

EMS location-allocation problem under uncertainties

Abstract

Emergencies, especially those considered routine (i.e., occurring on a daily basis), pose great threats to health, life, and property. Immediate response and treatment can greatly mitigate these threats. This research is conducted to optimize the locations of ambulance stations, deployment of ambulances, and dispatch of vehicles under demand and traffic uncertainty, which are the main factors that influence emergency response time. The research problem is formulated as a dynamic scenario-based two-stage stochastic programming model, aiming to minimize the total cost while responding to as much demand as possible. The Sample Average Approximation is proposed to approximate the original problem using a limited number of scenarios, and a two-phase Benders Decomposition solution scheme is proposed to accelerate computation, especially when solving a large-sized problem. Numerical experiments using real-world emergency data are conducted to validate the performance of the solution method. The results demonstrate the effectiveness and efficiency of the proposed algorithm. We additionally conduct a sensitivity analysis to evaluate the influences of crucial parameters, including the response time standard, facility capacity, service capacity, and facility heterogeneity. The managerial insights derived from sensitivity analysis will provide valuable guidance for the design of an emergency response system in practice.

Keywords: Emergency medical services, Location-allocation problem, Stochastic program, Sample average approximation, Benders Decomposition

1. Introduction

An emergency situation can happen anywhere and anytime, and such situations pose great risks to people's health, lives, and property. A routine emergency, such as a heart attack, road accident, or residential fire, is a type of emergency that can happen on a daily basis and has a small scope of influence. According to research by Vos et al. (2020), this type of emergency is one of the leading causes of death worldwide. Therefore, the immediate response and treatment to a routine emergency are of great importance. In this context, treatment can be categorized as in-hospital treatment or pre-hospital treatment. In-hospital treatment depends heavily on the process design and operations

34 management of patient flows, which is not considered in this research. Readers who are
35 interested in efficiency and effectiveness in the context of in-hospital treatment are
36 referred to Kuo (2014) and Kuo et al. (2016). Pre-hospital treatment usually involves
37 an emergency medical services (EMS) department. When an emergency call is received,
38 the call handler determines the severity of the situation and then dispatches response
39 vehicles accordingly. As the patient survival rate is highly dependent on the response
40 time, defined as the time interval between the reception of a call and the arrival of the
41 response vehicle at the emergency site, the dispatched vehicles must arrive at the
42 emergency site within certain time limits (Bürger et al., 2018; Erkut, Ingolfsson, &
43 Erdoğan, 2008; Knight, Harper, & Smith, 2012). The response time is directly affected
44 by the locations of ambulance stations, number of available vehicles, and dispatching
45 decisions. Therefore, optimizing the factors that directly influence the response time is
46 essential to guarantee that efficient and high-quality EMS services are available for the
47 public.

48 The optimization problem mentioned above is called a location-allocation problem.
49 Such problems are challenging because they usually are formulated as mixed-integer
50 programming models containing many decision variables and constraints. The
51 inclusion of uncertainty further complicates this type of problem. This research
52 addresses a location-allocation problem with system congestion under demand and
53 traffic uncertainty. System congestion refers to a situation where the available response
54 vehicles are insufficient to respond to demand. In this article, system congestion is
55 captured by vehicle availability, i.e., the number of vehicles that can respond to demand,
56 and is influenced by the dispatch of vehicles to demands that overlap in time with the
57 current demand. The uncertainty in demand includes the number of emergencies,
58 locations, occurrence times, numbers of ambulances needed, and the service time
59 needed, which directly influence optimization decisions. The service time is the time
60 interval between the arrival of a vehicle at the emergency site and its return to the station.
61 Traffic uncertainty is mainly represented by travel time, which influences the coverage
62 set of each demand, i.e., the subset of facility sites that can cover demand within time
63 standard. These uncertainties are captured by scenarios, and the problem is formulated
64 as a dynamic two-stage stochastic model. In the first stage, the model determines the
65 optimal locations of ambulance stations and deployment of ambulances without
66 considering realized uncertainties. In the second stage, resource decisions on vehicle

67 dispatching are made according to the realized uncertainties, first-stage decisions, and
68 status of available vehicles. The objective of the model is to minimize the total cost
69 while responding to as much demand as possible. The total cost comprises the station
70 set-up cost, vehicle purchasing cost, demand fulfillment cost, and penalties for failing
71 to respond to demand within the required time standard and failing to satisfy demand.
72 We first apply Sample Average Approximation (SAA) to approximate the original
73 problem using samples. We next propose a two-phase Benders Decomposition solution
74 scheme (TPBD) to accelerate the computation. To evaluate the performance of our
75 proposed methods, we conduct numerical experiments using real-world emergency data.
76 We also conduct a sensitivity analysis to evaluate the influences of crucial parameters,
77 including the response time standard, facility capacity, service capacity, and facility
78 heterogeneity. We obtain managerial insights that will provide valuable guidance for
79 the design of an emergency response system in practice.

80 The rest of this paper is organized as follows. Section 2 reviews relevant literature
81 and illustrates contribution. The problem description is presented in Section 3 where
82 the deterministic model is first introduced and then extended to consider random
83 demands and travel time. Section 4 introduces solution approach. We conduct
84 numerical experiments in Section 5. Section 6 concludes this research.

85

86 **2. Literature Review**

87 The research in routine EMS can be divided into many branches, such as dispatching,
88 location, deployment, assignment, and relocation. Dispatching problem decides which
89 ambulances are dispatched to serve each demand. Location problem determines where
90 to set up stations that can host ambulances. Deployment problem optimizes the number
91 of ambulances hosted at each station. Assignment problem assigns serving stations to
92 demand. Relocation problem identifies where ambulances are moved to. Readers who
93 are interested in these topics can refer to Aringhieri, Bruni, Khodaparasti, and van Essen
94 (2017) and Bélanger, Ruiz, and Soriano (2019). Due to complexity of dispatching and
95 location problem, a lot of research investigates these two problems separately. However,
96 as all of the branches are essential components for an efficient EMS system, research
97 that makes a combination of these branches gains popularity in recent decades. The
98 combination further complicates the problem, especially when uncertainty is
99 considered. Therefore, in this section, we first introduce dispatching and location

100 problem separately. Then we comprehensively review research that integrates several
101 types of decisions. Finally, we analyze research gaps in existing literature and illustrate
102 how this research fills the gaps.

103

104 2.1 Dispatching Problem

105 According to Lee (2012), dispatching can be divided into two types: call-initiated
106 and server-initiated. Call-initiated dispatching is to choose an appropriate ambulance to
107 respond to an emergency call, while server-initiated dispatching is to decide the demand
108 in the waiting list to be served when a server is available. Regardless of the type of
109 dispatching, decision makers usually adopt certain dispatching policy when making
110 dispatching decisions. The most commonly used dispatching strategy is nearest
111 available policy, which means that an emergency call is most likely to be served by
112 available vehicles that are closest to it (Dean, 2008; Lee, 2011; Zarkeshzadeh, Zare,
113 Heshmati, and Teimouri, 2016). In addition to nearest available policy, several other
114 rules are proposed based on characteristics of the problem. Lee (2011) adopts the
115 concept of preparedness, a quantitative function of the number of available ambulances
116 and call rate, when designing dispatching algorithm for ambulance services. Lee (2012,
117 2013) define centrality that reflects the density of emergency calls and propose a
118 dispatching policy based on this definition. McLay and Mayorga (2013a) proposes
119 Markov decision process that dispatches distinguished ambulances (i.e., ambulances
120 with different response and service time) to prioritized demand and at the same time
121 considers estimation error of patient priority to calculate the optimal dispatching
122 policies. The purpose is to maximize the expected coverage of true high-risk calls. It is
123 extended by McLay and Mayorga (2013b) to consider both efficiency and equity.
124 Efficiency is represented by expected coverage of high-priority demand. Equity is
125 captured by four types of equity constraints, two of which reflect customer equity and
126 the remaining two reflect server equity. Sudtachat, Mayorga, and McLay (2014) further
127 extends the problem by dispatching two types of ambulances to three-priority-level
128 demand. Zarkeshzadeh, Zare, Heshmati, and Teimouri (2016) develops a weighted
129 hybrid method that combines centrality, nearest neighbor, and first-in-first-out into one
130 model to take advantage of each method. After developing different strategies,
131 performance evaluation is also important. Haghani, Tian, and Hu (2004) uses
132 simulation to evaluate three response strategies, namely the first called first served

133 strategy, the nearest origin assignment strategy, and the flexible assignment strategy
134 that uses real-time traffic information. Bandara, Mayorga, and McLay (2014) also tests
135 different response strategies to find the optimal dispatching strategy for EMS systems
136 considering demand priority. The above models are formulated under deterministic
137 environment. Jenkins, Robbins, and Lunday (2021) optimizes dispatch of military
138 medical evacuation assets considering uncertain demand where the uncertainty is
139 represented by scenarios. The problem is formulated as a discounted, infinite-horizon
140 Markov decision process model and solved by two approximate dynamic programming
141 methods.

142

143 2.2 Location Problem

144 Research on EMS location problem has a long history which can date back to 1970s.
145 Most of the studies are extensions of two classic coverage models: location set covering
146 problem (LSCP) and maximal covering location problem (MCLP). LSCP is first
147 proposed by Toregas, Swain, ReVelle, and Bergman (1971), which minimizes the
148 number of facilities to cover all demands. The requirement of a mandatory coverage of
149 all demand points may be impossible to implement under some situations, such as
150 budget shortage. Church and ReVelle (1974) proposes MCLP to maximize the demand
151 coverage under facility number constraint. When designing fire protection system
152 where two different types of equipment have to be considered, Schilling, Elzinga,
153 Cohon, Church, and ReVelle (1979) changes definition of coverage in MCLP and
154 proposes FLEET model that requires demand be covered only when it is simultaneously
155 within distance standard of two types of equipment. Daskin and Stern (1981) extends
156 location problem to consider system congestion, which allows demands to be covered
157 by multiple locations so that even the nearest vehicles are engaged, other vehicles
158 within coverage radius can serve demands. Gendreau, Laporte, and Semet (1997)
159 considers redundant coverage and proposes double coverage model (DCM), in which
160 two response distances are considered. All demands are required to be covered within
161 the larger response distance and at the same time a proportion of demands must be
162 covered within smaller response distance.

163 The models in above articles are deterministic, which can obtain optimality or near-
164 optimality in simplified assumptions of the real-life practices. However, as the
165 operational environment keeps changing, the deterministic inputs may cause biased

166 results. Therefore, probability is added into the model to represent system instability.
167 Daskin (1983) is one of the early research that adopts busy fraction, i.e., the probability
168 that a server (i.e., ambulance) cannot respond to the demand within time requirements,
169 and proposes a model called MEXCLP, which maximizes the expected coverage.
170 ReVelle and Hogan (1988, 1989a, 1989b) embed busy fraction into chance constraints
171 that require that the probability of the demand being responded is no less than a
172 reliability level, resulting in a problem called maximum availability location problem
173 (MALP). Sorensen and Church (2010) combines MALP with MEXCLP and compares
174 this new model with the two original models in a range of test problems. Liu, Li, Liu,
175 and Patel (2016) combines MALP with DCM to maximize coverage of demand at
176 guaranteed service reliability in a primary distance standard and at the same time to
177 ensure a full coverage in a secondary distance standard. When traffic situation is
178 uncertain, Goldberg and Paz (1991) considers the distribution of the travel time when
179 locating ambulance stations to maximize the expected coverage. The distribution
180 determines the probability that the demand could be responded by the vehicle at certain
181 station within the threshold time. Schmid and Doerner (2010) develops multi-period
182 model to take into account time dependent speed. Berman, Hajizadeh, and Krass (2013)
183 uses MEXCLP where travel time uncertainty is represented by scenarios with certain
184 probability to maximize expected coverage. El Itani, Abdelaziz, and Masri (2019)
185 considers the combination of MEXCLP and MALP and proposes a bi-objective model
186 that simultaneously maximizes expected coverage and minimizes expected cost when
187 paying for external ambulances is allowed. In addition to system efficiency, some
188 research considers equity of the system. Chanta, Mayorga, Kurz, and McLay (2011)
189 defines the concept of envy to model equity when location ambulance stations.
190 Customer envy is calculated based on the distance between demand area and stations
191 on a pre-determined preference list. The objective is to minimize weighted envy where
192 demand density and vehicle availability obtained through queuing theory are two
193 weights used in objective function. Chanta, Mayorga, and McLay (2014)
194 simultaneously considers efficiency and equity by developing a bi-objective covering
195 location model where the former is captured by the first objective and the latter is
196 represented by one of three second objectives at a time.

197

198 2.3 Hybrid Problem

199 The combination of several problems will complicate research. In early days, the
200 combination is usually done at the same level, e.g., strategic location, deployment and
201 assignment, or operational dispatching, deployment, and relocation are combined in
202 one research. Ingolfsson, Budge, and Erkut (2008) optimizes the deployment of
203 ambulances to stations and the assignment to demand under random pre-travel delay,
204 travel time, and vehicle availability. The randomness in the pre-travel delay and travel
205 time is captured by the deviation from the mean time. Beraldi and Bruni (2009)
206 innovatively incorporates joint probabilistic chance constraints into the traditional two-
207 stage stochastic programming model to explore base station location, fleet size, and
208 ambulance assignment problem for EMS under demand uncertainty. van den Berg and
209 Aardal (2015) extends the MEXCLP into a multi-period version with the goal to
210 maximize the expected coverage, minimize start-up cost, and minimize penalty for
211 relocation throughout the day. Degel, Wiesche, Rachuba, and Werners (2015) also uses
212 time-dependent parameters to obtain the more precise and practical solutions for
213 maximum demand coverage. To improve the coverage level and system efficiency, the
214 relocation and additional flexible stations are considered in the model. Liu, Li, and
215 Zhang (2019) uses two-stage distributionally robust model with joint chance constraints
216 to optimize location and deployment considering two types of uncertainties related with
217 demand. Boutilier and Chan (2020) uses two-stage robust optimization model, which
218 makes sure that the worst-case solution is also optimal, to determine location and
219 routing of emergency response vehicles in low- and middle-income countries under
220 demand and travel time uncertainty.

221 Vehicles are dispatched to satisfy certain demand. To make sure that enough demands
222 are served, the number of ambulances deployed at each station is often jointly optimized
223 with dispatching. Bertsimas and Ng (2019) uses both stochastic and robust two-stage
224 models to solve ambulance dispatching and deployment problem under demand
225 uncertainty. It requires that the number of vehicles dispatched do not exceed the total
226 number deployed. The uncertainty set of robust optimization is calculated based on
227 data-driven approach. When a vehicle is dispatched, the location of the vehicle is empty,
228 reducing protection for surrounding areas. This phenomenon is especially severe in
229 areas with high demand density. An effective method to improve the situation is to
230 relocate vehicles from other less busy stations. Nasrollahzadeh, Khademi, and Mayorga
231 (2018) optimizes real-time ambulance dispatching and relocation, which is formulated

232 as an infinite-horizon Markov decision process. The model is solved by approximate
233 dynamic programming. Park, Waddell, and Haghani (2019) optimizes dispatch of
234 emergency vehicles in freeway under randomness of requests. Different from research
235 that only looks at past and current demand information, this article further looks ahead
236 a short-term future demand based on incident distribution to dispatch and relocate
237 vehicles. A dynamic programming based method is proposed to solve the problem.

238 As researchers have deeper understanding about hybrid problem, more complicated
239 combinations (e.g., location and dispatching) are taken into account. Toro-Díaz,
240 Mayorga, Chanta, and McLay (2013) integrates mixed-integer programming model for
241 location and dispatching and hypercube queuing model for system congestion
242 considering fixed priority list for each demand area. Nickel, Reuter-Oppermann, and
243 Saldanha-da-Gama (2016) uses two-stage stochastic programming model to optimize
244 ambulance location, fleet deployment, and the number of vehicles dispatched from
245 stations to serve demands under demand uncertainty. Boujemaa et al. (2018) extends
246 the problem to consider two types of vehicles. Nelas and Dias (2020) proposes a new
247 integer linear programming model that allows vehicle substitution and considers system
248 congestion. Bélanger et al. (2020) proposes a recursive simulation-optimization
249 framework that iterates between an integer programming model and a discrete event
250 simulation model. The integer programming model determines optimal ambulance
251 location and dispatching list for each demand area under given response probability.
252 The discrete event simulation model dispatches vehicles and updates response
253 probability under solution obtained from integer programming model. Peng, Delage,
254 and Li (2020) extends the problem to multiperiod and proposes envelop constraints to
255 guarantee coverage under extreme scenarios. Yoon, Albert, and White (2021) improves
256 solution technique for two-stage stochastic programming model.

257

258 2.4 Research Gap and Contribution

259 As location and dispatch belong to different decision levels, the combination of two-
260 level problems is computationally intensive. Consequently, little relevant hybrid
261 research is available. This article addresses location, dynamic real-time dispatching,
262 and fleet deployment, which are rarely explored in combination in the literature.
263 Location and fleet deployment problems are usually formulated as mixed-integer
264 programming models, while dispatching problems are usually solved by drawing on

265 queuing theory, which considers dispatch policies and preference lists. However, we
266 use an innovative mixed-integer programming model to formulate the hybrid problem.
267 We further incorporate system congestion into the model when selecting vehicles to be
268 dispatched. In the literature, busy fraction and queuing theory are the most commonly
269 used methods to model system congestion. Busy fraction is generally assumed to be
270 fixed, independent, and exogenous and thus cannot reflect the dynamic and endogenous
271 characteristics of system congestion. When dispatch is modeled by Markov decision
272 process, to simplify the computation, it usually adopts certain assumptions (e.g., arrival
273 and service process) and dispatch policies (e.g., nearest available, preparedness, or
274 centrality). In our research, we aim to establish a dynamic and endogenous dispatch
275 strategy without excessive assumptions or preset dispatching policies. Rather than
276 using busy fraction or queuing theory, we represent system congestion by the functions
277 of a parameter indicating the overlap between demands, deployment decisions, and
278 previous dispatch decisions. The dispatch policy is not pre-defined but rather
279 determined by an objective function. We also incorporate both demand and travel time
280 uncertainties into the model, whereas most research considers neither or only one of
281 these variables. Demand and travel time uncertainties are usually represented by
282 scenarios and the time-dependent travel time, respectively. We adopt the scenario
283 method to represent demand uncertainty, as accurate demand information in each
284 scenario is helpful to model vehicle availability. We represent travel time uncertainty
285 by combining scenario and multi-period methods. We first generate scenarios, each of
286 which represents traffic information from one day. We then divide each scenario into
287 24 equal 1-hour segments and calculate the segment-dependent travel time. This
288 approach to travel time uncertainty is helpful for calculating the actual coverage set of
289 each demand, which is essential for calculating the total cost. Finally, we propose a new
290 algorithm to accelerate computation, especially for large-size problems.

291

292 **3. Problem Description and Formulation**

293 This research deals with ambulance location-allocation problem with vehicle
294 availability under demand and traffic uncertainty. We first present a deterministic model
295 in Section 3.1 to make it easier to understand the reasoning and logic behind the model.
296 Then in Section 3.2, a scenario-based two-stage stochastic model is proposed to deal
297 with uncertainties.

298

299 3.1 Deterministic Model

300 In this section, we assume that the information about demand and traffic situation is
301 known a priori. We denote all demands by an ordered set $I = \{1, \dots, |I|\}$ where
302 demand $i - 1$ occurs before demand i , $i \in I \setminus \{1\}$. The whole research region is
303 divided into several zones, each of which is represented by its centroid. The location of
304 demand i is the centroid of the zone where i occurs. As the model is deterministic,
305 the occurrence time t_i , the number of required vehicles d_i , and the required service
306 time l_i (time interval between the arrival of a vehicle at the location of demand i and
307 its return to station) are also given.

308 The set of candidate facilities is denoted by J . Note that in this research, facility, site,
309 and station refer to the same thing and are used interchangeably. The location of each
310 facility j , $j \in J$ is given. One decision of this research is to determine which sites to
311 open, represented by a binary variable x_j , which equals 1 if facility j is open and 0
312 otherwise. Once facility j is open, it will incur a set-up cost f_j . The number of
313 vehicles deployed at open station j is denoted by y_j which is a decision variable.
314 Each vehicle is purchased at a cost h and can serve any demand, but will be penalized
315 if it is located outside the coverage set N_i of demand i . The coverage set of demand
316 i is the set of facilities that can cover demand i within response time standard R ;
317 $N_i \subseteq J$. Each vehicle has the workload limit, the maximum number of demands the
318 vehicle can serve. The purpose of this limit is to balance the workload between stations.

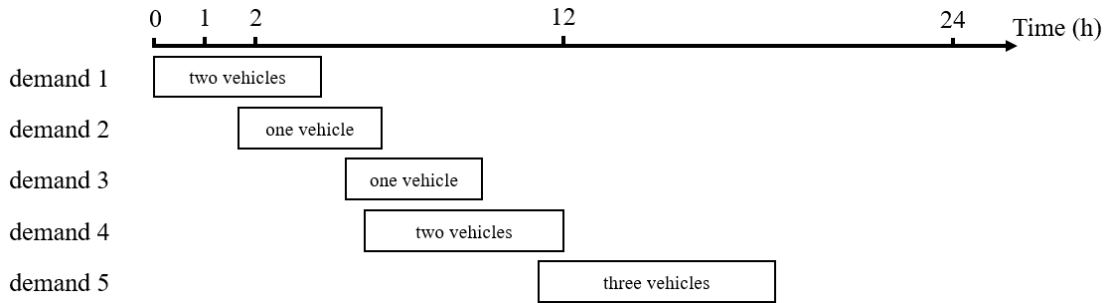
319 We adopt first-come-first-serve policy for demand service and allow demands to be
320 partially served. The first-come-first-serve policy is widely used in the EMS response,
321 which means if there is an overlap between the service time of two demands requiring
322 the vehicle from the same station, the one that comes first will be served by the vehicle
323 at this station, while the one that comes later will experience one of the three situations:
324 1) It will be served by vehicles from the same station if there are enough available
325 vehicles left; 2) It will be served by available vehicles from another station; 3) It will
326 be partly served by vehicles from the same station and the rest is served by available
327 vehicles from another station. This can be illustrated by an example in Figure 1, which
328 shows 5 demands. The left side, right side, and length of rectangle represent the
329 occurrence time, finish time, and time interval from occurrence to completion,
330 respectively. The number of vehicles needed is also shown in the figure. These demands

331 will be served by vehicles from two stations, each of which hosts two vehicles. We
 332 assume that station 1 is preferred to all the demands than station 2. The number of
 333 available vehicles to each demand and the dispatching decisions are shown in Table 1.
 334 Demand 1 happens first, so two vehicles from station 1 are dispatched to demand 1.
 335 When demand 2 occurs, these two vehicles are still engaged in the last service, so the
 336 vehicle from the less preferred station 2 is dispatched. The same rule is applicable to
 337 demand 3 and 4. When demand 5 occurs, only two vehicles are available, but it requires
 338 three vehicles. Thus, two vehicles are dispatched and demand 5 is partially served. To
 339 calculate the vehicle availability, we need to define a binary parameter $\delta_{i'j}$, which
 340 indicates whether there is overlapping time between demand i and i' for vehicles at
 341 station j to serve them.

$$342 \quad \delta_{i'j} = \begin{cases} 0, & \text{if } t_i \geq t_{i'} + T_{i'j} + l_{i'} \\ 1, & \text{otherwise} \end{cases}, \quad \forall i \in I, j \in J, i' \in \{1, 2, \dots, i-1\}. \quad (1)$$

343 Notice that when $\delta_{i'j} = 0$, demand i and i' are disjoint. The service to i' will not
 344 influence the available vehicles to i . However, when $\delta_{i'j} = 1$, the allocation decision
 345 of i' has influence on vehicle availability to i .

346



347

348 Figure 1. One example of demand in one day period

349

350 Table 1. One example of vehicle allocation

	Station 1		Station 2	
	Available	Allocated	Available	Allocated
d1	2	2	2	0
d2	0	0	2	1
d3	2	1	1	0
d4	1	1	1	1
d5	1	1	1	1

351

352 The goal of the research is to identify the optimal location of stations, the number
 353 and deployment of vehicles, and the allocation of demand at a minimum total cost. All
 354 the notations for the deterministic model are introduced in Table 2 and the model is
 355 given as follows:

356

357

Table 2. Notations for deterministic model

Sets	
I	The set of a sequence of demand, where the demand $i - 1$ happens before i
J	The set of candidate facility sites
N_i	The coverage set of demand i , i.e., the set of facility sites that can cover demand i within time standard ($N_i = \{j \in J, T_{ij} \leq R\}$)
Parameters	
f_j	The fixed cost of opening station j
h	The purchasing cost of a vehicle
c_{ij}	The unit transportation cost for vehicle at station j to serve demand i
c^{lost}	The unit penalty of demand unsatisfaction
γ	The unit penalty for violating the response time standard
T_{ij}	The travel time for vehicle at station j to serve demand i
R	The response time standard
Q_j	The maximal number of vehicles that can be hosted at station j
d_i	The number of vehicles needed for demand i
μ_j	The maximum number of demands each vehicle at station j could serve
t_i	The occurrence time of demand i
l_i	The service time of demand i , the time needed after vehicle arriving at the emergency site until the vehicle goes back to station again
$\delta_{i'j}$	1 if there is overlapping time between demand i and i' for vehicles at station j to serve them, 0 otherwise
Decision variables	
x_j	1 if a station is set up at site j , 0 otherwise
y_j	The number of vehicles hosted at location j
z_{ij}	The number of vehicles at location j that are dispatched to demand i

358

$$\begin{aligned}
359 \quad [M1] \quad \text{Min} \quad & \sum_{j \in J} f_j x_j + \sum_{j \in J} h y_j + \sum_{i \in I} \sum_{j \in J} c_{ij} T_{ij} z_{ij} + \sum_{i \in I} \sum_{j \in J \setminus N_i} \gamma (T_{ij} - R) z_{ij} + \\
360 \quad & c^{lost} (\sum_{i \in I} d_i - \sum_{i \in I} \sum_{j \in J} z_{ij}) \tag{2}
\end{aligned}$$

361 subject to

$$362 \quad y_j \leq Q_j x_j, \quad \forall j \in J \tag{3}$$

$$363 \quad \sum_{j \in J} z_{ij} \leq d_i, \quad \forall i \in I \tag{4}$$

$$364 \quad z_{ij} \leq y_j - \sum_{i'=1}^{i-1} \delta_{i'ij} z_{i'j}, \quad \forall i \in I, j \in J \tag{5}$$

$$365 \quad \sum_{i \in I} z_{ij} \leq \mu_j y_j, \quad \forall j \in J \tag{6}$$

$$366 \quad x_j \in \{0,1\}, \quad \forall j \in J \tag{7}$$

$$367 \quad y_j \in \mathbb{Z}_0^+, \quad \forall j \in J \tag{8}$$

$$368 \quad z_{ij} \in \mathbb{Z}_0^+, \quad \forall i \in I, j \in J. \tag{9}$$

369 The objective function (2) minimizes total cost, which consists of station set-up cost,
370 vehicle purchasing cost, demand fulfillment cost, the penalty for not responding the
371 demand in required time standard, and penalty for demand unsatisfaction. For the third
372 term, we only account for the vehicles located outside the coverage set because only
373 these vehicles cannot respond to demand within time standard. Constraints (3) require
374 that vehicles can only be located at open station and the number cannot exceed the
375 station capacity. Constraints (4) state that demand can be partially satisfied. Constraints
376 (5) require that only vehicles available at station when demand occurs are available to
377 serve it. The second term on the right-hand side is the number of unavailable vehicles
378 when demand i occurs, which is determined by $\delta_{i'ij}$ and $z_{i'j}$. $\delta_{i'ij}$ indicates
379 whether different demands would have overlap if they are allocated to the same station,
380 which is illustrated by Equation (1). If when demand i occurs, any service to demand
381 i' that occurs before demand i has already been completed, there is no overlap
382 between i and i' and $\delta_{i'ij} = 0$. If when demand i occurs, there are demands being
383 served, the overlap exists and $\delta_{i'ij} = 1$. $z_{i'j}$ is the number of vehicles at location j
384 that are dispatched to demand i' that occurs before demand i . When $\delta_{i'ij} = 0$, i.e.,
385 there is no overlap between demand i and i' , whatever the decision $z_{i'j}$ is, the
386 available vehicles at station j will not be influenced. When $\delta_{i'ij} = 1$, i.e., there is
387 overlap between demand i and i' , if $z_{i'j} = 0$, the available vehicles at station j will
388 not be influenced. If $z_{i'j} > 0$, $z_{i'j}$ number of the vehicles are engaged in demand i'
389 when i occurs, thus the available vehicles at station j will be reduced by $z_{i'j}$. Then
390 the total number of unavailable vehicles because of the preexisting demands is
391 $\sum_{i'=1}^{i-1} \delta_{i'ij} z_{i'j}$. Constraints (6) state that the demands served by each station cannot

392 exceed the workload of the vehicle at this station. Constraints (7) to (9) set the domain
393 of decision variables.

394 **3.2 Stochastic Model**

395 One of the important limitations of deterministic model is the assumption that all the
396 parameters are known in advance. However, in practice, it is hard to know what will
397 happen in the future, especially the demand and traffic condition. For this reason, we
398 have to develop a model that could take the uncertainty into consideration. The model
399 developed in this section is a scenario-based two-stage stochastic programming model.
400 In the first stage, the model determines the optimal location of ambulance stations, fleet
401 size, and the deployment of ambulances without considering the realization of
402 uncertainties. In the second stage, recourse decisions on ambulances dispatching are
403 made based on scenarios, first-stage decisions, and state of available vehicles. We
404 denote by S the set of scenarios. Each scenario $s \in S$ contains the information of
405 demand and traffic situation during a one-day period and is associated with a probability
406 of occurrence p_s . The information includes a sequence of demands that happen during
407 the day, their location and occurrence time, the travel time between candidate stations
408 and demand sites when the demand occurs, and the service time. Under the changing
409 scenarios, the value of some parameters and variables may change accordingly. The
410 demand coverage set N_i is different because the demand location and the travel time
411 between demand and candidate locations change. Parameter $\delta_{ii'j}$ and decision
412 variables z_{ij} and σ_i also change. The definition of $\delta_{ii'j}$ under scenario s is given
413 as follows:

$$414 \delta_{ii'j}^s = \begin{cases} 0, & \text{if } t_i^s \geq t_{i'}^s + T_{i'j}^s + l_{i'}^s, \\ 1, & \text{otherwise} \end{cases}, \forall i \in I, j \in J, i' \in \{1, 2, \dots, i-1\}, s \in S. \quad (10)$$

415 The objective function turns to calculate the minimum expected total cost. The
416 additional parameters and variables for stochastic model are listed in Table 3 and the
417 scenario-based two-stage stochastic programming model is as follows

Table 3. Additional notations in stochastic model

Sets	
S	The set of scenarios
I^s	The set of a sequence of demand under scenario s
N_i^s	The coverage set of demand i under scenario s
Parameters	
T_{ij}^s	The travel time for vehicle at station j to serve demand i under scenario s
d_i^s	The number of vehicles needed for demand i under scenario s
t_i^s	The occurrence time of demand i under scenario s
l_i^s	The service time of demand i under scenario s
$\delta_{i'ij}^s$	1 if there is overlapping time between demand i and i' for vehicles at station j to serve them under scenario s , 0 otherwise
Decision variables	
z_{ij}^s	The number of vehicles at location j that are dispatched to demand i under scenario s

419

$$420 \quad \text{Min } \sum_{j \in J} f_j x_j + \sum_{j \in J} h y_j + \sum_{s \in S} p_s \left(\sum_{i \in I^s} \sum_{j \in J} c_{ij} T_{ij}^s z_{ij}^s + \sum_{i \in I^s} \sum_{j \in J \setminus N_i^s} \gamma (T_{ij}^s - \right. \\ 421 \quad \left. R) z_{ij}^s + c^{lost} (\sum_{i \in I^s} d_i^s - \sum_{i \in I^s} \sum_{j \in J} z_{ij}^s) \right) \quad (11)$$

422 subject to (3), (7), (8)

$$423 \quad \sum_{j \in J} z_{ij}^s \leq d_i^s, \quad \forall i \in I^s, s \in S \quad (12)$$

$$424 \quad z_{ij}^s \leq y_j - \sum_{i'=1}^{i-1} \delta_{i'ij}^s z_{i'j}^s, \quad \forall i \in I^s, j \in J, s \in S \quad (13)$$

$$425 \quad \sum_{i \in I^s} z_{ij}^s \leq \mu_j y_j, \quad \forall j \in J, s \in S \quad (14)$$

$$426 \quad z_{ij}^s \in \mathbb{Z}_0^+, \quad \forall i \in I^s, j \in J, s \in S. \quad (15)$$

427

428 4. Solution Approach

429 The scenario-based two-stage stochastic programming model is challenging to solve
430 because in real life the demand and traffic situation are changing all the time, resulting
431 in a large number of scenarios, which makes the problem computationally intractable.

432 In this section, we will first introduce SAA, which is used to approximate the original
433 problem with samples. Then a Benders Decomposition based approach is proposed to
434 speed up the computation.

435

436

437 4.1 Sample Average Approximation

438 The SAA method is a Monte Carlo simulation-based approach to solve stochastic
 439 optimization problems. The basic idea of this approach is to approximate the true
 440 distribution by empirical distribution obtained from samples. The sample is represented
 441 by S' , which is a finite set of scenarios sampled from
 442 S with the same probability of occurrence, i.e., $S' \subseteq S$, $|S'|$ is the sample size, and
 443 each scenario in S' has the same probability $1/|S'|$. The SAA formulation is given as
 444 follows:

$$445 \text{ [SAA] Min } \sum_{j \in J} f_j x_j + \sum_{j \in J} h_j y_j + \sum_{s \in S'} \frac{1}{|S'|} \left(\sum_{i \in I^s} \sum_{j \in J} c_{ij} T_{ij}^s z_{ij}^s + \sum_{i \in I^s} \sum_{j \in J \setminus N_i^s} \gamma (T_{ij}^s - \right. \\ 446 \left. R) z_{ij}^s + c^{lost} (\sum_{i \in I^s} d_i^s - \sum_{i \in I^s} \sum_{j \in J} z_{ij}^s) \right) \quad (16)$$

447 subject to (3), (7), (8), (12)–(15) in which S is replaced by S' .

448 When using SAA to solve the problem, one essential procedure is to determine the
 449 number of scenarios in S' . The solution quality will be improved with the increase of
 450 sample size, while the model will become computationally intractable. We need to
 451 strike a balance between precision and computational tractability. Algorithm 1
 452 describes the procedure to evaluate the solution quality of SAA under given sample
 453 size, which includes the calculation of confidence intervals (CIs) for lower bound,
 454 upper bound, and optimality gap under given sample size. We will discuss the choice
 455 of $|S'|$ in Section 5 using real-world emergency data.

456

Algorithm 1 Estimate $(1 - \tau)$ -CI for lower bound, upper bound, and optimality gap of two-stage stochastic program

1. Generate a set of scenarios S' .
2. Solve the SAA with S' and obtain the optimal first-stage solution x^*, y^* .
3. **for** $m = 1, 2, \dots, M$ **do**
4. Generate a set of new independent scenarios S_m , $|S_m| = |S'|$.
5. Solve SAA with S_m and obtain the objective value v_m .
6. Generate a set of new independent scenarios S'_m , $|S'_m| \gg |S_m|$.
7. Evaluate the quality of the first-stage solution x^*, y^* on scenarios in S'_m . Input x^*, y^* into SAA, obtaining cost v_{x^*, y^*}^m .
8. Let $g_m := v_{x^*, y^*}^m - v_m$.
9. **end for**

10. Estimate $(1 - \tau)$ -CI for lower bound

11. Let $L := \frac{1}{M} \sum_{m=1}^M v_m$ and $S_L := \frac{1}{M-1} \sum_{m=1}^M (v_m - L)^2$.

12. The $(1 - \tau)$ -CI for lower bound is $\left[L - \frac{t_{M-1, \frac{\tau}{2}} \sqrt{S_L}}{\sqrt{M}}, L + \frac{t_{M-1, \frac{\tau}{2}} \sqrt{S_L}}{\sqrt{M}} \right]$, $t_{M-1, \frac{\tau}{2}}$ is the t-value

obtained from t-distribution with degrees of freedom $M - 1$ and confidence level $1 - \tau$.

13. Estimate $(1 - \tau)$ -CI for upper bound

14. Let $U := \frac{1}{M} \sum_{m=1}^M v_{x^*, y^*}^m$ and $S_U := \frac{1}{M-1} \sum_{m=1}^M (v_{x^*, y^*}^m - U)^2$.

15. The $(1 - \tau)$ -CI for upper bound is $\left[U - \frac{t_{M-1, \frac{\tau}{2}} \sqrt{S_U}}{\sqrt{M}}, U + \frac{t_{M-1, \frac{\tau}{2}} \sqrt{S_U}}{\sqrt{M}} \right]$.

16. Estimate $(1 - \tau)$ -CI for optimality gap

17. Let $G := \frac{1}{M} \sum_{m=1}^M g_m$ and $S_G := \frac{1}{M-1} \sum_{m=1}^M (g_m - G)^2$.

18. The $(1 - \tau)$ -CI for optimality gap is $\left[0, G + \frac{t_{M-1, \frac{\tau}{2}} \sqrt{S_G}}{\sqrt{M}} \right]$.

457

458 4.2 Benders Decomposition Based Solution Scheme

459 Benders decomposition is efficient in dealing with complicated large-scale two-stage
460 stochastic programming problem. The basic idea is to decompose the scenario-based
461 two-stage stochastic programming model into a master problem and a number of
462 subproblems. Then these two types of problems are solved iteratively, adding additional
463 constraints to the master problem. Usually the second-stage model under each scenario
464 is regarded as an independent subproblem, see Peng, Delage, and Li (2020). In this case,
465 many Benders cuts, one for each subproblem, are added at each iteration to accelerate
466 the convergence. According to Adulyasak, Cordeau, and Jans (2015), adding too many
467 cuts at each iteration can lead to a worse performance because of the time taken to solve
468 the master problem. Therefore, in this research, we first separate first- and second-stage
469 models and then aggregate all the scenarios in the SAA for the second-stage model,
470 resulting in a master problem MP and a subproblem $SP(x, y)$ that depend on the
471 first-stage solutions as follows:

$$472 [MP] \text{ Min } \sum_{j \in J} f_j x_j + \sum_{j \in J} h_j y_j + \theta \quad (17)$$

473 subject to (3), (7), (8)

$$474 \theta \geq Q(x, y) \quad (18)$$

$$475 [SP(x, y)] \quad Q(x, y) = \min_{s \in S'} \frac{1}{|S'|} (\sum_{i \in I^s} \sum_{j \in J} c_{ij} T_{ij}^s z_{ij}^s + \sum_{i \in I^s} \sum_{j \in J \setminus N_i^s} \gamma (T_{ij}^s -$$

476 $R) z_{ij}^s + c^{lost} (\sum_{i \in I^s} d_i^s - \sum_{i \in I^s} \sum_{j \in J} z_{ij}^s)) \quad (19)$

477 subject to

$$478 \quad \sum_{j \in J} z_{ij}^s \leq d_i^s, \quad \forall i \in I^s, s \in S' \quad (20)$$

$$479 \quad z_{ij}^s \leq y_j - \sum_{i'=1}^{i-1} \delta_{i'ij}^s z_{i'j}^s, \quad \forall i \in I^s, j \in J, s \in S' \quad (21)$$

$$480 \quad \sum_{i \in I^s} z_{ij}^s \leq \mu_j y_j, \quad \forall j \in J, s \in S' \quad (22)$$

$$481 \quad z_{ij}^s \in \mathbb{Z}_0^+, \quad \forall i \in I^s, j \in J, s \in S' \quad (23)$$

482 where S' represents the sample used in SAA.

483 Both master problem and subproblem are integer optimization problems. Repeatedly
484 solving these two types of problems is time-consuming. What is more, the coefficient
485 matrix in constraints is not necessarily totally unimodular, making the problem more
486 difficult to solve. Thus, we do not apply the classical Bender decomposition method
487 (Benders, 1962) to solve this problem. As an alternative, we propose a two-phase
488 Benders Decomposition based method. In phase 1, we relax all the integer constraints
489 in MP and $SP(x, y)$, obtaining relaxed version called RMP and $RSSP(x, y)$
490 respectively, where θ in Eq. (35) is an underestimator of $Q_r(x, y)$ (the subscript “r”
491 means “relaxed”) and successively approximates the shape of $Q_r(x, y)$ by adding cuts
492 derived from dual problem $DRSP(x, y)$. We iterate between RMP and $RSSP(x, y)$,
493 adding Benders cuts (35) into RMP until no Benders cut is added, the gap between
494 upper bound UB and lower bound LB is below a predetermined tolerance ϵ , or a
495 preset time is reached.

$$496 \quad [RMP] \text{ Min } \sum_{j \in J} f_j x_j + \sum_{j \in J} h_j y_j + \theta \quad (24)$$

497 subject to (3)

$$498 \quad \theta \geq Q_r(x, y) \quad (25)$$

$$499 \quad 0 \leq x_j \leq 1, \quad \forall j \in J \quad (26)$$

$$500 \quad y_j \geq 0, \quad \forall j \in J \quad (27)$$

$$501 \quad [RSP(x, y)] \quad Q_r(x, y) = \min_{s \in S'} \frac{1}{|S'|} (\sum_{i \in I^s} \sum_{j \in J} c_{ij} T_{ij}^s z_{ij}^s + \sum_{i \in I^s} \sum_{j \in J \setminus N_i^s} \gamma (T_{ij}^s -$$

$$502 \quad R) z_{ij}^s + c^{lost} (\sum_{i \in I^s} d_i^s - \sum_{i \in I^s} \sum_{j \in J} z_{ij}^s)) \quad (28)$$

503 subject to (20)–(22)

$$504 \quad z_{ij}^s \geq 0, \quad \forall i \in I^s, j \in J, s \in S'. \quad (29)$$

$$505 \quad [DRSP(x, y)] \quad D_r(x, y) = \max_{s \in S'} (\sum_{i \in I^s} \lambda_{is}^1 d_i^s + \sum_{i \in I^s} \sum_{j \in J} \lambda_{ijs}^2 y_j + \sum_{j \in J} \lambda_{js}^3 \mu_j y_j +$$

$$506 \quad \frac{c^{lost}}{|S'|} \sum_{i \in I^s} d_i^s) \quad (30)$$

507 subject to

$$508 \quad \lambda_{is}^1 + \lambda_{ijs}^2 + \sum_{i'=i+1}^{|I^s|} \delta_{i'i}^s \lambda_{ijs}^2 + \lambda_{js}^3 \leq \frac{1}{|S'|} (c_{ij} T_{ij}^s + \mathbb{I}_{j \in J \setminus N_i^s} \gamma (T_{ij}^s - R) - c^{lost}), \quad \forall i \in$$

$$509 \quad I^s, j \in J, s \in S' \quad (31)$$

$$510 \quad \lambda_{is}^1 \leq 0, \quad \forall i \in I^s, s \in S' \quad (32)$$

$$511 \quad \lambda_{ijs}^2 \leq 0, \quad \forall i \in I^s, j \in J, s \in S' \quad (33)$$

$$512 \quad \lambda_{js}^3 \leq 0, \quad \forall j \in J, s \in S' \quad (34)$$

513 where $\mathbb{I}_{j \in J \setminus N_i^s}$ is an indicator function. If $j \in J \setminus N_i^s$, the value is 1, otherwise, it is 0.

$$514 \quad \theta \geq \sum_{s \in S'} (\sum_{i \in I^s} \bar{\lambda}_{is}^1 d_i^s + \sum_{i \in I^s} \sum_{j \in J} \bar{\lambda}_{ijs}^2 y_j + \sum_{j \in J} \bar{\lambda}_{js}^3 \mu_j y_j + \frac{c^{lost}}{|S'|} \sum_{i \in I^s} d_i^s) \quad (35)$$

515 where $(\bar{\lambda}_{is}^1, \bar{\lambda}_{ijs}^2, \bar{\lambda}_{js}^3)$ are optimal solutions of $DRSP(x, y)$ given optimal master
516 problem solution (\bar{x}, \bar{y}) . These expressions are added to RMP to tighten the lower
517 bound.

518 In phase 2, we restore the integrality constraints in master problem and keep all the
519 Benders cuts generated in phase 1. We input the optimal solution of MP to
520 $RSSP(x, y)$ to evaluate whether additional cuts are needed. If no Benders cut is added
521 or a preset time is reached, solve $SP(x, y)$ under given optimal solution of MP . The
522 whole algorithm is presented as follows:

523

Algorithm 2 Benders Decomposition based solution approach (TPBD)

19. **Phase 1 procedure:**

20. **Input** A tolerance $\epsilon_1 \geq 0$ and maximum phase 1 run time TL_1 .

21. **Initialize** $UB_1 = \infty, LB_1 = 0$.

22. **for** $i = 1, 2, \dots$ **do**

23. Solve RMP to obtain optimal solution (\bar{x}_r, \bar{y}_r) and the optimal objective value $lobj$.

24. Update $LB := \max(LB, lobj)$.

25. Solve $RSP(\bar{x}_r, \bar{y}_r)$ to obtain $Q_r(\bar{x}_r, \bar{y}_r)$.

26. **if** $\theta < Q_r(\bar{x}_r, \bar{y}_r)$ **then**

27. Add optimality cut (35) to RMP .

28. **end if**

29. Calculate $uobj = \sum_{j \in J} f_j \bar{x}_{rj} + \sum_{j \in J} h_j \bar{y}_{rj} + Q_r(\bar{x}_r, \bar{y}_r)$.

30. Update $UB := \min(UB, uobj)$

31. **if** No Benders cut is added or $UB - LB \leq \epsilon_1$ or $runtime \geq TL_1$ **then**

32. Break.

33. **end if**

34. **end for**

35. **Phase 2 procedure:**

36. **Input** Maximum phase 2 run time TL_2 .

37. **Initialize** $UB_2 = \infty$, $LB_2 = LB_1$.

38. Keeping all Benders cuts generated in phase 1 procedure.

39. **for** $i = 1, 2, \dots$ **do**

40. Solve MP to obtain optimal solution (\bar{x}, \bar{y}) and the optimal objective value $lobj^*$.

41. Update $LB := \max(LB, lobj^*)$.

42. Solve $RSP(\bar{x}, \bar{y})$ to obtain $Q_r(\bar{x}, \bar{y})$.

43. **if** $\theta < Q_r(\bar{x}_r, \bar{y}_r)$ **then**

44. Add optimality cut (35) to MP .

45. **end if**

46. **if** no Benders cut is added or $runtime \geq TL_2$ **then**

47. Solve $SP(\bar{x}, \bar{y})$ to obtain $Q(\bar{x}, \bar{y})$.

48. **end if**

49. Calculate $uobj^* = \sum_{j \in J} f_j \bar{x}_j + \sum_{j \in J} h_j \bar{y}_j + Q(\bar{x}, \bar{y})$.

50. Update $UB := \min(UB, uobj^*)$.

51. Break.

52. **end for**

53. **Return** UB and corresponding optimal solution (\bar{x}, \bar{y}) .

524

525 **5. Numerical Experiments**

526 To evaluate the performance of our method, we conducted numerical experiments
527 using real-world emergency data. We first determined how many scenarios are needed
528 to obtain a high level of approximation precision while at the same time make the model
529 computationally tractable in SAA. Then we compared the computational performances
530 of SAA with TPBD under the same sample. To show the robustness of solution method,
531 we conducted out-of-sample analysis. Next, we evaluated the benefit of using stochastic
532 programming approach over deterministic model. Finally, we conducted sensitive
533 analysis to show how the value of some crucial parameters will influence the optimal
534 objective value of our model, which yields some valuable managerial insights. All the
535 experiments were carried out on a computer with i9-12900K CPU, 3.20 GHz processing

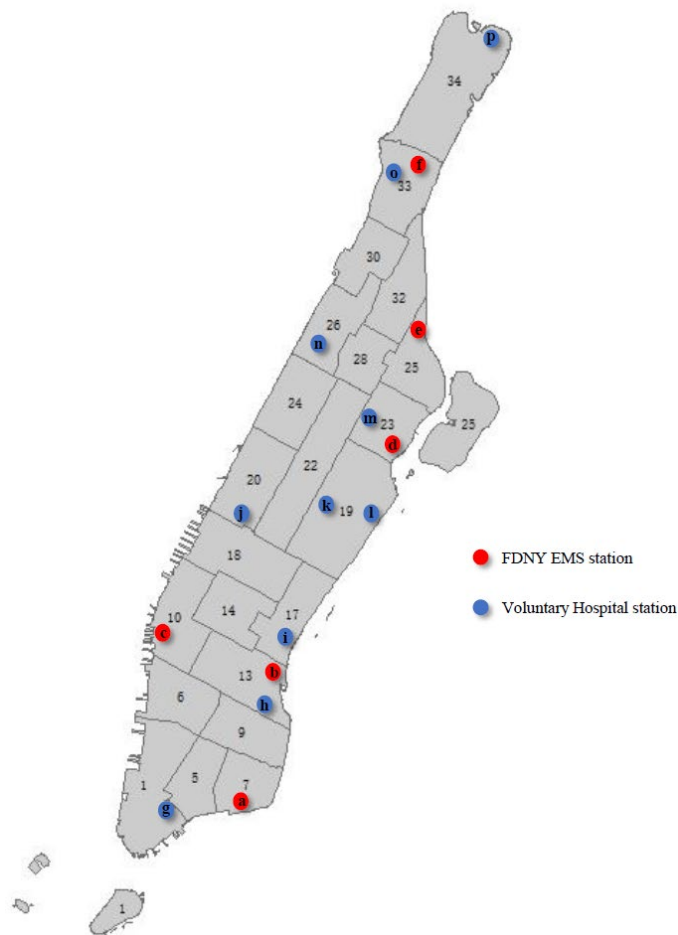
536 speed and 32 GB of memory. The model and the algorithm were implemented in C++
537 programming and solved by CPLEX 12.10.

538

539 **5.1 Parameter Setting**

540 We used the emergency incident data of Manhattan, which is provided by Fire
541 Department of New York City (FDNY). The data spans from the time the incident is
542 created to the time the incident is closed in the system, including the incident datetime,
543 incident location, response time, response police precinct, etc. We finally chose the data
544 of year 2011, which includes 225634 emergency call logs occurring at 22 police
545 precincts, as the data of the other years has a large number of missing values. The 22
546 police precincts were regarded as the demand areas. The candidate facility site was
547 obtained from the website of New York City Fire Department Bureau of Emergency
548 Medical Services (FDNY EMS), which is divided into four sectors, but nearly all the
549 incidents are served by two sectors: FDNY EMS municipal (FEM) and Voluntary
550 Hospital EMS (VHE). The former controls 70% of the ambulances in the New York
551 City 911 System and serves 63% of ambulance tours while the latter controls the rest
552 of the ambulances and serves 37% of the ambulance tours. FEM now operates 6 stations
553 in Manhattan and VHE provides emergency services through 10 stations. Totally, 16
554 stations were used as the candidate facility sites in this section. The demand area and
555 candidate stations are shown in Figure 2, where the numbers and letters are the indices
556 of the police precinct and candidate station, respectively. The red and blue circles
557 represent the location of stations operated by FEM and VHE, respectively. In order to
558 reflect the fact that the 6 stations operated by FEM work as the main emergency
559 facilities, the fixed cost was set lower and facility and service capacity were higher than
560 those of stations operated by VHE. The vehicle travel speed data was calculated using
561 the 2011 Yellow Taxi Trip Data, which includes trip distance and trip duration. We
562 divided each day into 24 equal segments, i.e., each segment is one hour, and we
563 calculated the average speed of each segment as the speed at which vehicles responded
564 to emergency calls occurring during that segment. The response time standard was set
565 to 9 minutes according to National Fire Protection Association benchmark. The fixed
566 cost to set up a station, including land cost, construction cost, material cost, etc., is
567 amortized according to service life of buildings to \$1500 and \$4500 per day for stations
568 belong to FEM and VHE, respectively. The reason for cost difference between FEM

569 and VHE stations is that FEM stations are more convenient and flexible for vehicle
570 deployment and dispatching. When an emergency occurs, FEM stations are preferred
571 for service. Therefore, facility capacity and vehicle service capacity for FEM stations
572 are also greater than that of VHE stations. Vehicle purchase cost, including ambulance
573 cost, equipment cost, maintenance cost, etc., is amortized to \$300 per vehicle per day
574 according to service life of ambulance. Unit transportation cost, including labor cost,
575 fuel cost, etc., is \$30 per minute. The penalty cost is set to 5 times the transportation
576 cost because survivability is highly dependent on response time and the response that
577 is later than response time standard may have serious consequences to life and property.
578 The values of parameters used in numerical experiments are listed in Table 4.
579



580
581 Figure 2. The demand area and candidate stations
582

583

Table 4. The values of parameters used in numerical experiments

Parameter	Value
f_j	\$1500 per day if j belongs to FEM \$4500 per day if j belongs to VHE
h	\$300 per vehicle per day
c_{ij}	\$30 per minute
γ	\$150 per minute
Q_j	10 vehicles if j belongs to FEM 5 vehicles if j belongs to VHE
α	90%
μ_j	20 demands per day if j belongs to FEM 5 demands per day if j belongs to VHE

584

585 **5.2 Determination of Sample Size**

586 In this section, we will determine the sample size for the following numerical
587 experiments. We first introduce how a sample is generated. Then, we illustrate the
588 procedures and results of solution quality evaluation. Accordingly, we finally determine
589 the sample size.

590

591 **5.2.1 Sample Generation**

592 A sample is made up of a specific number of scenarios. A scenario is real emergency
593 rescue data for a day, including demand and travel speed data. The specific number of
594 scenarios is called sample size. We have emergency call logs and yellow taxi trip data
595 of year 2011, which spans over 365 days. We randomly generated a number between 1
596 and 365 and extracted all emergency call logs and yellow taxi trip data for the day
597 corresponding to this number as scenario data. We repeated this procedure in the
598 remaining data set until the number of scenarios equals the sample size.

599

600 **5.2.2 Solution Quality Evaluation**

601 We used Algorithm 1 to evaluate the performance of different sample sizes, which
602 were set to 1, 3, 5, 8, 10, 15, 20, 30, 50, 80, and 100. We first ran line 1 and 2 of
603 Algorithm 1 to obtain the optimal first-stage solution. Then we iterated line 4 to 8 for
604 10 times where $|S'_m|$ was set to 150. Finally, we calculated solution quality, which is

605 shown in Table 5. NS, PE, CI, CI-L, CI-U, and Time represent sample size, point
606 estimate, 95% confidence interval, ratio between 95% confidence interval for
607 optimality gap and point estimate of lower bound, ratio between 95% confidence
608 interval for optimality gap and point estimate of upper bound, and running time,
609 respectively. When sample size was set to 80 and 100, the computer ran out of memory
610 due to the large number of variables and constraints. Therefore, we do not show the
611 results of sample size 80 and 100. It can be seen from the table that when sample size
612 is 20, the error rate is below 1% with a high probability. Considering the trade-off
613 between solution quality and computational tractability, we set the sample size to 20 in
614 the following calculation.

615
616

Table 5. Solution quality of different sample sizes (95% CIs)

NS	Lower bound		Upper bound		Optimality gap		CI-L	CI-U	Time
	PE	CI	PE	CI	PE	CI			
1	188700	(182300, 195100)	191700	(177200, 206200)	3063	(0, 18340)	9.72%	9.57%	164
3	178900	(173500, 184300)	181000	(177000, 185000)	2141	(0, 9291)	5.19%	5.13%	674
5	186400	(184300, 188500)	189300	(184000, 194600)	2904	(0, 8149)	4.37%	4.30%	1380
8	183400	(181600, 185200)	185900	(182000, 189800)	2466	(0, 5760)	3.14%	3.10%	2679
10	180800	(175600, 186000)	182000	(178100, 185900)	1194	(0, 4330)	2.39%	2.38%	2648
15	180400	(175900, 184900)	181600	(176200, 187000)	1252	(0, 2197)	1.22%	1.21%	3835
20	180300	(176400, 184200)	181500	(178900, 184100)	1220	(0, 1404)	0.78%	0.77%	6587
30	183700	(180700, 186700)	184700	(180300, 189100)	952	(0, 1072)	0.58%	0.58%	11960
50	179400	(176500, 182300)	180200	(177700, 182700)	786	(0, 803)	0.45%	0.45%	35360

617

618 5.3 Numerical Performance of Our TPBD Algorithm

619 Since each scenario contains around 800 demands, with the increase of sample size,
620 SAA will become computationally inefficient or even run out of memory. In this section,
621 we compared the computational performances of SAA and TPBD under different
622 samples, which are shown in Table 6. The first column represents the sample size. *Total*
623 *Cost* is the objective value. *Time* is the computation time. It can be seen that TPBD
624 obtains the same results as SAA. When the sample size is small, SAA gets the optimal
625 solution in shorter time. As the sample size increases, TPBD becomes more efficient.
626 If we set the time limit to one hour, TPBD can get the optimal solutions under all the

627 sample sizes. However, SAA can only finish the computation for sample sizes less than
 628 50. The results suggest that TPBD is an efficient algorithm, especially for large-size
 629 problem. Therefore, in the following numerical experiments, we will use TPBD to solve
 630 the problem unless otherwise specified.

631 Table 6. Comparison of computational performances between SAA and TPBD

NS	Total Cost		Time	
	SAA	TPBD	SAA	TPBD
1	157800	157800	1.21	6.27
3	180100	180100	9.19	23.41
5	179700	179700	27.75	64.05
8	188600	188600	40.05	111.80
10	172500	172500	51.41	154.50
15	170400	170400	127.80	257.10
20	183000	183000	842.40	550.80
30	183800	183800	1037.00	620.50
50	186900	186900	4942.00	1355.00
80	178900	178900	9207.00	2907.00
100	180200	180200	11650.00	3588.00

632

633 5.4 Robustness Evaluation

634 Optimal solutions obtained through TPBD are based on generated scenarios. It is
 635 highly possible that in practice the realized demand is not a member of used samples,
 636 resulting in poor performance of solution approach. To test the robustness of TPBD
 637 method, we conduct out-of-sample analysis, which evaluates performance of optimal
 638 solutions using out-of-sample data. The procedures are shown in Algorithm 3. We
 639 divide all scenarios into two sets: in-sample scenario S^{in} and out-of-sample scenario
 640 S^{out} , where $S^{in} \cup S^{out} = S$. All scenarios in S^{in} are input into TPBD to obtain
 641 optimal first-stage solution x^* and y^* . Then N scenarios are randomly selected from
 642 S^{out} , each of which is input into deterministic model [M1] to calculate optimal
 643 dispatching decisions. There are two possible outcomes: optimal solution exists and
 644 optimal solution does not exist. In the first case, calculate coverage level, i.e., the
 645 percentage of demand covered by opening stations within response time standard,
 646 response level, i.e., the proportion of demand that is responded within response time
 647 standard, and demand loss rate, i.e., the percentage of demand that is unserved.
 648 Robustness is measured by robustness level, the proportion of scenarios where optimal

649 solution can be found, $(1 - \tau)$ -CI of coverage level, response level, and demand loss
650 rate.
651

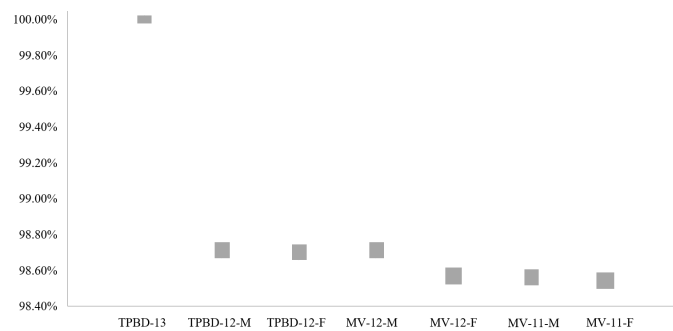
Algorithm 3 Robustness evaluation

1. Generate a set of scenarios S^{in} .
 2. Solve the TPBD problem with S^{in} and obtain the optimal first-stage solution x^* and y^* .
 3. **for** $n = 1, 2, \dots, N$ **do**
 4. Generate a new independent scenario s_n from S^{out} , $s_n \in S^{out}$.
 5. Solve deterministic model [M1] with s_n , x^* and y^* .
 6. If optimal solution can be found, calculate coverage level CL_n , response level RL_n , and demand loss rate DL_n .
 7. **end for**
 8. Calculate the number of iterations N_o where optimal solution can be found.
 9. Calculate robustness level $BL = \frac{N_o}{N}$.
 10. Let $ACL = \frac{\sum_{n=1}^{N_o} CL_n}{N_o}$ and $S_{ACL} = \frac{1}{N_o-1} \sum_{n=1}^{N_o} (CL_n - ACL)^2$.
 11. Let $ARL = \frac{\sum_{n=1}^{N_o} RL_n}{N_o}$ and $S_{ARL} = \frac{1}{N_o-1} \sum_{n=1}^{N_o} (RL_n - ARL)^2$.
 12. Let $ADL = \frac{\sum_{n=1}^{N_o} DL_n}{N_o}$ and $S_{ADL} = \frac{1}{N_o-1} \sum_{n=1}^{N_o} (DL_n - ADL)^2$.
 13. The $(1 - \tau)$ -CI of coverage level is $\left[ACL - \frac{t_{N_o-1, \frac{\tau}{2}} \sqrt{S_{ACL}}}{\sqrt{N_o}}, ACL + \frac{t_{N_o-1, \frac{\tau}{2}} \sqrt{S_{ACL}}}{\sqrt{N_o}} \right]$.
 14. The $(1 - \tau)$ -CI of response level is $\left[ARL - \frac{t_{N_o-1, \frac{\tau}{2}} \sqrt{S_{ARL}}}{\sqrt{N_o}}, ARL + \frac{t_{N_o-1, \frac{\tau}{2}} \sqrt{S_{ARL}}}{\sqrt{N_o}} \right]$.
 15. The $(1 - \tau)$ -CI of demand loss rate is $\left[ADL - \frac{t_{N_o-1, \frac{\tau}{2}} \sqrt{S_{ADL}}}{\sqrt{N_o}}, ADL + \frac{t_{N_o-1, \frac{\tau}{2}} \sqrt{S_{ADL}}}{\sqrt{N_o}} \right]$
-

652
653 We ran TPBD and deterministic model [M1] 10 times separately. There are two types
654 of solutions for TPBD according to the number of opening stations: one with 12 opening
655 stations and the other with 13 opening stations. For the solution with 13 opening stations,
656 the number of deployed vehicles is the same. We denoted this type of solution by
657 TPBD-13. For the solution with 12 opening stations, the number of deployed vehicles
658 varies. Therefore, we selected two solutions with the fewest and most vehicles, denoting
659 them by TPBD-12-F and TPBD-12-M, respectively. For deterministic model, there are
660 two types of solutions according to the number of opening stations: one with 11 opening

661 stations and the other with 12 opening stations. The number of deployed vehicles under
 662 both solution types varies. We took the same measure as TPBD with 12 opening stations
 663 and finally we got four solution types, namely MV-12-F, MV-12-M, MV-11-F, and MV-
 664 11-M. Totally, we got seven different optimal first-stage solutions, which were put into
 665 Algorithm 3 to calculate robustness level and 95% CI of coverage level, response level,
 666 and demand loss rate with N set to 150. Figure 3 and Figure 4 show 95% CI of
 667 coverage level and response level, respectively. Results show that compared with
 668 deterministic model, stochastic model can obtain solutions that cover and respond to
 669 more emergency calls within response time standard with high probability. Even though
 670 these solutions are tested in out-of-sample data, at least 98.68% and 95.55% of the
 671 emergency demands can be covered and responded to in time with high probability.
 672 While solutions of deterministic model can provide a maximum 98.69% and 94.52%
 673 of coverage and responses respectively within response time standard with high
 674 probability. Most of the time, the response level is below 90%. All solution types do
 675 not have demand unsatisfied, except MV-12-F and MV-11-F having demand loss rate
 676 0.2% and 0.6% respectively with a high probability. Robustness level is 100% for all
 677 solution types, indicating that solutions obtained by TPBD can find optimal value under
 678 every scenario in out-of-sample analysis. All these results suggest that TPBD is a robust
 679 solution approach.

680



681

682

Figure 3. 95% CI of coverage level

683

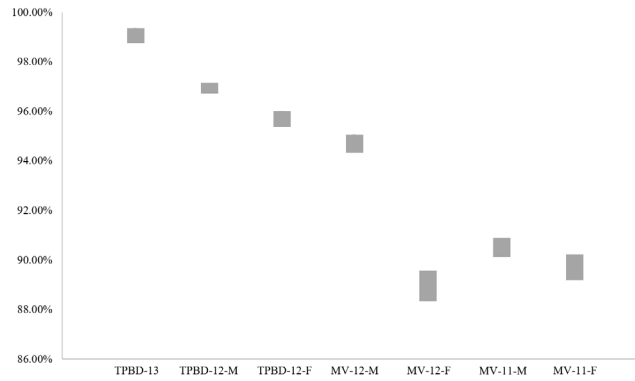


Figure 4. 95% CI of response level

684

685

686

687 5.5 The Benefit of Stochastic Programming

688 We evaluate the benefit of stochastic programming by comparing the performance of
689 the mean value solution with the stochastic solution under the same sample. The mean
690 value solution was obtained by taking the mean value of the sample into the
691 deterministic model [M1]. For a sample with 20 scenarios, we first ignored the date of
692 the emergency calls and summed up all the emergency calls that happened at the same
693 demand area and during the same time segment. There were 24 time segments, each
694 representing one hour of the 24 hours. Therefore, we obtained demand distribution over
695 22 demand areas \times 24 time segments. Then we divided the distribution by 20, i.e., the
696 number of scenarios, and obtained the average number of emergency calls occurred at
697 each demand area during each time segment, which is the mean value of a sample. We
698 put the average data into the deterministic model [M1] and obtained the first-stage
699 solution. Then we solved the second-stage problem by fixing the first-stage solution in
700 TPBD under the same sample, the result of which was compared with that of the TPBD
701 using the same scenarios. The results are shown in the Table 7. The first and second
702 row are performance of mean value solution and stochastic solution, respectively.
703 Improvement indicates the percentage cost saving and coverage improvement of
704 stochastic model. The second column is the total cost. FC, PC, TRC, and Penalty
705 represent fixed cost of opening stations, purchasing cost of ambulances, transportation
706 cost of demand fulfillment, and penalty for overtime, respectively. CL is the coverage
707 level, i.e., the percentage of demand covered by opening stations within response time
708 standard. RL is the response level, i.e., the percentage of demand being served within
709 response time standard. DL is the demand loss rate, i.e., the percentage of demand that
710 is unserved. The last column is the computation time.

711 The results show that stochastic programming could reduce the total cost by 7.93%,
 712 which is achieved by opening more facilities and purchasing more vehicles to reduce
 713 the demand fulfillment cost and delay penalty. As there are more facilities and vehicles,
 714 the coverage level and response level are improved by 4.67% and 5.25%, respectively.
 715 Therefore, stochastic programming could achieve a better coverage with less cost
 716 compared with deterministic one.

717

718 Table 7. The performance of the mean value solution and the stochastic solution under
 719 the same sample

	TC	FC	PC	TRC	Penalty	CL	RL	DL	Time
MV	192900	31500	19500	120400	21500	94.30%	88.98%	0.00%	15.26
TPBD	177600	36000	24000	110800	6800	98.70%	93.65%	0.00%	280
Improvement	7.93%	-14.29%	-23.08%	7.97%	68.37%	-4.67%	-5.25%	0.00%	-1736.83%

720

721

722 5.6 Sensitive analysis

723 We conducted the sensitive analysis to evaluate the influence of the value of crucial
 724 parameters on the optimal value of the stochastic programming. The following
 725 parameters are considered: response time standard, facility capacity, service capacity,
 726 and facility heterogeneity.

727

728 5.6.1 Impact of response time standard

729 Table 8 reports the effect of response time standard on cost, the number of facilities
 730 opened, the number of vehicles purchased, coverage level, and response level. RT
 731 means the value of response time standard. We can see that with the release of response
 732 time standard, total cost will reduce to certain value and then stay unchanged. Because
 733 when response time standard is relaxed, more demands can be served within time
 734 requirement, greatly reducing penalty. At the same time, with the enlarge of coverage
 735 radius, stations and vehicles can cover further demands without bearing penalty,
 736 reducing the number of stations and vehicles needed and thus reducing set-up and
 737 purchasing cost. When response time standard is raised to certain level, demands are
 738 fully covered and can be responded without any penalty under the required service level.
 739 The best combination of costs components is achieved so that further increasing

740 response time standard will not influence the system performance.

741 The results suggest that increasing response time standard under threshold can reduce
742 total cost but does not have any effect when response time standard is already beyond
743 the threshold. Besides, response time standard highly determines emergency
744 survivability. When authorities set up response time standard, it is important to make a
745 balance between cost and survivability.

746

747

Table 8. The influence of response time standard

RT	Total Cost	Fixed Cost	Purchase Cost	Transportation Cost	Penalty	Coverage Level	Response Level
1	644500	45000	27600	115500	456400	1.85%	1.71%
3	424300	45000	27600	115500	236200	34.48%	24.38%
5	281100	45000	27600	115600	92900	73.24%	57.34%
7	214600	45000	27000	116300	26300	90.43%	81.96%
9	188400	36000	24900	120700	6800	98.79%	94.29%
11	178400	27000	22500	127100	1800	100.00%	98.12%
13	174000	18000	20100	135900	0	100.00%	100.00%
15	174000	18000	20100	135900	0	100.00%	100.00%
17	174000	18000	20100	135900	0	100.00%	100.00%

748

749 5.6.2 Impact of facility capacity and vehicle service capacity

750 The impact of facility capacity is reflected by Table 9. The first and second column
751 are the capacity of stations operated by FEM and VHE, respectively. We assume that
752 fixed cost of station does not change with capacity. We can see that increasing the
753 capacity of stations operated by FEM does not change the model performance while
754 increasing the capacity of stations operated by VHE first decreases total cost and when
755 the capacity increases to 20 vehicles, the optimal value does not change anymore.
756 Because current station capacity of FEM is enough to achieve the required service level,
757 making increased capacity redundant. However, increasing station capacity of VHE,
758 which is originally set to a low value, allows more demands to be served by closer
759 vehicles, and therefore reducing the number of stations, transportation cost, and penalty.
760 It is worth noting that even though the number of stations is reduced, the coverage level
761 stays the same and what's more, response level is improved, which means that
762 restricting station capacity of VHE will waste money building redundant stations and

763 forcing demands to be served by further vehicles.

764 Results in Table 9 suggest that there is a threshold for station capacity where when
 765 the value is below the threshold, increasing capacity can reduce total cost because more
 766 vehicles can be used to serve demands, thus reducing station set-up cost, transportation
 767 cost, and penalty, while when the value is above the threshold, increasing capacity does
 768 not influence the performance of the system because existing vehicles can already
 769 achieve the required service level at the minimal cost and additional capacity dose not
 770 contribute to better service, thus being redundant.

771

772 Table 9. The impact of facility capacity on optimal value

CM	CV	TC	FC	NM	NV	PC	VM	VV	TRC	Penalty	CL	RL
10	5	182400	36000	6	6	24600	52	30	115800	6000	98.53%	94.13%
15	5	182400	36000	6	6	24600	52	30	115800	6000	98.53%	94.13%
20	5	182400	36000	6	6	24600	52	30	115800	6000	98.53%	94.13%
25	5	182400	36000	6	6	24600	52	30	115800	6000	98.53%	94.13%
30	5	182400	36000	6	6	24600	52	30	115800	6000	98.53%	94.13%
10	5	182400	36000	6	6	24600	52	30	115800	6000	98.53%	94.13%
10	10	165800	27000	6	4	25800	47	39	109600	3400	98.53%	96.31%
10	15	160500	27000	6	4	28500	43	52	104100	900	98.53%	98.11%
10	20	160400	27000	6	4	28800	43	53	103900	700	98.53%	98.16%
10	25	160400	27000	6	4	28800	43	53	103900	700	98.53%	98.16%
10	30	160400	27000	6	4	28800	43	53	103900	700	98.53%	98.16%

773

774 Table 10 shows impact of service capacity, where the first and second column are
 775 the vehicle service capacity of stations operated by FEM and VHE, respectively. Table
 776 10 has the same results as Table 9, because increasing vehicle service capacity under
 777 fixed facility capacity has a similar effect to increasing facility capacity under fixed
 778 vehicle service capacity, both of which increase demands served by each station.

779 Results of Table 9 and 10 inform us that there are thresholds for facility and vehicle
 780 service capacity. Setting both capacities to their thresholds could achieve maximal
 781 coverage with minimal cost.

782 Table 10. The impact of service capacity on optimal value

SM	SV	TC	FC	NM	NV	PC	VM	VV	TRC	Penalty	CL	RL
----	----	----	----	----	----	----	----	----	-----	---------	----	----

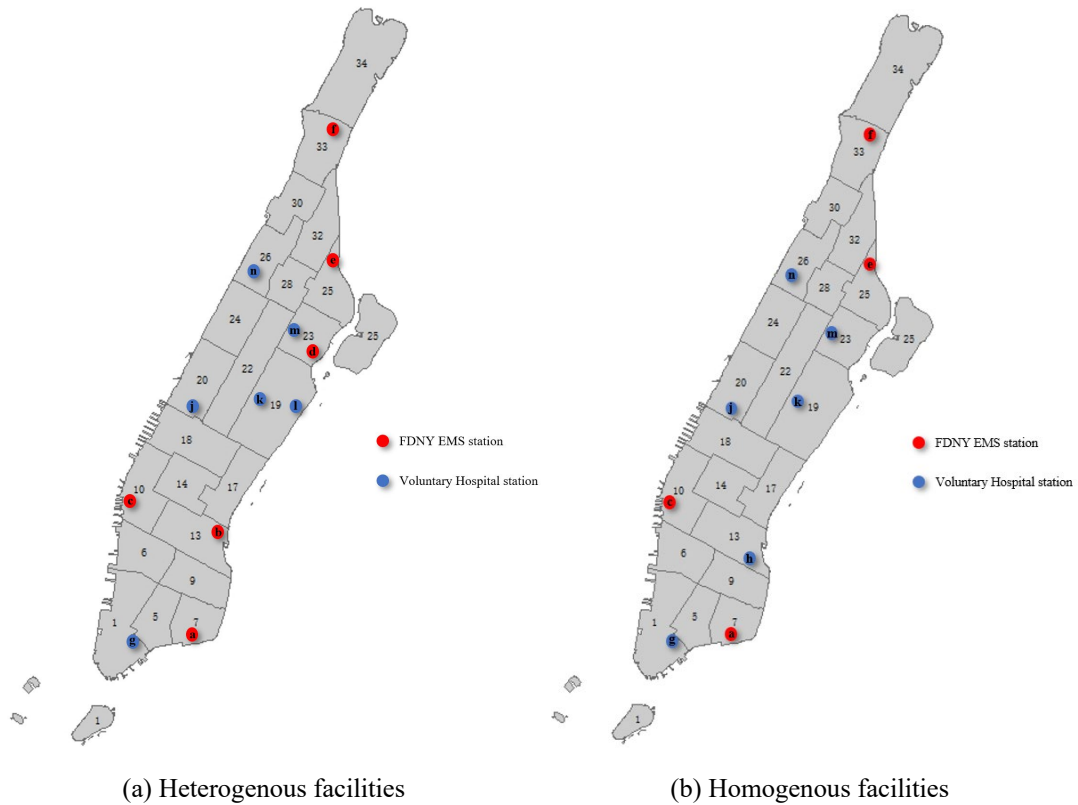
20	5	180900	36000	6	6	24300	51	30	114400	6200	98.23%	93.88%
20	10	160800	31500	6	5	21300	46	25	105400	2600	98.23%	96.74%
20	15	153200	31500	6	5	19800	41	25	98400	3500	98.23%	97.23%
20	20	153200	31500	6	5	19800	41	25	98400	3500	98.23%	97.23%
20	25	153200	31500	6	5	19800	41	25	98400	3500	98.23%	97.23%
20	5	180900	36000	6	6	24300	51	30	114400	6200	98.23%	93.88%
25	5	180900	36000	6	6	24300	51	30	114400	6200	98.23%	93.88%
30	5	180900	36000	6	6	24300	51	30	114400	6200	98.23%	93.88%
35	5	180900	36000	6	6	24300	51	30	114400	6200	98.23%	93.88%
40	5	180900	36000	6	6	24300	51	30	114400	6200	98.23%	93.88%

783

784 5.6.3 Impact of facility heterogeneity

785 To reflect the fact that stations of FEM control 70% of the ambulances in the New
786 York City 911 System and serve 63% of ambulance tours while stations of VHE control
787 the rest of the ambulances and serve 37% of the ambulance tours, we differentiated
788 these two types of stations in terms of cost and capacity. In this section, we will compare
789 results of problem setting with facility heterogeneity to those of problem setting where
790 all facilities are assumed homogenous in term of cost and capacity.

791 We ran 10 times of each problem setting. The optimal location solutions are the same
792 for 10 iterations, which are shown in Figure 5. Figure 5(a) and 5(b) show optimal
793 location solutions for heterogenous and homogenous facilities, respectively. We can see
794 that when heterogeneity is considered, all the 6 stations operated by FEM, i.e., the red
795 circles, are selected for the following reasons: first, these stations are cheaper than the
796 rest of the stations; second, they can serve more demands; third, they are sparsely
797 distributed throughout the whole region. When the 6 stations cannot meet the
798 requirements, stations in the place not covered by FEM or with high-density population,
799 i.e., stations *g, j, k, l, m, and n*, are selected to fill the vacancy. This suggests that with
800 heterogeneity, municipal stations are the best choice to serve emergency calls and when
801 demands exceed their service capacity, stations that can fill the vacancy are preferred.



804 Figure 5. Comparison of optimal location solutions for heterogenous and homogenous
 805 facilities

806

807 If we assume that all facilities are the same, results will change. The number of
 808 selected stations of FEM and VHE are 4 and 6, respectively. Stations *b*, *d*, and *l* are
 809 replaced by station *h*. The coverage level and response level are shown in Table 11,
 810 which indicate better performances of location solutions for homogenous facilities

811

812 Table 11. The coverage level and response level of optimal location solutions for
 813 heterogenous and homogenous facilities

814

	CL	RL
Homogenous	98.72%	98.13%
Heterogenous	98.70%	93.65%

815

816 Results of comparison indicate that even though stations of FEM are preferred by
 817 New York City 911 System for various reasons, increasing utilization of VHE stations
 818 will improve system performances.

819 6. Conclusion

820 In this research, we propose a dynamic scenario-based two-stage stochastic
821 programming model to optimize location of ambulance stations, deployment of
822 ambulances, and vehicle dispatching under uncertain demand and traffic situation. The
823 objective is to minimize total cost composed of station set-up cost, vehicle purchasing
824 cost, demand fulfillment cost, penalty for overtime under service level requirements,
825 and penalty for demand unsatisfaction. We apply Sample Average Approximation to
826 approximate the original problem and propose a two-phase Benders Decomposition
827 algorithm to accelerate computation. Numerical experiments are conducted to evaluate
828 the performance of the solution method. Results show that proposed solution method is
829 effective and efficient. Besides, it is proved to be robust in out-of-sample analysis. The
830 comparison between stochastic model and deterministic model reveals that stochastic
831 program could achieve a better coverage with less cost compared with deterministic
832 one. We also conduct sensitive analysis to evaluate the influence of the value of crucial
833 parameters on the optimal value of the stochastic programming model. Results suggest
834 that there exist thresholds for response time standard, facility capacity, and vehicle
835 service capacity. Increasing each of them under threshold can reduce total cost but does
836 not have any effect when the threshold is reached. Response time standard highly
837 determines emergency survivability. Authorities should make a balance between cost
838 and survivability when setting up response time standard. For facility capacity and
839 vehicle service capacity, it is recommended to set both to their thresholds in order to
840 achieve maximal coverage with minimal cost. Facility heterogeneity will also influence
841 problem solutions and performances. When considering heterogeneity, municipal
842 stations are the best choice to serve emergency calls and when demands exceed their
843 service capacity, stations that can fill the vacancy are preferred. When all facilities are
844 homogeneous, better performances can be achieved. Therefore, even though stations of
845 FEM are preferred by New York City 911 System for various reasons, increasing
846 utilization of VHE stations will improve system performances.

847 Future research can be extended in four directions. First, emergency demands can be
848 prioritized with the level of priority representing severity of emergency. We can require
849 that high priority calls be served first. Second, we can use multi-type vehicles to serve
850 demands, such basic life support ambulance and advanced life support ambulance. Each
851 type of vehicle can only serve specified demands. Third, we allow vehicles not to

852 respond immediately to a demand because a more severe emergency may happen in the
853 near future, resulting in no vehicle response. Fourth, as vehicles may not be able to
854 serve all demands, we can consider demand queuing.

855 **References**

- 856 Aringhieri, R., Bruni, M. E., Khodaparasti, S., & van Essen, J. T. (2017). Emergency
857 medical services and beyond: Addressing new challenges through a wide
858 literature review. *Computers & Operations Research*, 78, 349-368.
- 859 Adulyasak, Y., Cordeau, J. F., & Jans, R. (2015). Benders decomposition for production
860 routing under demand uncertainty. *Operations Research*, 63(4), 851-867.
- 861 Bandara, D., Mayorga, M. E., & McLay, L. A. (2014). Priority dispatching strategies
862 for EMS systems. *Journal of the Operational Research Society*, 65(4), 572-587.
- 863 Bélanger, V., Lanzarone, E., Nicoletta, V., Ruiz, A., & Soriano, P. (2020). A recursive
864 simulation-optimization framework for the ambulance location and dispatching
865 problem. *European Journal of Operational Research*, 286(2), 713-725.
- 866 Bélanger, V., Ruiz, A., & Soriano, P. (2019). Recent optimization models and trends in
867 location, relocation, and dispatching of emergency medical vehicles. *European*
868 *Journal of Operational Research*, 272(1), 1-23.
- 869 Benders, J. F. (1962). Partitioning procedures for solving mixed-variables programming
870 problems. *Numerische Mathematik*, 4(1), 238-252.
- 871 Beraldi, P., & Bruni, M. E. (2009). A probabilistic model applied to emergency service
872 vehicle location. *European Journal of Operational Research*, 196(1), 323-331.
- 873 Beraldi, P., Bruni, M. E., & Conforti, D. (2004). Designing robust emergency medical
874 service via stochastic programming. *European Journal of Operational Research*,
875 158(1), 183-193.
- 876 Berman, O., Hajizadeh, I., & Krass, D. (2013). The maximum covering problem with
877 travel time uncertainty. *IIE Transactions*, 45(1), 81-96.
- 878 Bertsimas, D., & Ng, Y. (2019). Robust and stochastic formulations for ambulance
879 deployment and dispatch. *European Journal of Operational Research*, 279(2),
880 557-571.
- 881 Boujemaa, R., Jebali, A., Hammami, S., Ruiz, A., & Bouchriha, H. (2018). A stochastic
882 approach for designing two-tiered emergency medical service systems. *Flexible*
883 *Services and Manufacturing Journal*, 30(1), 123-152.
- 884 Boutilier, J. J., & Chan, T. C. (2020). Ambulance emergency response optimization in
885 developing countries. *Operations Research*, 68(5), 1315-1334.
- 886 Bürger, A., et al. (2018). The effect of ambulance response time on survival following
887 out-of-hospital cardiac arrest: an analysis from the German resuscitation

888 registry. *Deutsches Ärzteblatt International*, 115(33-34), 541.

889 Chanta, S., Mayorga, M. E., Kurz, M. E., & McLay, L. A. (2011). The minimum p-envy
890 location problem: a new model for equitable distribution of emergency
891 resources. *IIE Transactions on Healthcare Systems Engineering*, 1(2), 101-115.

892 Chanta, S., Mayorga, M. E., & McLay, L. A. (2014). Improving emergency service in
893 rural areas: a bi-objective covering location model for EMS systems. *Annals of
894 Operations Research*, 221(1), 133-159.

895 Church, R., & ReVelle, C. (1974). The maximal covering location problem. Paper
896 Presented at the Papers of the Regional Science Association.

897 Daskin, M. S. (1983). A maximum expected covering location model: formulation,
898 properties and heuristic solution. *Transportation Science*, 17(1), 48-70.

899 Daskin, M. S., & Stern, E. H. (1981). A hierarchical objective set covering model for
900 emergency medical service vehicle deployment. *Transportation Science*, 15(2),
901 137-152.

902 Dean, S. F. (2008). Why the closest ambulance cannot be dispatched in an urban
903 emergency medical services system. *Prehospital and Disaster Medicine*, 23(2),
904 161-165.

905 Degel, D., Wiesche, L., Rachuba, S., & Werners, B. (2015). Time-dependent ambulance
906 allocation considering data-driven empirically required coverage. *Health Care
907 Management Science*, 18(4), 444-458.

908 El Itani, B., Abdelaziz, F. B., & Masri, H. (2019). A bi-objective covering location
909 problem: case of ambulance location in the Beirut area, Lebanon. *Management
910 Decision*, 57(2), 432-444.

911 Erkut, E., Ingolfsson, A., & Erdoğan, G. (2008). Ambulance location for maximum
912 survival. *Naval Research Logistics*, 55(1), 42-58.

913 Fire Department of New York City. EMS Incident Dispatch Data.
914 [https://data.cityofnewyork.us/Public-Safety/EMS-Incident-Dispatch-
915 Data/76xm-jjuj](https://data.cityofnewyork.us/Public-Safety/EMS-Incident-Dispatch-Data/76xm-jjuj).

916 Gendreau, M., Laporte, G., & Semet, F. (1997). Solving an ambulance location model
917 by tabu search. *Location Science*, 5(2), 75-88.

918 Goldberg, J., & Paz, L. (1991). Locating emergency vehicle bases when service time
919 depends on call location. *Transportation Science*, 25(4), 264-280.

920 Haghani, A., Tian, Q., & Hu, H. (2004). Simulation model for real-time emergency

921 vehicle dispatching and routing. *Transportation Research Record*, 1882(1),
922 176-183.

923 Ingolfsson, A., Budge, S., & Erkut, E. (2008). Optimal ambulance location with random
924 delays and travel times. *Health Care Management Science*, 11(3), 262-274.

925 Jenkins, P. R., Robbins, M. J., & Lunday, B. J. (2021). Approximate dynamic
926 programming for military medical evacuation dispatching policies. *INFORMS*
927 *Journal on Computing*, 33(1), 2-26.

928 Knight, V. A., Harper, P. R., & Smith, L. (2012). Ambulance allocation for maximal
929 survival with heterogeneous outcome measures. *Omega*, 40(6), 918-926.

930 Kuo, Y. H. (2014). Integrating simulation with simulated annealing for scheduling
931 physicians in an understaffed emergency department. *HKIE Transactions*, 21(4),
932 253-261.

933 Kuo, Y. H., Rado, O., Lupia, B., Leung, J. M., & Graham, C. A. (2016). Improving the
934 efficiency of a hospital emergency department: a simulation study with
935 indirectly imputed service-time distributions. *Flexible Services and*
936 *Manufacturing Journal*, 28(1), 120-147.

937 Laporte, G., & Louveaux, F. V. (1993). The integer L-shaped method for stochastic
938 integer programs with complete recourse. *Operations Research Letters*, 13(3),
939 133-142.

940 Lee, S. (2011). The role of preparedness in ambulance dispatching. *Journal of the*
941 *Operational Research Society*, 62(10), 1888-1897.

942 Lee, S. (2012). The role of centrality in ambulance dispatching. *Decision Support*
943 *Systems*, 54(1), 282-291.

944 Lee, S. (2013). Centrality-based ambulance dispatching for demanding emergency
945 situations. *Journal of the Operational Research Society*, 64(4), 611-618.

946 Liu, K., Li, Q., & Zhang, Z. H. (2019). Distributionally robust optimization of an
947 emergency medical service station location and sizing problem with joint
948 chance constraints. *Transportation Research Part B: Methodological*, 119, 79-
949 101.

950 Liu, Y., Li, Z., Liu, J., & Patel, H. (2016). A double standard model for allocating limited
951 emergency medical service vehicle resources ensuring service reliability.
952 *Transportation Research Part C: Emerging Technologies*, 69, 120-133.

953 McLay, L. A., & Mayorga, M. E. (2013a). A model for optimally dispatching

954 ambulances to emergency calls with classification errors in patient priorities.
955 *IIE Transactions*, 45(1), 1-24.

956 McLay, L. A., & Mayorga, M. E. (2013b). A dispatching model for server-to-customer
957 systems that balances efficiency and equity. *Manufacturing & Service*
958 *Operations Management*, 15(2), 205-220.

959 Nasrollahzadeh, A. A., Khademi, A., & Mayorga, M. E. (2018). Real-time ambulance
960 dispatching and relocation. *Manufacturing & Service Operations Management*,
961 20(3), 467-480.

962 Nelas, J., & Dias, J. (2020). Optimal emergency vehicles location: An approach
963 considering the hierarchy and substitutability of resources. *European Journal of*
964 *Operational Research*, 287(2), 583-599.

965 New York City Fire Department Bureau of Emergency Medical Services.
966 <https://emspac.org/sectors/municipal/>.

967 Nickel, S., Reuter-Oppermann, M., & Saldanha-da-Gama, F. (2016). Ambulance
968 location under stochastic demand: A sampling approach. *Operations Research*
969 *for Health Care*, 8, 24-32.

970 Park, H., Waddell, D., & Haghani, A. (2019). Online optimization with look-ahead for
971 freeway emergency vehicle dispatching considering availability. *Transportation*
972 *Research Part C: Emerging Technologies*, 109, 95-116.

973 Peng, C., Delage, E., & Li, J. (2020). Probabilistic envelope constrained multiperiod
974 stochastic emergency medical services location model and decomposition
975 scheme. *Transportation Science*, 54(6), 1471-1494.

976 ReVelle, C., & Hogan, K. (1988). A reliability-constrained siting model with local
977 estimates of busy fractions. *Environment and Planning B: Planning and Design*,
978 15(2), 143-152.

979 ReVelle, C., & Hogan, K. (1989a). The maximum availability location problem.
980 *Transportation Science*, 23(3), 192-200.

981 ReVelle, C., & Hogan, K. (1989b). The maximum reliability location problem and α -
982 reliable p-center problem: derivatives of the probabilistic location set covering
983 problem. *Annals of Operations Research*, 18(1), 155-173.

984 Schilling, D., Elzinga, D. J., Cohon, J., Church, R., & ReVelle, C. (1979). The
985 TEAM/FLEET models for simultaneous facility and equipment siting.
986 *Transportation Science*, 13(2), 163-175.

987 Schmid, V., & Doerner, K. F. (2010). Ambulance location and relocation problems with
988 time-dependent travel times. *European Journal of Operational Research*,
989 207(3), 1293-1303.

990 Sorensen, P., & Church, R. (2010). Integrating expected coverage and local reliability
991 for emergency medical services location problems. *Socio-Economic Planning*
992 *Sciences*, 44(1), 8-18.

993 Sudtachat, K., Mayorga, M. E., & McLay, L. A. (2014). Recommendations for
994 dispatching emergency vehicles under multitiered response via simulation.
995 *International Transactions in Operational Research*, 21(4), 581-617.

996 Toregas, C., Swain, R., ReVelle, C., & Bergman, L. (1971). The location of emergency
997 service facilities. *Operations Research*, 19(6), 1363-1373.

998 Toro-Díaz, H., Mayorga, M. E., Chanta, S., & McLay, L. A. (2013). Joint location and
999 dispatching decisions for emergency medical services. *Computers & Industrial*
1000 *Engineering*, 64(4), 917-928.

1001 van den Berg, P. L., & Aardal, K. (2015). Time-dependent MEXCLP with start-up and
1002 relocation cost. *European Journal of Operational Research*, 242(2), 383-389.

1003 Vos, T., et al. (2020). Global burden of 369 diseases and injuries in 204 countries and
1004 territories, 1990–2019: a systematic analysis for the Global Burden of Disease
1005 Study 2019. *The Lancet*, 396(10258), 1204-1222.

1006 [Yellow Taxi Trip Data \(2011\). https://data.cityofnewyork.us/Transportation/2011-](https://data.cityofnewyork.us/Transportation/2011-Yellow-Taxi-Trip-Data/jr6k-xwua)
1007 [Yellow-Taxi-Trip-Data/jr6k-xwua.](https://data.cityofnewyork.us/Transportation/2011-Yellow-Taxi-Trip-Data/jr6k-xwua)

1008 Yoon, S., Albert, L. A., & White, V. M. (2021). A stochastic programming approach for
1009 locating and dispatching two types of ambulances. *Transportation Science*,
1010 55(2), 275-296.

1011 Zarkeshzadeh, M., Zare, H., Heshmati, Z., & Teimouri, M. (2016). A novel hybrid
1012 method for improving ambulance dispatching response time through a
1013 simulation study. *Simulation Modelling Practice and Theory*, 60, 170-184.

1014