{\rho_j^n}{n!} + \frac{\rho_j^{s_j}}{s_j!} \sum_{n=s_j+1}^{K_j} \left(\frac{\rho_j}{s_j}\right)^{n-s_j}\right]^{-1}. \qquad (4)$$

203 Note that the probability $p_{nj}$ at each facility $j$ is a function of $\lambda_j$ and $s_j$. The

204 notation $p_{Kj}$ is the probability of balking owing to a limited space. It allows a

205 facility's arrival rate to exceed its service rate, without unbounded grow of queue length.

206 The effective arrival rate, i.e., the number of clients who could access the service, is

207 denoted by $\bar{\lambda}_j$. There is,

$$208 \qquad \bar{\lambda}_j = \lambda_j (1 - p_{Kj}), \quad \forall j \in \mathbf{M}. \qquad (5)$$

209 **2.1 The user utility function**

210     Clients are assumed to patronize a facility that maximizes their individual utility, i.e.,

211 minimizes their generalized costs. Therefore, it is critical to understand how clients

212 make their choices. Let us now present our user choice modeling, which essentially

213 establishes a utility function depending on the attractiveness of a facility that they are

214 aware of. Let $U_{ij}$ denote the observed utility of users from demand node $i$ receiving

215 the service at facility location $j$. It mainly comprises four components: (i) $u_j$, a

216 constant attraction of location $j$, which might include intrinsic factors such as parking

217 convenience, practitioner reputation, service quality, etc.; (ii) $t_{ij}$, the shortest path travel

218 time from origin node $i$ to destination facility $j$; (iii) $w_j(\lambda_j, s_j)$, the average dwell

219 time at location $j$ including queueing time and service time, which is a function of

220 arrival rate $\lambda_j$ and server number $s_j$; and (iv) $p_{Kj}(\lambda_j, s_j)$, the probability of unmet

221 service (i.e., balking) due to physical constraint. Note that $w_j$ and $p_{Kj}$ are

222 continuous functions with respect to $\lambda_j$ and $s_j$.

223     As it is an $M/M/s_j/K_j$ queueing system at facility $j$, for any $s_j \geq 1$, the

224 average dwell time $w_j(\lambda_j, s_j)$ could be given by the following set of equations

225 according to the classical queueing theory:

$$w_j(\lambda_j, s_j) = \frac{L_j}{\overline{\lambda}_j}, \quad \forall j \in \mathbf{S}, \tag{6}$$

226

$$L_j = \sum_{n=s_j}^{K_j} (n - s_j) p_{nj} + \rho_j(1 - p_{Kj}), \quad \forall j \in \mathbf{S}, \tag{7}$$

227

228 where $L_j$ is the average length of the queue in terms of client number, $\overline{\lambda}_j$ is the

229 effective arrival rate according to Eq. (5), $p_{nj}$ is the probability of having $n$ clients

230 at the facility according to Eq. (3), and $\rho_j$ is the intensity of service as defined

231 previously. Eq. (6) is the famous Little's formula in queueing theory.

232     The way of integrating utility could be various. Following the conventional way in

233 the literature, we assume a linear additive functional form of $U_{ij}$ to incorporate the

234 above four components with different weights. It is a standard assumption in the utility

235 theory. In addition, it is also reasonable to assume that $U_{ij}$ is positively associated

236 with benefit $u_j$ but negatively associated with cost $t_{ij}$, $w_j(\lambda_j, s_j)$, and $p_{Kj}(\lambda_j, s_j)$.

237 In this framework, $U_{ij}$ is given by (Dan and Marcotte, 2019):

$$U_{ij} = u_j - \beta_1 t_{ij} - \beta_2 w_j(\lambda_j, s_j) - \beta_3 p_{Kj}(\lambda_j, s_j), \quad \forall i \in \mathbf{N}, \ j \in \mathbf{S}, \tag{8}$$

238

239 where $\beta_1$ and $\beta_2$ denote the coefficients of the travel time and queueing time

240 respectively, and $\beta_3$ is interpreted as the price of service inaccessibility. In practice,

241 parameters $\beta_1$, $\beta_2$, and $\beta_3$ can be estimated empirically using realistic surveys. The

242 different weights on travel time and waiting time could be possible and are allowed,

243 given the different perceptions of clients for them. The definition of real values for these

244 parameters is outside the scope of this paper. Note that besides these specific parts, the

245 utility function can also be extended to incorporate other observable attributes, such as

246 the parking cost and service price, depending on available data.

247     The users interacts with each other until no one person could increase his utility by

248 unilaterally changing his facility choice, which is known as Nash equilibrium state.

249 Mathematically it is important to note the interdependency between the arrival rate $\lambda_j$

250 and the expected waiting time $w_j(\lambda_j, s_j)$ and the probability of balking $p_{Kj}(\lambda_j, s_j)$.

251 According to our modelling framework, $\lambda_j$ is the sum of $x_{ij}$, which depends on $U_{ij}$,

252 which further depends on $w_j(\lambda_j, s_j)$ and $p_{Kj}(\lambda_j, s_j)$. That is, the value of $\lambda_j$

253 depends on itself indirectly. Since we consider a network of competitive facilities, it
254 implies that we need address a Nash equilibrium problem to determine demand
255 allocation $x_{ij}$ given facility locations and associated capacities. Specifically, it is
256 better known as user equilibrium problem.

## 2.2 The user equilibrium model

258 It is assumed that clients always choose the facility with the highest observed utility.
259 The clients are assumed to re-evaluate their utilities after several times of visits. They
260 could also learn from others by social network for example. Therefore, they are
261 assumed to know about the queues and capacities of the facilities to make near-optimal
262 decisions. The competition between clients will reach a user equilibrium state finally.
263 Let $\bar{U}_i$ denote the highest utility of clients at demand node $i$, i.e.,

$$\bar{U}_i = \max_{j \in M} U_{ij}, \quad \forall i \in \mathbf{N}. \tag{9}$$

265 Given the determined location $\mathbf{S}$ and capacities $s_j$, $\forall j \in \mathbf{S}$, no client wants to
266 change her/his facility choice at user equilibria. Therefore, the equilibrium condition
267 can be characterized by the following complementarity system

$$U_{ij}^* = u_j - \beta_1 t_{ij} - \beta_2 w_j(\lambda_j^*, s_j) - \beta_3 p_{Kj}(\lambda_j^*, s_j) \begin{cases} = \bar{U}_i^* & \text{if } x_{ij}^* > 0 \\ \leq \bar{U}_i^* & \text{if } x_{ij}^* = 0 \end{cases}, \forall i \in \mathbf{N}, \ j \in \mathbf{S}, \tag{10}$$

269 where $U_{ij}^*$ and $\bar{U}_i^*$ denote the utility of clients from demand node $i$ visiting
270 preventive health facility $j$ and the highest utility of clients from demand node $i$ at
271 user equilibrium state, respectively. Moreover, it should be noted that

$$\lambda_j^* = \sum_{i \in \mathbf{N}} x_{ij}^*, \quad \forall j \in \mathbf{S},$$

273 where $\lambda_j^*$ denotes the arrival rate of clients at facility $j$ at user equilibrium state, and
274 $x_{ij}^*$ denotes the allocated number of clients from demand node $i$ to facility location
275 $j$ at user equilibrium state.
276 The equilibrium condition (10) means that if there is a client flow from demand node
277 $i$ to facility location $j$, then $U_{ij}^*$, the utility of users from node $i$ to facility $j$, must
278 be equal to the highest utility $\bar{U}_i^*$; otherwise, it is no more than the highest. It implies
279 that each user patronizes the facility with the highest observed utility. Accordingly, at
280 equilibrium state, users issued from a common origin node will experience identical
281 utilities, thus achieving a well-known Nash equilibrium state. They cannot improve

282     their utility by changing facility choice.

283      To find $\lambda_j^*$ and implicit $x_{ij}^*$ in Eq. (10) given determined location $\mathbf{S}$, we can solve

284 the following equivalent nonlinear mathematical programming with symmetric Jacobin

285 matrix of utility function:

$$\max_{\mathbf{x}} \; Z(\mathbf{x}\,|\,\mathbf{S}) = \sum_{i\in\mathbf{N}}\sum_{j\in\mathbf{S}} \int_0^{\lambda_j} U_{ij}(\omega, s_j)\,d\omega \tag{11}$$

287     subject to

$$\sum_{j\in\mathbf{S}} x_{ij} = h_i, \quad \forall j\in\mathbf{S}, \tag{12}$$

$$x_{ij} \geq 0, \quad \forall i\in\mathbf{N},\, j\in\mathbf{S}, \tag{13}$$

290     where

$$\lambda_j = \sum_{i\in\mathbf{N}} x_{ij}, \quad \forall i\in\mathbf{N},\, j\in\mathbf{S}. \tag{14}$$

292     **Theorem 1.** *Given the determined location $\mathbf{S}$, the mathematical programming (11)*

293 *-(14) is equivalent to equilibrium condition (10).*

294     **Proof.** In order to prove that the mathematical programming is equivalent to Eq. (10),

295 we reformulate the model as a Lagrange function with nonnegative constraints only,

296 i.e.,

$$F = Z(\mathbf{x}\,|\,\mathbf{S}) - \sum_{i\in\mathbf{N}} w_i \left(\sum_{j\in\mathbf{S}} x_{ij} - h_i\right)$$
$$s.t. \quad x_{ij} \geq 0, \quad \forall i\in\mathbf{N}, \quad j\in\mathbf{S}, \tag{15}$$

298     where $w_i$ is a Lagrange multiplier of constraint (12).

299      According to Karush–Kuhn–Tucker (KKT) conditions, the optimal conditions of

300 this Lagrange function are given by

$$x_{ij}\frac{\partial F}{\partial x_{ij}} = 0, \quad \forall i\in\mathbf{N},\, j\in\mathbf{S}, \tag{16}$$

$$\frac{\partial F}{\partial x_{ij}} \leq 0, \quad \forall i\in\mathbf{N},\, j\in\mathbf{S}, \tag{17}$$

$$\frac{\partial F}{\partial w_i} = 0, \quad \forall i\in\mathbf{N}, \tag{18}$$

$$x_{ij} \geq 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}. \tag{19}$$

It is straightforward to find that Eq. (18) is equivalent to Eq. (12). Eqs. (16) and (17) imply that

$$
\begin{aligned}
&\text{if } x_{ij} > 0, \ \frac{\partial F}{\partial x_{ij}} = 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}, \\
&\text{if } x_{ij} = 0, \ \frac{\partial F}{\partial x_{ij}} \leq 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}.
\end{aligned}
\tag{20}
$$

Note that since we have

$$
\begin{aligned}
\frac{\partial L}{\partial x_{ij}} &= \frac{\partial}{\partial \lambda_j}[\sum_{i \in \mathbf{N}}\sum_{j \in \mathbf{S}}\int_0^{\lambda_j} U_{ij}(\omega, s_j)d\omega]\frac{\partial \lambda_j}{\partial x_{ij}} - \frac{\partial}{\partial x_{ij}}(\sum_{i \in \mathbf{N}} w_i(\sum_{j \in \mathbf{S}} x_{ij} - h_i)) \\
&= U_{ij} - w_i,
\end{aligned}
\tag{21}
$$

Eq. (20) can be further rewritten as follows:

$$
\begin{aligned}
&\text{if } x_{ij} > 0, \ U_{ij} - w_i = 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}, \\
&\text{if } x_{ij} = 0, \ U_{ij} - w_i \leq 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}.
\end{aligned}
\tag{22}
$$

It can be also reformulated in the following complementary form:

$$(U_{ij} - w_i)x_{ij} = 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}, \tag{23}$$

$$U_{ij} - w_i \leq 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}, \tag{24}$$

$$x_{ij} \geq 0, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}. \tag{25}$$

It can be seen that Eq. (22) means that if there is client flow, i.e., $x_{ij} > 0$, the utility $U_{ij}$ will be equal to $w_i$, and if there is no client flow, i.e., $x_{ij} = 0$, the utility $U_{ij}$ is no more than $w_i$. Therefore, the Lagrange multiplier $w_i$ can be interpreted as the highest utility $\bar{U}_i^*$ incurred by clients at demand node $i$. Hence, Eq. (22) is equivalent to Eq. (10). Therefore, we can conclude that the solution of the mathematical programming (11)-(14) satisfies the equilibrium condition (10). The proof of the theorem is complete.

## 2.3 The bilevel programming model

The entire problem considered here is a bilevel decision structure where the upper level problem is the determination of facility locations and associated capacities by urban planners, and the lower level problem is the determination of equilibrium flows of users from demand nodes to facility locations given the upper level decisions. Note that the equilibrium flows $x_{ij}$, as well as the arrival rate $\lambda_j$ are not decision variables. They are determined endogenously by the lower level model. The decision variables are the location variables $y_j$ and associated capacities $s_j$ in the upper level model. Once the values of these variables are fixed, all of the remaining auxiliary variables and parameters can be computed.

There usually is a limited investment budget to support the establishment and operation of the preventive health facilities in practice. This budget constraint can be used to incorporate the cost differences of establishing and operating facilities at different locations of an urban area. The budget is set to be $B$. Let $c_j^f$ be the fixed cost of establishing a facility at location $j \in \mathbf{M}$ and $c^v$ be the unit operation cost of adding a server to a facility that is identical for each location. In addition, for cost effectiveness, we assume that facilities cannot be operated unless the number of their clients exceeds a minimum threshold $R_{\min}$. Moreover, the number of servers at facility $j$ cannot exceed an upper bound $\hat{s}_j$ due to physical condition. The value of $\hat{s}_j$ is typically given by the urban planner on the basis of specific conditions and may differ from location to location.

Bi-objective optimization is adopted in the upper level model where the upper level is urban planners and the lower level is facility users. In order to formulate service network design problem considering accessibility-based equity, the maximum probability of balking $p_{Kj}$, $\forall j \in \mathbf{M}$, is minimized. It is regarded as the primary objective. As there are possible chances that the primary objective is always zero for unsaturated flows, a secondary objective is introduced to minimize the maximum waiting time at a facility in order to reach equitable queueing. Therefore, the upper level model of service network design problem can be formulated as follows:

$$\text{Primary Objective} \quad \min E_1(\mathbf{S}) = \max\left\{p_{Kj}, \forall j \in \mathbf{M}\right\} \tag{26}$$

$$\text{Secondary Objective} \quad \min E_2(\mathbf{S}) = \max\left\{w_j(\lambda_j, s_j), \forall j \in \mathbf{M}\right\} \tag{27}$$

subject to

$$s_j \geq y_j, \quad \forall j \in \mathbf{M}, \tag{28}$$

$$s_j \leq \hat{s}_j y_j, \quad j \in \mathbf{M}, \tag{29}$$

$$\sum_{i \in \mathbf{N}} x_{ij} = \lambda_j, \quad \forall j \in \mathbf{M}, \tag{30}$$

$$x_{ij} \leq y_j, \quad \forall i \in \mathbf{N}, \quad j \in \mathbf{M}, \tag{31}$$

$$\bar{\lambda}_j = \lambda_j (1 - p_{Kj}), \quad \forall j \in \mathbf{M}, \tag{32}$$

$$\lambda_j \geq R_{\min} y_j, \quad \forall j \in \mathbf{M}, \tag{33}$$

$$\sum_{j \in \mathbf{M}} c_j^f y_j + c^v \sum_{j \in \mathbf{M}} s_j \leq B, \tag{34}$$

$$y_j \in \{0,1\}, \quad s_j \in \mathbf{Z}^+, \quad \forall j \in \mathbf{M}. \tag{35}$$

where $x_{ij}$ is determined by the following lower level model after the location variables $y_j$ and associated capacity variables $s_j$ are determined:

$$\max_{\mathbf{x}} \; Z(\mathbf{x} \mid \mathbf{S}) = \sum_{i \in \mathbf{N}} \sum_{j \in \mathbf{S}} \int_0^{\lambda_j} U_{ij}(\omega, s_j) d\omega \tag{36}$$

subject to

$$\sum_{j \in M} x_{ij} = h_i, \quad \forall j \in \mathbf{S} \tag{37}$$

$$x_{ij} \geq 0, \quad \forall i \in \mathbf{N}, \quad j \in \mathbf{S}. \tag{38}$$

The primary objective function (26) is to minimize the maximum probability of balking and the secondary objective function (27) is to minimize the maximum queueing time. They are both Min-Max optimization problems so as to reach service equity, which is robust for any level of demands. Constraints (28) ensure the assignment of at least one server to each open facility. Constraints (29) limit the number of servers not exceeding $\hat{s}_j$. Constraints (30) define the arrival rate $\lambda_j$. Constraints (31) ensures that clients can only obtain the service from open facilities. Constraints (32) are the definition of effective arrival rates. Constraints (33) stipulate that the arrival rate at an open facility must satisfy the minimum workload requirement. Constraint (34) is the budget and Constraints (35) define the feasible domain of decision variables $y_j$ and $s_j$.

**3. Solution method**

Since the bi-objective bilevel programming model is highly nonlinear and contains

382　integer decision variables, it poses big challenges to solve the model exactly. Therefore,

383　the focus of this study is to adopt efficient and effective heuristic algorithms that have

384　many successful applications for facility network design problem (Auerbach and Kim,

385　2021; Chambari, et al, 2011; Ershadi and Shemirani, 2021; Zhang and Atkins, 2019).

386　The bilevel framework is carefully followed by our solution method. For the upper level

387　location problem, a meta-heuristic, generic algorithm (GA) with elite strategy, is

388　proposed to find the optimal locations and associated capacities. However, we

389　definitely believe that the more advanced heuristics will improve the computation

390　efficiency. For the lower level allocation problem, we adopt method of successive

391　averages (MSA) to solve the user equilibrium model. This demand allocation algorithm

392　determines the equilibrium flows of users to facilities after the upper level decisions are

393　confirmed. Thus, the demand allocation algorithm serves as an embedded module for

394　the facility location algorithm. For the ease of easier understanding, we describe the

395　demand allocation algorithm first.

396　**3.1 Demand allocation algorithm for the lower level model**

397　　Given the upper level facility decisions $\mathbf{S}$ and $s_j$, $\forall j \in \mathbf{S}$, the lower level

398　problem of the user choice model is to find the equilibrium flows. The adopted

399　algorithm is a kind of iterative method, known as MSA. Let $k$ be the iteration index

400　and $K$ be a maximum iteration number. In addition, let $\varepsilon$ be a predetermined error

401　tolerance parameter, and $\theta_k \in (0,1)$, $k = 1,\ldots,K$, be a step-length parameter at

402　iteration $k$. The specific computation steps are listed below:

403　　***Step 0 (Initialization)***: Set the values of $\varepsilon$ and $K$; initiate $k = 0$; set initial

404　allocation

405
$$x_{ij}^0 = \frac{h_i}{|\mathbf{S}|}, \forall i \in \mathbf{N}, j \in \mathbf{S} .$$

406　　***Step 1 (Calculation of utility)***: Update $k := k+1$; calculate $\lambda_j$, $\forall j \in \mathbf{S}$, from Eq.

407　(2); calculate the shortest path travel time $t_{ij}$, $\forall i \in \mathbf{N}$, $j \in \mathbf{S}$, using Dijkstra's

408　algorithm; calculate probability of balking $p_{Kj}(\lambda_j, s_j)$ from Eq. (3), effective arrival

409　rate $\bar{\lambda}_j$ from Eq. (5), waiting time $w_j(\lambda_j, s_j)$ from Eq. (6); calculate $U_{ij}$, $\forall i \in \mathbf{N}$,

410　$j \in \mathbf{S}$, from Eq. (8); find $\bar{U}_i$, $\forall i \in \mathbf{N}$, from Eq. (9).

411　　***Step 2 (All-or-nothing allocation)***: Set flow $x_{ij}'$ by all-or-nothing rule as follows,

412　i.e., allocate all clients from the same demand node to the most attractive facility, i.e.,

413
$$x'_{ij} = \begin{cases} h_i & \text{if } U_{ij} = \bar{U}_i \\ 0 & \text{if } U_{ij} < \bar{U}_i \end{cases}, \quad \forall i \in \mathbf{N}, j \in \mathbf{S}.$$

414 **Step 3 (Generation of search direction)**: Define $d_{ij} = x'_{ij} - x_{ij}^{k-1}$, $\forall i \in \mathbf{N}$, $j \in \mathbf{S}$, as

415 a search direction.

416 **Step 4 (Flow update)**: Update client flow $x_{ij}^k = x_{ij}^{k-1} + \theta_k d_{ij}$, $\forall i \in \mathbf{N}$, $j \in \mathbf{S}$, where

417 $\theta_k$ is the step-length parameter given by,

418
$$\theta_k = \frac{1}{k+1}.$$

419 **Step 5 (Stopping criteria)**: If the relative difference between $x_{ij}^k$ and $x_{ij}^{k-1}$ is equal

420 or less than $\varepsilon$, or $k \geq K$, set $x_{ij} := x_{ij}^k$ and stop; otherwise, go to Step 1. The relative

421 error is defined as,

422
$$\frac{\| x_{ij}^k - x_{ij}^{k-1} \|}{\| x_{ij}^{k-1} \|} \leq \varepsilon, \forall i \in \mathbf{N}, j \in \mathbf{S}.$$

423 **Step 6 (Return results)**: Return the incumbent solution to the upper level model,

424 including equilibrium flows, balking probabilities, and waiting times.

425 The suggested method in each iteration identifies a new search direction for $x_{ij}$ in

426 Step 3 and then updates $x_{ij}$ by a step-length in Step 4. The procedure continues until

427 one of the stopping conditions in Step 5 is met. The step-length $\theta_k$ in each iteration is

428 determined in advance. There are a variety of ways to set $\theta_k$. To achieve convergence,

429 $\theta_k$ should decrease with $k$ and locate between zero and one. Here we set $\theta_k$ as the

430 reciprocal of the iteration number $(k+1)$ as usual. It is worth noting that the $x_{ij}^k$

431 updated in Step 4 may result in an arrival rate at a facility exceeds the capable space

432 allowed. At this situation, excess clients will be denied to access health services and

433 become lost demand.

434 **3.2 Facility location algorithm for the upper level model**

435 We develop a genetic algorithm with elite strategy to solve the upper level problem,

436 because it is one of the most popular meta-heuristics for addressing combinatorial

437 optimization problems with many successful applications. It has the ability to explore

438 other parts of the feasible space while avoiding local optima. Although it is time-

439 consuming and the global optima is not guaranteed mathematically, it is still widely

440 used for nonlinear programming problems.

441 In genetic algorithms, each chromosome represents a solution to the problem, and

442 the quality of a solution is measured by a fitness value. Note that since it is a bi-objective

443 optimization problem, the primary objective works as the first fitness value, and the
444 secondary objective is adopted when there are equal primary objectives. In this study,
445 an integer coding technique is employed to define a chromosome. Each chromosome is
446 made up of several genes that are nonnegative integer numbers. Each gene corresponds
447 to a candidate location in $M$ , and its value represents the number of servers. If there
448 is no server at a location, the facility is not opened at that location. The following is
449 how we implement the genetic algorithm with elite strategy:

450 **Step 0 (Initialization)**: Set the used parameters, including the population size $P$ , the
451 maximum number of generations $G$ , the crossover probability $p_c$ , the mutation
452 probability $p_m$ , the label of generation $g = 1$, and the fraction of elite $p_e$ .

453 **Step 1 (Generation of initial population):** Randomly generate $P$ feasible solutions
454 as an initial population of chromosomes, scattering the entire range of possible solutions.
455 If one chromosome is not feasible according to the constraints, generate another one
456 until a feasible solution is found.

457 **Step 2 (Calculation of fitness value)**: For each chromosome in the population, the
458 value of fitness is generated that is the objective function value. It is used to evaluate
459 the quality of each chromosome in the population. Note that there are two objective
460 functions in the upper level model. One is primary objective and the other is secondary
461 objective. Therefore, there are two fitness values in order.

462 **Step 3 (Generation of new population)**:

463 **Step 3.1 (Selection)**: According to the values of fitness evaluated in Step 2, the best
464 fraction $p_e$ is labeled for elites, and the worst fraction $p_e$ is discarded. A stratified
465 sequencing method is used here where the primary objective value is sorted first and
466 the secondary objective value is sorted next.

467 **Step 3.2 (Crossover)**: The remaining $(1 - p_e)P$ chromosomes are used for
468 crossover operation. These chromosomes are matched in pairs randomly. The
469 probability of carrying out the crossover is $p_c$ . If the two parent chromosomes are
470 chosen for crossover, a gene location is randomly identified to across over to generate
471 two off-springs as new chromosomes. If newborn chromosomes are not feasible
472 according to constraints in the upper level model, try another gene location until they
473 are feasible.

474 **Step 3.3 (Mutation)**: A chromosome is determined for mutation with probability
475 $p_m$ . Randomly choose two genes with at least one positive, and interchange their
476 values. If the new chromosome is not feasible, try another two gene locations until a

477 feasible off-spring is generated.

478       ***Step 3.4 (Elitism)***: Form a new generation. After genetic operations, there are still

479 $(1 - p_e)P$ feasible chromosomes. The labeled $p_e P$ elites are added to ensure the

480 population size $P$. This allows the best chromosomes from the current generation to

481 carry over the next generation unaltered. It guarantees that the solution quality will not

482 decrease from one generation to the next. Update the notation of generation be

483 $g := g + 1$.

484     ***Step 4 (Stopping criterion)***: If the maximum number of generations $G$ is achieved,

485 i.e., $g \geq G$, terminate the iteration process and output the results. Otherwise, turn to

486 Step 2.

487 **4. Computational experiments**

488 **4.1 An illustrative case**

489     We conduct a computational experiment to assess the performance of proposed

490 model and algorithm with Sioux Falls network. This network has been widely used for

491 validation in the network design problems. It is a medium sized network as depicted in

492 Fig. 1. The network consists of 24 nodes and 76 links. In the computational experiments,

493 it is assumed that there are 8 population nodes and 8 potential locations in the region.

494 Therefore, there are a total number of 64 origin-destination (O-D) pairs. The travel time

495 and length of each link are given in Table 1. The link length can be converted to the

496 link travel time, by assuming a constant link travel speed of 30 miles/hour. Recognizing

497 that clients are only a very small part of road travelers, the facility choice and route

498 choice of clients is assumed not to affect road travel times. That is, the link travel times

499 are constants. This is different with classical road network design problems as the

500 congestions take place in facilities other than roads. The preventive health demand data,

501 i.e., the number of clients per hour (clients/hr), are listed in Table 2. The demand for

502 preventive health services is fixed at origin zones while the trip distribution is not fixed

503 and it is determined by facility planning. The clients have freedom to choose their

504 favorite facilities.

505

506                    **Fig. 1.** The Sioux Falls test network

507

508                     ***Insert Table 1 here***

509

510                     ***Insert Table 2 here***

511

Based on the proposed model and solution method, the following parameter values are used in the case study.

*Problem parameters*

· the service rate of each server $\mu = 6$ clients/hr;

· the constant facility attraction $u_j = 0$;

· the coefficient to travel time $\beta_1 = 1$ and that to waiting time $\beta_2 = 1$;

· the price of service inaccessibility $\beta_3 = 1$;

· the maximum number of servers $\hat{s}_j = 10$;

· the fixed establishment cost $c_j^f = 0$;

· the unit cost of a server $c^v = 1$;

· the budget $B = 40$;

· the minimum workload $R_{\min} = 10$ clients/hr;

*Method of successive averages parameters*

· the maximum iteration number $K = 100$;

· the error tolerance $\varepsilon = 0.01$;

*Genetic algorithm parameters*

· the population size $P = 100$;

· the maximum number of generations $G = 20$;

· the crossover probability $p_c = 0.5$;

· the mutation probability $p_m = 0.2$;

· the fraction of elite $p_e = 0.1$.

The algorithms are coded using a free open-source language R 3.6.3. All runs are performed at a personal computer with 3.6 gigahertz Intel i7-4790 CPU and 16 gigabytes RAM. The genetic algorithm stopped after 1.42 hours for this case study. The evolutionary process begins to be stable after 19 generations as shown in Fig. 2. It can be concluded that the final results are satisfying solutions. The selected locations to set up preventive health facilities are nodes 3, 9, 16, 19, and 23. Their associated number of servers are 6, 9, 8, 7, and 10 correspondingly. The service quality of each facility is shown in Table 3. It shows that the maximum probability of balking is 0.132 in node 3 and the maximum waiting time is 1.22 hours in node 23. It can be concluded that the service quality among all facilities is quasi-equal in terms of balking probability and waiting time. The accessibility-based service equity is achieved which is the policy goal.

544

545                  **Fig. 2.** The evolutionary process of genetic algorithm

546

547                  ***Insert Table 3 here***

548      The demand allocation at equilibrium state is presented in Table 4. It shows that

549 clients from the same demand node usually patronize the same facility such as nodes 1,

550 2, 4, 5, 13, and 14, even if they are free to head for different facilities. However, the

551 clients are possible to be assigned to more than one facility such as nodes 19 and 23 if

552 their utilities are quasi-equal.

553

554                  ***Insert Table 4 here***

555

556 **4.2 Sensitivity analysis**

557      It is always beneficial to do a sensitivity analysis which could provide valuable

558 managerial insights. We conduct a sensitivity analysis with varying budget control here,

559 which is also a cost-benefit analysis in economics. The budget is increased from 30 to

560 60 at step-length 5. The results are shown in Fig. 3 where the horizontal axis is budget

561 and the vertical axis is maximum probability of balking among all of facility locations.

562 At the very beginning, the maximum probability of balking is 37.4% with budget 30. It

563 is a low level of service that is difficult to accept. It is no doubt that the probability of

564 balking decreases with budget. The maximum probability of balking is decreased to

565 2.2% with budget 45. The maximum waiting time is 1.67 hours at this time. Whether

566 the budget is good enough depends on the policy makers. The probability of balking

567 will continue to decrease until zero. After budget 50, the customers will not be denied

568 to access facilities, which are unsaturated flows. The service network is not that

569 congested. There are enough vacancies for clients. Since then, the secondary objective,

570 minimizing the maximum waiting time, will play an important role as the primary

571 objective will not move forward and keep zero. Therefore, the proposed methods are

572 robust for capacitated facility location problems.

573

574                **Fig. 3.** A sensitivity analysis with varying budget

575

576      It is also interesting to do a sensitivity analysis on demand with given infrastructure

577 investment budget. The demand is fixed in the short run, but it can change with time in

578  the long run. In order to investigate the benefit of investment budget, a demand
579  expansion coefficient is adopted varying from 0.7 to 1.3 at step-length 0.1. The budget
580  is set to be 40. The results are shown in Fig. 4 where the horizontal axis is varying
581  demand and the vertical axis is maximum probability of balking. At the very beginning,
582  there is no balking and the clients will not be denied to access health service because
583  demand is insufficient. The maximum waiting time is 0.546 hours when demand
584  expansion coefficient is 0.7. If the demand becomes even less ,the probability of facility
585  idleness will increase, which means there is a waste of investment. It is undoubted that
586  the maximum probability of balking will increase with demand. When demand
587  expansion coefficient is 1.3, the maximum probability of balking will increase to 35.3%.
588  If the probability is unacceptable, more infrastructure investment is needed.
589
590  **Fig. 4.** A sensitivity analysis with varying demand
591

592  **5. Conclusions**

593  Preventive health services can detect serious diseases at early stage and make a lot
594  of savings on health expenditures. It is critical for urban sustainable development as
595  shown by the current COVID-19 pandemic. The authorities realize to improve the level
596  of preventive health services to avoid expensive social cost. Noticed the disparity in the
597  accessibility to health facilities among different zones, this study proposes a bilevel
598  programming model to improve the accessibility-based service equity for health
599  infrastructure planning problems. The facilities are capacitated where a customer
600  observes the queue on arrival and leaves if there are no vacancies. In the upper level
601  model, a bi-objective programming model is adopted to descript the urban planner
602  where the primary objective is to minimize the maximum probability of balking and
603  the secondary objective is to minimize the maximum queueing time subject to an
604  investment budget. In the lower level model, a deterministic user equilibrium model is
605  adopted to descript facility users, which is formulated as an equivalent mathematical
606  programming problem. The user utility is defined to include travel time, queueing time,
607  and the probability of balking for capacitated health facilities. The solution method is
608  designed to correspond to the bilevel decision framework where a genetic algorithm
609  with elite strategy is adopted for the upper level model and the method of successive
610  averages is used for the lower level model. Note that a stratified sequencing method is
611  used in genetic algorithm where the primary objective value is sorted first and the

612     secondary objective value is sorted next.

613     To validate the proposed methods, we conduct several computational experiments

614 and derive some interesting managerial insights. We find that the proposed methods are

615 efficient and effective. These methods can reach a satisfactory solution in a reasonable

616 computation time. It is found that although the clients from the same demand node may

617 visit more than one facility, they usually visit the same facility. This is caused by

618 deterministic user equilibrium. The results also indicate that the service quality is quasi-

619 equal in terms of balking probability. The bi-objective method is robust for any level of

620 budget and demand. The sensitivity analysis with varying budget shows that the

621 maximum probability of balking decreases with budget. However, the marginal benefit

622 is decreasing. There is an optimal budget beyond which further increment of investment

623 will not offset its benefits. On the other hand, the sensitivity analysis with varying

624 demand shows that more investment is desired for expanded demand in order to

625 maintain a certain level of service.

626     This study would be improved in several ways in the near future. First, we would

627 like to find a more realistic case with empirical data to show how our methods can be

628 applied in practice. Second, the user utility function will be extended to include other

629 observable attributes, such as the parking time, the service quality, the service price, etc.

630 The formulation of user choice behavior would be more realistic. Third, there are some

631 other advanced heuristics used for similar problems, such as vibration damping

632 optimization algorithm and cutting plane algorithm. It would be interesting to conduct

633 a comparison of results across different heuristics. Last but not the least, the unobserved

634 utility will be incorporated in terms of a random term. Then the stochastic user

635 equilibrium can be adopted to substitute the deterministic user equilibrium.

636

637 **Declarations of interest**

638     None

639 **Data Availability Statement**

640     Some or all data, models, or code that support the findings of this study are available

641 from the corresponding author upon reasonable request.

646    Polytechnic University (ZVTK).

647

648    **References**

649    Abouee-Mehrizi, H., Babri, S., Berman, O., Shavandi, H. (2011) Optimizing capacity,
650    pricing and location decisions on a congested network with balking. *Mathematical*
651    *Methods of Operations Research* 74, 233-255.

652    Auerbach, J. D., Kim, H. (2021). Local network connectivity optimization: an
653    evaluation of heuristics applied to complex spatial networks, a transportation case study,
654    and a spatial social network. *PeerJ Computer Science* 7, e605.

655    Berman, O., Drezner, Z. (2006) Location of congested capacitated facilities with
656    distance-sensitive demand. *IIE Transactions* 38, 213-221.

657    Berman, O., Krass, D., Wang, J. (2006) Locating service facilities to reduce lost
658    demand. *IIE Transactions* 38, 933-946.

659    Castillo, I., Ingolfsson, A., Sim, T. (2009) Social optimal location of facilities with
660    fixed servers, stochastic demand, and congestion. *Production and Operations*
661    *Management* 18, 721-736.

662    Chambari, A., Rahmaty, S. H., Hajipour, V., Karimi, A. (2011) A bi-objective model
663    for location-allocation problem within queuing framework. *International Journal of*
664    *Computer and Information Engineering* 5, 555-562.

665    Dan, T., Marcotte, P. (2019) Competitive facility location with selfish users and
666    queues. *Operations Research* 67, 479-497.

667    Davari, S., Kilic, K., Naderi, S. (2016) A heuristic approach to solve the preventive
668    health care problem with budget and congestion constraints. *Applied Mathematics and*
669    *Computation* 276, 442-453.

670    Dogan, K., Karatas, M., Yakici, E. (2020) A model for locating preventive health care
671    facilities. *Central European Journal of Operations Research* 28, 1091-1121.

672    Erdoğan, G., Stylianou, N., Vasilakis, C. (2019) An open source decision support
673    system for facility location analysis. *Decision Support System* 125, 113116.

674    Ershadi, M.M., Shemirani, H.S. (2021) Using mathematical modeling for analysis of
675    the impact of client choice on preventive healthcare facility network design.
676    *International Journal of Healthcare Managemen* 14, 588-602.

677    Filippi, C., Guastaroba, G., Huerta-Muñoz, D.L., Speranza, M.G. (2021) A kernel
678    search heuristic for a fair facility location problem. *Computers & Operations Research*
679    132, 105292.

Gu, W., Wang, X., McGregor, S.E. (2010) Optimization of preventive health care facility locations. *International Journal of Health Geographics* 9, 17.

Kucukyazici, B., Zhang, Y., Ardestani-Jaafari, A., Song, L. (2020) Incorporating patient preferences in the design and operation of cancer screening facility networks. *European Journal of Operational Research* 287, 616-632.

Marianov, V. (2003) Location of Multiple-Server Congestible Facilities for Maximizing Expected Demand, when Services are Non-Essential. *Annals of Operations Research* 123, 125-141.

Marianov, V., Rios, M., Barros, F.J. (2005) Allocating servers to facilities, when demand is elastic to travel and waiting times. *RAIRO - Operations Research* 39, 143-162.

Marianov, V., Ríos, M., Icaza, M.J. (2008) Facility location for market capture when users rank facilities by shorter travel and waiting times. *European Journal of Operational Research* 191, 32-44.

Mousazadeh, M., Torabi, S.A., Pishvaee, M.S., Abolhassani, F. (2018) Accessible, stable, and equitable health service network redesign: A robust mixed possibilistic-flexible approach. *Transportation Research Part E: Logistics and Transportation Review* 111, 113-129.

Pourrezaie-Khaligh, P., Bozorgi-Amiri, A., Yousefi-Babadi, A., Moon, I. (2022) Fix-and-optimize approach for a healthcare facility location/network design problem considering equity and accessibility: A case study. *Applied Mathematical Modelling* 102, 243-267.

Risanger, S., Singh, B., Morton, D., Meyers, L.A. (2021) Selecting pharmacies for COVID-19 testing to ensure access. *Health Care Management Science* 24, 330-338.

Tao, Z., Cheng, Y., Dai, T., Rosenberg, M.W. (2014) Spatial optimization of residential care facility locations in Beijing, China: Maximum equity in accessibility. *International Journal of Health Geographics* 13, 33.

Tao, Z., Cheng, Y., Dai, T., Zheng, Q. (2016) Application and validation of gravity p-median model in facility location research. *System Engineering Theory and Practice* 36, 1600-1608.

Verter, V., Lapierre, S.D. (2002) Location of Preventive Health Care Facilities. *Annals of Operations Research* 110, 123-132.

Wang, F., Tang, Q. (2013) Planning toward equal accessibility to services: a quadratic programming approach. *Environment and Planning B: Planning and Desig* 40, 195-

714  212.

715  Zhang, W., Cao, K., Liu, S., Huang, B. (2016) A multi-objective optimization
716  approach for health-care facility location-allocation problems in highly developed cities
717  such as Hong Kong. *Computers, Environment and Urban Systems* 59, 220-230.

718  Zhang, Y., Atkins, D. (2019) Medical facility network design: User-choice and
719  system-optimal models. *European Journal of Operational Research* 273, 305-319.

720  Zhang, Y., Berman, O., Marcotte, P., Verter, V. (2010) A bilevel model for preventive
721  healthcare facility network design with congestion. *IIE Transactions* 42, 865-880.

722  Zhang, Y., Berman, O., Verter, V. (2009) Incorporating congestion in preventive
723  healthcare facility network design. *European Journal of Operational Research* 198,
724  922-935.

725

726  **List of Tables**

727

728  **Table 1** Network characteristics for the Sioux Falls network

| Link | Length (*mile*) | Travel time (*hr*) | Link | Length (*mile*) | Travel time (*hr*) |
|---|---|---|---|---|---|
| 1,3 | 3.6 | 0.12 | 33,36 | 3.6 | 0.12 |
| 2,5 | 2.4 | 0.08 | 34,40 | 2.4 | 0.08 |
| 4,14 | 3.0 | 0.10 | 37,38 | 1.8 | 0.06 |
| 6,8 | 2.4 | 0.08 | 39,74 | 2.4 | 0.08 |
| 7,35 | 2.4 | 0.08 | 41,44 | 3.0 | 0.10 |
| 9,11 | 1.2 | 0.04 | 42,71 | 2.4 | 0.08 |
| 10,31 | 3.6 | 0.12 | 45,57 | 2.4 | 0.08 |
| 12,15 | 2.4 | 0.08 | 46,67 | 2.4 | 0.08 |
| 13,23 | 3.0 | 0.10 | 49,52 | 1.2 | 0.04 |
| 16,19 | 1.2 | 0.04 | 50,55 | 1.8 | 0.06 |
| 17,20 | 1.8 | 0.06 | 53,58 | 1.2 | 0.04 |
| 18,54 | 1.2 | 0.04 | 56,60 | 2.4 | 0.08 |
| 21,24 | 6.0 | 0.20 | 59,61 | 2.4 | 0.08 |
| 22,47 | 3.0 | 0.10 | 62,64 | 3.6 | 0.12 |
| 25,26 | 1.8 | 0.06 | 63,68 | 3.0 | 0.10 |
| 27,32 | 3.0 | 0.10 | 65,69 | 1.2 | 0.04 |
| 28,43 | 3.6 | 0.12 | 66,75 | 1.8 | 0.06 |
| 29,48 | 3.0 | 0.10 | 70,72 | 2.4 | 0.08 |

| 30,51 | 4.8 | 0.16 | 73,76 | 1.2 | 0.04 |

729

730

**Table 2** Demand and facility data for the Sioux Falls network

| Population node $i$ | Demand $h_i$ (clients/hr) | Facility location $j$ | Capable space $K_j$ (clients) |
|---|---|---|---|
| 1 | 41 | 3 | 50 |
| 2 | 33 | 7 | 70 |
| 4 | 23 | 9 | 60 |
| 5 | 29 | 11 | 60 |
| 13 | 41 | 16 | 80 |
| 14 | 35 | 19 | 70 |
| 15 | 43 | 21 | 50 |
| 20 | 26 | 23 | 80 |

731

732

**Table 3** The network design scheme and their service level

| Facility location | Server number | Balking probability | Effective arrival rate (clients/hr) | Waiting time (hr) |
|---|---|---|---|---|
| 3 | 6 | 0.132 | 36 | 1.21 |
| 9 | 9 | 0.096 | 54 | 1.12 |
| 16 | 8 | 0.118 | 48 | 1.10 |
| 19 | 7 | 0.091 | 42 | 1.20 |
| 23 | 10 | 0.128 | 60 | 1.22 |

733

734

**Table 4** The demand allocation at equilibrium state

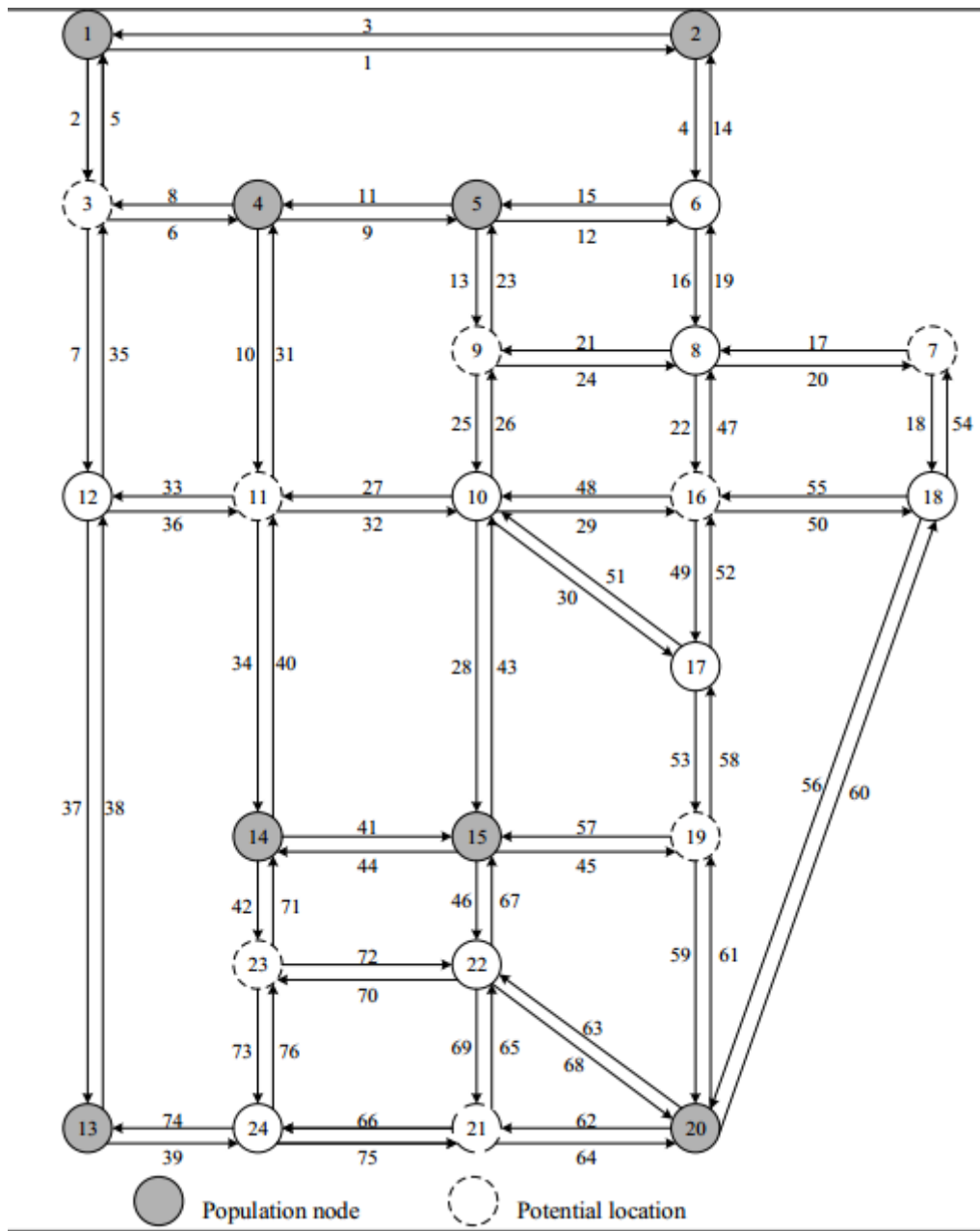| Population node | Selected facility location | | | | |
|---|---|---|---|---|---|
| | 3 | 9 | 16 | 19 | 23 |
| 1 | **31.10** | 3.77 | 1.94 | 1.34 | 2.55 |
| 2 | 2.07 | 5.52 | **21.77** | 2.56 | 1.08 |
| 4 | 2.48 | **16.96** | 1.45 | 1.10 | 1.10 |
| 5 | 1.79 | **21.43** | 2.22 | 1.79 | 1.37 |
| 13 | 3.16 | 1.94 | 1.34 | 2.55 | **31.71** |
| 14 | 0.63 | 3.26 | 1.68 | 5.36 | **24.27** |
| 15 | 0.13 | 5.89 | **11.01** | 21.90 | 3.97 |
| 20 | 0.08 | 0.87 | **13.08** | **9.54** | 2.84 |

735
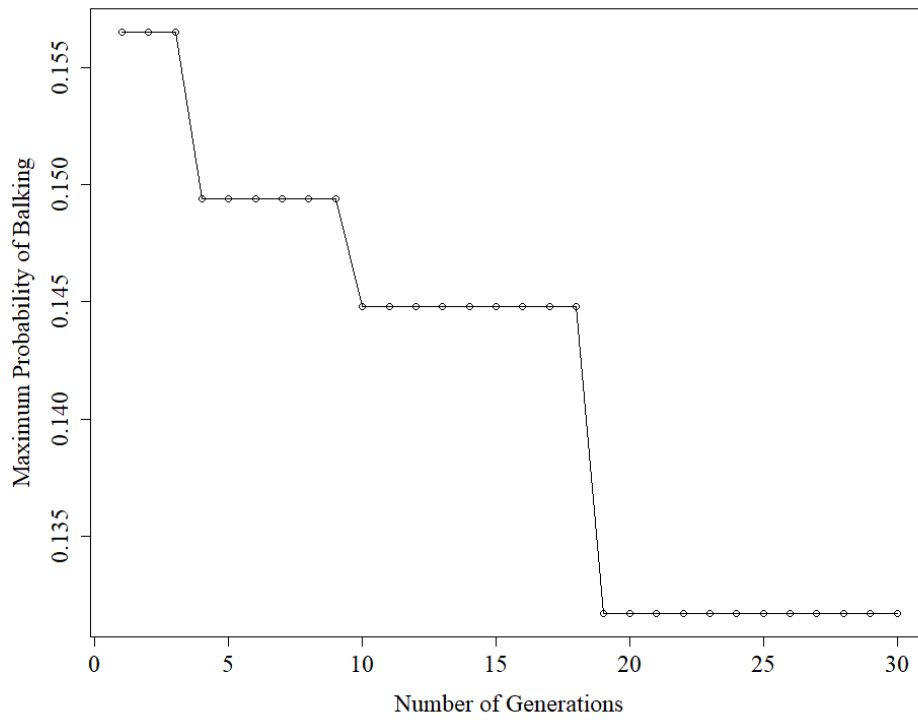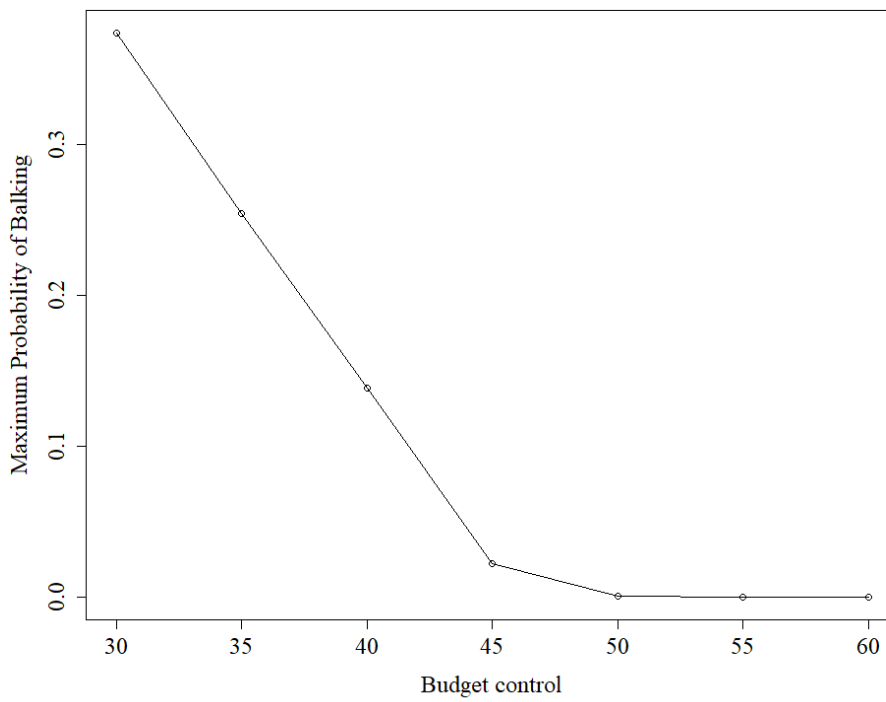
## List of Figures



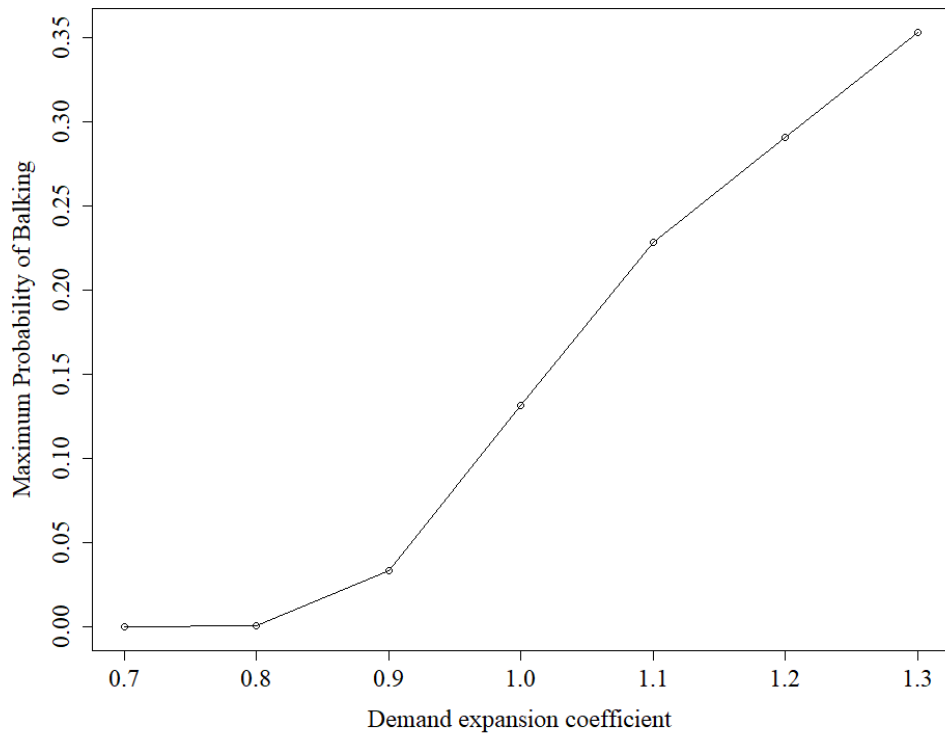**Fig. 1.** The Sioux Falls test network

739

**Fig. 2.** The evolutionary process of genetic algorithm



741

**Fig. 3.** A sensitivity analysis with varying budget

743

**Fig. 4.** A sensitivity analysis with varying demand

745