

**Manuscript version: Author's Accepted Manuscript**

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

**Persistent WRAP URL:**

<http://wrap.warwick.ac.uk/134032>

**How to cite:**

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

**Copyright and reuse:**

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

**Publisher's statement:**

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: [wrap@warwick.ac.uk](mailto:wrap@warwick.ac.uk).

## ROBUST CAPACITY PLANNING FOR ACCIDENT AND EMERGENCY SERVICES

ELVAN GÖKALP<sup>1</sup>

**Abstract.** Accident and emergency departments (A&E) are the first place of contact for urgent and complex patients. These departments are subject to uncertainties due to the unplanned patient arrivals. After arrival to an A&E, patients are categorized by a triage nurse based on the urgency. The performance of an A&E is measured based on the number of patients waiting for more than a certain time to be treated. Due to the uncertainties affecting the patient flow, finding the optimum staff capacities while ensuring the performance targets is a complex problem. This paper proposes a robust-optimization based approximation for the patient waiting times in an A&E. We also develop a simulation optimization heuristic to solve this capacity planning problem. The performance of the approximation approach is then compared with that of the simulation optimization heuristic. Finally, the impact of model parameters on the performances of two approaches is investigated. The experiments show that the proposed approximation results in good enough solutions.

**Keywords:** Health-care modelling and Capacity planning and Accident and Emergency Services and Queuing theory and Simulation Optimization and Robust Optimization.

### INTRODUCTION

Healthcare is one of the largest sectors affecting millions of lives worldwide. The aging phenomenon, increased rates of long-term conditions and public access to the healthcare generated a dramatic growth in the demand [37]. On the other hand, the capacity has not risen sufficiently to match this growth due to the inflexibility and scarcity of resources. Inevitably, the healthcare managers are under a significant pressure to improve the existing capacity and resource allocation policies. The pressure for an efficient service is accompanied by the challenges inherited from the nature of the healthcare services. The most important one of these challenges is the variation in the demand and the time required to treat patients. These uncertainties result in long waiting times to receive the service in the peak demand sessions. The delays in the treatment risk the patient lives, and thus, should be avoided.

The service delays is especially critical on the patient outcomes in accident and emergency departments (A&E). This unit is the first point of contact for most of the complex and life-threatening cases such as heart attack, stroke or loss of consciousness. Along with the overall increase in healthcare service demand, emergency unit arrivals have risen by 28%

---

December 2, 2019.

<sup>1</sup> Warwick Manufacturing Group, The University of Warwick, Coventry, CV4 7AL, UK.

between 2002 and 2017 in England [44]. To improve the efficiency in A&E's, the UK government has introduced the 4-hr policy in 2004, requiring all patients in A&E to be treated within 4 hours of their arrival. The target has been then reduced to 98% of patients in 2005 [42] and further to 95% in 2010 because the national levels fell below the previous targets [44]. In February 2018, the 4-hr target has been suspended [45] after the national levels gradually fell to 83%. Although the financial sanctions are suspended, the hospitals are still obliged to report their A&E waiting times.

The UK government is planning to put alternative measures for emergency waiting times such as tighter waiting hour targets for serious cases [57]. In an A&E, the severity of cases are initially assessed by a triage nurse who also categorizes them. While urgent patients are prioritized to receive the main treatment, rest of the patients are served based on the first-come first-served (FCFS) policy. Although prioritization improves the patient outcomes, the levels of the human and physical resources have a significant effect on the waiting times. The physical resources such as the number of cubicles where the triage and treatment takes place in A&E cannot be changed much. On the other hand, the level of human resources i.e. triage nurses and doctors is a tactical decision and easier to adjust based on the long-term demand projections.

A significant challenge for A&E resource planning is the daily and seasonal variation in the demand. Besides, the time required to treat patients vary significantly; the urgent cases take significantly more time than the non-urgent ones. These variations result in very long waiting times in the high-demand sessions. Finding the optimum staffing capacity to reduce these *worst-case* waiting times is a difficult problem due to the uncertainties. Additionally, the uncertainty in arrival and waiting times may not follow a known distribution and therefore the underlying assumptions for classical (queuing) approaches are violated. This paper proposes a novel approach based on robust optimization and queuing theory to find the optimum staffing capacities in an A&E to keep the worst-case waiting times below a certain threshold. The proposed approach does not require any assumption regarding the distribution of uncertainties. The contributions of the paper are:

- (1) The healthcare service in A&E is modelled using an approximation for the maximum waiting times combining the robust optimization and queuing theory, assuming that the arrival and service times can follow any distribution, and separating queues for urgent and non-urgent cases with prioritization. The results show that the resulting model can be solved to optimality.
- (2) As a benchmark to the proposed optimization approach, we also implement a simulation optimization (SO) based heuristic to find the optimum staff capacities. The performances of the solutions obtained by the optimization model and the SO based heuristic are then compared. The impact of different model parameters on the solutions are also investigated.

The paper is organized as follows. The next section introduces the literature related to capacity planning problems in healthcare and accident and emergency department modelling. Section 2 provides the problem description and underlying assumptions along with the optimization model formulation. Section 3 presents the worst-case waiting time approximation and the structural analysis of the optimization model. Section 4 details the proposed SO heuristic. In Section 5, we

introduce the design of experiments and results. Finally, Section 6 summarizes the study and provide several future directions.

## 1. RELATED LITERATURE

This section presents the related literature to capacity planning problem for A&E. First, we provide an overview of the capacity planning studies in healthcare. Then, we only focus on the A&E modelling and categorize the literature based on the modelling and solution techniques.

### 1.1. CAPACITY PLANNING IN HEALTHCARE

Queuing theory, a modelling approach to obtain performance measures in queuing systems, has been widely applied for the capacity planning of healthcare services; a related review can be found in [21]. Creemers et al. [16] use built-in queuing formulas to find the number of servers, i.e., capacity level, required to achieve a certain degree of performance. Hulshof et al. [32] use the queuing theory to model the elective patient admission process and study the resource allocation problem for hospitals with uncertain treatment pathways. They consider different queues for different types of services with time-dependent capacity levels of resources. Similarly, Cochran et al. [13] apply the queuing theory to test various capacity design alternatives to be used in real time Hospital Emergency Departments when the capacity cannot meet the demand. Bretthauer et al. [10] consider the capacity planning problem for healthcare operations with blocking between different units. Castillo et al. [11] determine capacity and location of healthcare facilities using queuing models with exponential service times and Poisson arrivals. By considering time-varying demands in hospitals, Green et al. [26] analyze the staffing requirement in hospitals based on queuing analysis. Mingzhu [40] develops a queuing network analysis considering multiple patient types to find optimum number of servers in an outpatient clinic. The main drawback of queuing models comes from their intractability due to nonlinear formulations of performance metrics under certain distribution assumptions for arrival and service processes.

Simulation is an alternative approach to model the service systems when the queuing formulations are not useful due to their complexities. Harper et al. [30] introduce a discrete-event simulation model to analyze the operations management of an intensive care unit and use the data generated by the simulation approach to solve the stochastic optimization model which computes the optimum number of nurses required to achieve the service targets. De Angelis et al. [17] consider SO to determine the capacity of a transfusion centre under multiple objectives: cost minimization to achieve a fixed waiting time and minimization of waiting time under a limited budget. The queuing system is modelled with a discrete-event simulation and the objective functions are approximated by function fitting with data generated by the simulation model. Similarly, Alfonso et al. [2] model processes in a blood collection unit with a simulation-based approach. They evaluate possible blood-collection server configurations from a cost-effectiveness perspective. Although simulation is very useful to model complex systems, it can only provide approximate solutions that are affected by the bias of data generation.

Optimization models in healthcare capacity planning focus not only on single hospital or department but also the interconnection between departments and hospitals, which usually has significant effects on the overall performance. Several studies focus on this interconnection in different capacity planning problems modelled for networks of hospitals or departments [4, 5, 19, 24, 25, 28, 38, 50, 54, 55, 56]. Flessa [19] develops a model to allocate resources in the preventive and curative services in hospitals. Govind et al., Gunes et al., Santibanez et al. and Stummer et al. [25, 28, 54, 55] focus on the location and number of beds in hospitals within a network to minimize operation cost and maximize patient utility.

Pehlivan et al. [50] develop a mixed-integer optimization model to determine the capacity of maternity facilities in a network in view of uncertain patient arrivals and service times. The objective is to minimize the number of refused admissions which is formulated by using available queuing formulations. They assume the interarrival and service times are exponentially distributed. On the other hand, Asaduzzaman et al. [5] develop a queuing model to find the optimum capacities of neonatal centres to minimize refusal and overflow probabilities. They also assume exponential interarrival and service times. To the best of our knowledge, the proposed approximation in this paper has not been utilized before for a capacity planning problem.

## 1.2. ACCIDENT AND EMERGENCY MODELLING

### *Simulation Modelling and Optimization:*

Mohiuddin et al. [41] identify and review 19 studies related to simulation modelling for emergency departments in the UK. Another comprehensive review of simulation modelling studies in emergency departments for normal and disaster conditions can be found in Gul et al. [27]. Among 106 reviewed papers, only few studies consider an optimization approach [1, 18, 22, 51, 59]. Fruggiero et al. [22] uses ant-colony optimization along with a simulation model to optimize the resources in an emergency department. Weng et al. [59], Ibrahim et al. [33] and Rico et al. [51] use OptQuest, a SO engine [49], to optimize nurse and physician numbers and nurse allocation in an influenza outbreak, respectively. Ghanes et al. [23] also uses SO to find the staffing levels in A&E with the objective of minimizing the patient length of stay. Chen and Wang [12] aim to find optimum number of staffing in A&E minimizing patient length of stay and the medical resources wasted by SO.

The most similar study to ours is Ahmed et al. [1] which presents a SO approach for capacity planning of an emergency department in Kuwait. They consider triage and a prioritized service queue. They consider stochastic constraints in a discrete SO problem maximizing the throughput. They also model and solve the optimization problem where the objective function is total cost of the staff and the constraints are the waiting times. This second problem description is very similar to ours. Their heuristic first identifies the feasible set of solutions and then finds the best solution among those based on random sampling. In each iteration, they compare the objective value of the new solution to the previous one and accept the new one if it supersedes in certain number of iterations. However, they have not provided any performance results for their heuristic results.

### *Mathematical Modelling and Queuing Theory:*

Queuing theory has been mostly been applied into staff and bed optimization, ambulance deployment problems; for a detailed review, readers are referred to [36]. Emergency departments are also modelled by queuing theory [60]. Reviews of modelling and queuing theory studies for A&E can be found in Saghafian et al., Costa et al. and Hu et al. [15, 31, 53]. These papers have considered the average instead of maximum waiting times. Few papers model separate queues based on patients' severity [13]. For example, Cochran and Roche [13] develop a queuing network model for an A&E that uses separate queues for low and high acuity patients. Optimum capacity (either in terms of staff or waiting area limit) for each step of patient flow in A&E is computed. They use an approximate waiting time formulation [3] and target waiting times and utilization rates of each step to set up the capacities assuming that the arrival times are exponentially distributed.

Mayhew et al. [39] develop a queuing model for A&E department assuming that the arrivals and service times are exponentially distributed. They have divided the arrivals into two and added a triage step before the treatment. They have compared the predicted overall departure time from A&E with the real departure times in A&E's obtained from the NHS UK. They have used the model to test whether the 4-hour target is achievable if part of the A&E service is carried in other units. They have not carried optimization or capacity planning.

## 2. CAPACITY PLANNING MODEL FOR A&E

This section first describes the underlying problem and our assumptions and then introduces the mathematical formulation.

### 2.1. PROBLEM DESCRIPTION AND ASSUMPTIONS

We model the activities in a typical (major) A&E department in the UK [48] for a finite planning horizon. With small modifications, the model can be applied to any other emergency department. As a patient arrives to the A&E, s/he is put into an FCFS queue for triage. A triage nurse categorizes the patient as discharge, a type 1 (urgent) or type 2 (non-urgent) based on the medical assessment. Note that we do not model specific triage categories which may be more than 2. Instead, we only divide them based on whether the patient is urgent or not, as in Ahmed et al. [1]. Type 1 and 2 patients are placed into two separate FCFS queues for the treatment. Type 1 patients are given priority for treatment. After the patients are treated by an A&E doctor, they are either discharged, referred or admitted to the hospital. The time spent in the A&E from arrival until the disposal (discharge, referral or admission) should be less than 4 hours for all patients. We assume that other medical activities required for the treatment such as laboratory tests are included in the treatment duration. Also note that we do not consider the single specialty cases such as ophthalmology or dental. These patients go through a separate route than the other two categories in A&E [45] and constitute a small percentage (0.5%) of the 4 hour breaches.

The uncertainties affecting the waiting times are the arrival times, the triage and treatment durations. Although arrival times can follow an exponential distribution, there is no consensus in the literature about the service time distributions

in A&E; triangular [20], general [34], exponential [53] and uniform [1] distributions are used to model the A&E treatment duration.

The hospital management aims to find the staff capacities with minimum cost while satisfying the waiting time targets. One intuitive approach for this capacity planning problem would be to allocate resources proportionate to the demand rates for each service. However, that method would not match the actual workloads. The actual congestion peak times lag the arrival times as the number of service stations as in the A&E increases [34]. The approximation and heuristic approaches provided in this paper use queuing model to estimate the actual workloads, and thus provide a better performance than simple allocation of resources based on the demand rates [34].

Note that the arrival rates to an A&E can vary based on the time of the day. Here, we only consider a stable arrival rate. In other words, we approximate the varying arrival rate with its average value. The reason for this assumption is that the model can be easily extended to time-dependent arrival rates, by simply adding time indices to the arrival rates. In such a case, one would find the staff capacities for each time period. This can easily be done by following our approach separately for each time period. In other words, extension to time-dependent arrival rates would not affect the model complexity or structure and thus the main objectives of this paper.

Another modelling choice is related to the 'boarding'. This term refers to the cases where after medical treatment is completed, the patient may need to wait for a bed in the hospital (if admitted). This may create additional delay on the patient's length of stay. However, since A&E beds are highly utilized and expensive resources, some hospitals put these admitted patients into 'buffer' wards such as Critical Decision Units [43]. Besides, the 'boarding' process would require us to model all the bed utilization in all wards of the hospital which is beyond the scope of this paper.

## 2.2. PROBLEM FORMULATION

This section provides a mathematical formulation for the capacity planning problem described above. The number of nurses in the triage is denoted by  $x_1$ . The fixed (unit) cost of triage nurses and doctors are shown with  $c_1$  and  $c_2$ , respectively. We assume that the patients arrive to the A&E with the mean interarrival time  $1/\lambda$  and standard deviation  $\sigma$ . After registration, they wait in the triage queue and assessed by a triage nurse under the FCFS rule. The mean triage time is  $1/\mu$  with standard deviation  $\sigma_1$ . Maximum total time spent in the triage by any patient arrived during the planning horizon is  $W_1(x_1)$ . A certain percentage,  $\theta$ , of the patients are discharged after triage. The others are categorized as type 1 or type 2 each of which has a separate queue for treatment. The rate of type 1 patients among whole arrivals is  $\alpha$ . Patients wait in the treatment queues until they are seen by one of  $x_2$  number of A&E doctors. Type 1 cases have the priority over type 2 and the queue is preemptive: the treatment of a type 2 case is stopped when a type 1 arrival occurs at the same time. The average treatment time for type 1 and type 2 are  $1/\mu_{12}$  and  $1/\mu_{22}$  with standard deviations  $\sigma_{12}$  and  $\sigma_{22}$ , respectively.

The hospital management aims to keep the total time spent in the A&E below  $\overline{W}$ . Due to the uncertainty in the arrival and service times, the waiting times experienced by the patients vary. Therefore, the model should be robust against

the uncertainties in the arrival and service times. We denote the maximum time spent between triage and treatment as  $W_{12}(x_2)$  and  $W_{22}(x_2)$  for type 1 and 2 patients, respectively. An approximation for the maximum time spent in the A&E can be written as:

$$W_1(x_1) + W_{12}(x_2) + 1/\mu + 1/\mu_{12},$$

and

$$W_1(x_1) + W_{22}(x_2) + 1/\mu + 1/\mu_{22},$$

for type 1 and 2 patients, respectively. Figure 1 shows a summary of the A&E operations along with the notation used in the model.

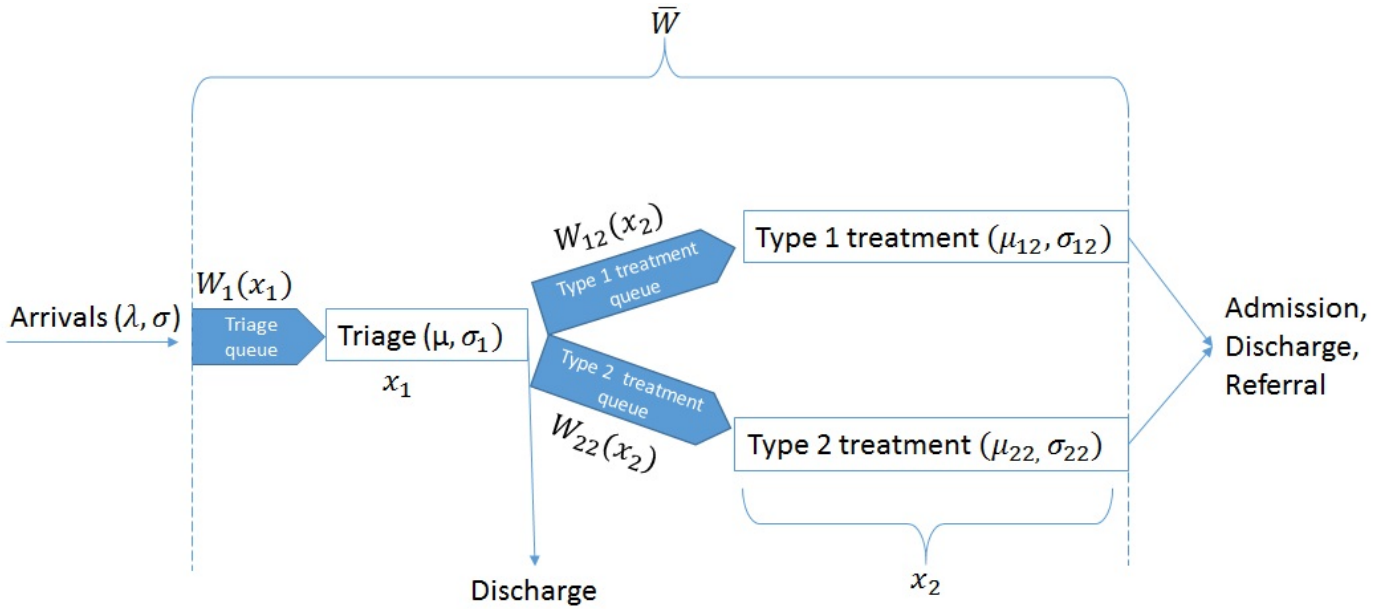


FIGURE 1. A description of the A&E service along with the notation used.

For a stable queue, the utilization rate (traffic intensity) should be smaller than 1 [35];  $\frac{\lambda}{x_1\mu} < 1$  for the triage queue. In other words, the total service rate ( $x_1\mu$ ) should be larger than the total arrival rate ( $\lambda$ ) such that the queue does not grow exponentially. This condition should be satisfied for the treatment queues as well,

$$x_2 > \lambda_{12}/\mu_{12}, \quad x_2 > \lambda_{22}/\mu_{22},$$

where  $\lambda_{12}$  and  $\lambda_{22}$  represent the arrival rates to the treatment queues of type 1 and type 2 patients, respectively.

The stochastic capacity planning model for the A&E can be formulated as follows:



$$AE : \quad \min \quad c_1 x_1 + c_2 x_2, \quad (1)$$

$$s.t. \quad \overline{W} \geq W_1(x_1) + W_{12}(x_2) + 1/\mu + 1/\mu_{12}, \quad (2)$$

$$\overline{W} \geq W_1(x_1) + W_{22}(x_2) + 1/\mu + 1/\mu_{22} \quad (3)$$

$$x_1 > \lambda/\mu, \quad (4)$$

$$x_2 > \lambda_{12}/\mu_{12}, \quad (5)$$

$$x_2 > (\lambda_{12} + \lambda_{22})/\mu_{22}, \quad (6)$$

$$x_1, x_2 \in Z^+. \quad (7)$$

As mentioned before, the government is planning to put waiting time targets for serious (type 1) patients only. Assuming that this policy is activated, the problem would be then modelled as:

$$AE^{red} : \quad \min \quad c_1 x_1 + c_2 x_2, \\ s.t. \quad (2), (4); (5); (7).$$

In the case of extension to a time-dependent arrival rate with  $\lambda_{12}(t)$  and  $\lambda_{22}(t)$ , the model variables would be differentiated for each time period, e.g.  $x_1(t)$  and  $x_2(t)$ , while the rest of the model would stay same. Since that extension does not affect the model, we continue with time-independent version in the rest of the paper.

In order to solve these models, we need to compute the waiting times for each patient arriving to the A&E. The exact computation of the waiting times in each scenario is difficult even with a fixed number of staff. The computational intractability due to combinatorial number of calculations has already been proven for a queuing system of multiple servers with exponential arrivals and general service time distribution as in Tijms et al. [58]. The next section approximates the maximum waiting times by using robust optimization principles and provides a tractable approximate model.

### 3. APPROXIMATION WITH ROBUST OPTIMIZATION

This section presents an approximation for the maximum waiting times in the A&E. Different approaches exist to approximate the maximum waiting time in a queuing system; for instance, see Gupta et al. [29]. However, these approximations usually do not lead to realistic results when the arrival process follows a distribution different from Poisson [7]. As an alternative approach, Bandi and Bertsimas [7] proposed to approximate the maximum waiting time in an FCFS queue when the arrival and service times follow an unknown distribution. Their approach is based on developing *uncertainty sets* for the uncertain arrival and service durations based on historical data. They used the central limit theorem which asserts the asymptotic results for a large set of independent and identically distributed random variables.

Readers are referred to [7] for more details regarding the approximation. According to Bandi and Bertsimas [7], the maximum waiting time in an FCFS queue with  $x_1$  servers  $\overline{W}(x_1)$ , arrival and service rates  $\lambda$  and  $\mu$ , arrival and service time variabilities  $\Gamma^a$  and  $\Gamma^s$  can be approximated as,

$$\overline{W}(x_1) = \frac{\lambda(\Gamma^a + \Gamma^s/\sqrt{x_1})^2}{4[1 - \lambda/(\mu x_1)]}. \quad (8)$$

The arrival and service time variabilities are set based on the desired conservativeness level. For example, they can be set as double or three times of the standard deviation of the corresponding uncertain parameter to cover most of the possible realizations.

We use (8) to approximate the maximum waiting times in the A&E. The variability in the arrival times to the A&E is denoted by  $\Gamma^a$  and the variabilities in triage duration and treatment duration for type 1 and type 2 patients are denoted by  $\Gamma^s, \Gamma_{12}^s, \Gamma_{22}^s$ , respectively. Based on [7], we set  $\Gamma^a = k \cdot \sigma_a$  and  $\Gamma^s = k \cdot \sigma_s$  with  $k > 0$ , where  $\sigma_a$  and  $\sigma_s$  are the standard deviations of the corresponding interarrival and service times, respectively. The parameter  $k$  is set based on the desired conservativeness level of the model: a larger  $k$  corresponds to a more conservative model against the uncertainties in the arrival and service times. We present the effect of different conservativeness levels on the results in Section 5.2. We approximate the maximum waiting times in the triage and type 1 treatment queues as follows:

$$\overline{W}_1(x_1) = \frac{\lambda(\Gamma^a + \Gamma^s/\sqrt{x_1})^2}{4[1 - \lambda/(\mu x_1)]}, \quad (9)$$

$$\overline{W}_{12}(x_2) = \frac{\lambda_{12}(\Gamma_{12}^a + \Gamma_{12}^s/\sqrt{x_2})^2}{4[1 - \lambda_{12}/(\mu_{12}x_2)]}. \quad (10)$$

Note that the computation of the waiting times of type 2 patients is more complicated. A type 2 patient in the treatment queue is always served after all existing type 1 patients are served. This would imply that type 2 patients always wait more than a type 1 patient. In the worst case, a type 2 patient would wait for the maximum waiting time for a type 1 patient in addition to the maximum waiting time in type 2 treatment queue. Then, an approximation for the maximum type 2 patient waiting time can be formulated as:

$$\overline{W}_{22}(x_2) = \frac{\lambda_{22}(\Gamma_{22}^a + \Gamma_{22}^s/\sqrt{x_2})^2}{4[1 - \lambda_{22}/(\mu_{22}x_2)]} + \overline{W}_{21}(x_2). \quad (11)$$

The model  $AE$  and  $AE^{red}$ , respectively, can be reformulated as:

$$\begin{aligned} AE_{rob} : \quad & \min \quad c_1 x_1 + c_2 x_2, \\ s.t. \quad & \overline{W} \geq \frac{1}{\mu} + \frac{1}{\mu_{12}} + \overline{W}_1(x_1) + \overline{W}_{12}(x_2), \end{aligned} \quad (12)$$

$$\overline{W} \geq \frac{1}{\mu} + \frac{1}{\mu_{22}} + \overline{W}_1(x_1) + \overline{W}_{22}(x_2), \quad (13)$$

4, 5, 67.

$$\begin{aligned}
AE_{rob}^{red} : \quad & \min \quad c_1 x_1 + c_2 x_2, \\
& s.t. \quad 12, 4, 5, 7.
\end{aligned}$$

which have linear objective functions and non-linear constraints. The next proposition states that the relaxed version of  $AE_{rob}$  has a convex feasible set. The same proposition and proof apply to  $AE_{rob}^{red}$ . Therefore, both models have global optimums [9].

**Proposition 1.** *For relaxed variables  $x_1, x_2 \in R^+$ , the model  $AE_{rob}$  has a convex feasible set.*

*Proof.* Let's first show the convexity of constraint (12). Let's define  $f(x_1, x_2) = \overline{W}_1(x_1) + \overline{W}_{12}(x_2)$ . For a multi-variate function to be convex, its Hessian matrix should be a positive semi-definite matrix. Let  $H = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$  denote the Hessian matrix of function  $f(x_1, x_2)$ , where its second order derivatives are denoted by  $a, b$  and  $c$ . For the function to be convex, all principal minors,  $a, c, (ac - b^2)$ , should be non-negative. Note that  $b = \frac{\partial^2 f(x, \alpha)}{\partial x \partial \alpha}$  is always zero, because function  $f(x_1, x_2)$  can be divided into two separate functions of variables  $x_1$  and  $x_2$ . Therefore, it is enough to show that  $a = \frac{\partial^2 f(x_1, x_2)}{\partial x_1^2}$  and  $c = \frac{\partial^2 f(x_1, x_2)}{\partial x_2^2}$  are non-negative. These derivatives can be written as:

$$a = \frac{\partial^2 \overline{W}_1(x_1)}{\partial x_1^2}, \quad c = \frac{\partial^2 \overline{W}_{12}(x_2)}{\partial x_2^2},$$

due to the separability of the maximum waiting time functions. We arrive the formulation of the second order derivative  $a$  after some intermediate calculations as:

$$a = \frac{\Gamma^a \Gamma^s (-m^2 + 6my^2 + 3y^4) + 4y^3 (m(\Gamma^a)^2 + (\Gamma^s)^2)}{2y^3(y^2 - m)^3},$$

where  $y = \sqrt{x_1}$ , and,  $m = \lambda/\mu$ . The denominator of  $a$  is always positive due to the traffic intensity condition  $\mu x_1 > \lambda$ . The absolute of the only negative term in the nominator,  $-m^2 \Gamma^a \Gamma^s$ , is always smaller than the second term of the nominator  $6m \Gamma^a \Gamma^s y^2$  because  $m < 1$  and  $x_1 \geq 1$ . Therefore, the second order derivative  $a$  is always positive. We omit the computations for  $c$  that follows the same structure as  $a$ . Since all principal minors of  $H$  are non-negative,  $f(x_1, x_2)$  is convex. The second constraint (13) possesses the same structure, and is also convex. Note that all other constraints, (4), (5), (6) are also convex. Because all constraints are convex, the relaxed version of model  $AE_{rob}$  has a convex feasible set.  $\square$

Since the relaxed version of  $AE_{rob}$  is a convex optimization problem with a global optimum, convex non-linear optimization solvers such as Bonmin [14] can find the global optimum solution efficiently following Bonami et al. [8].

#### 4. SIMULATION OPTIMIZATION HEURISTIC

The models  $AE$  and  $AE^{red}$  are non-linear integer programming models that are very difficult to solve with traditional optimization techniques. A possible solution approach for these models is SO that is based on simulating alternative solutions and comparing the simulated objective function values. An intuitive method for SO is to first enumerate all feasible solutions and simulate their performances. However, this method would require a long computation time, and therefore, more clever search strategies are needed. Our literature review shows that most of the SO studies employ built-in optimization packages within a commercial simulation modelling software such as OptQuest in Simul8.

To investigate the performance of the worst-case approximation and optimization approach presented in Section 3, we also design and implement an SO heuristic as a benchmark solution method. The performances of this heuristic and the optimization via a commercial non-linear integer solver are then compared.

Other than the stochastic waiting time constraints, our model has a set of deterministic constraints related to traffic intensity (7) and a monotonic objective function. Therefore, we do not need to search for the optimum solution randomly as in Ahmed et al. [1]; we can start the search from the smallest capacities satisfying the traffic intensity constraints and increase these capacities incrementally until the stochastic constraints are satisfied with a certain confidence rate.

For this purpose, we first develop a simulation model of the A&E operations described in Section 2.1 and implement it on Matlab. The planning period of the simulation is set to  $T$  minutes, while the time unit is one minute. One iteration of the simulation model comprises of  $j$  runs of  $n$  scenarios. In each iteration, the SO heuristic searches for a better solution based on the simulation outputs. Starting from the minimum possible levels of the capacity variables (satisfying the traffic intensity constraints), in each iteration, we increment the capacity variable that has the largest potential to decrease the maximum waiting times. To identify the variable with the largest potential improvement, we use the approximate maximum waiting time formulation (8). Note that the objective function (total cost) increases by  $c_1$  and  $c_2$  with one unit increase in  $x_1$  and  $x_2$ , respectively. In other words, increasing  $x_2$  by one would result in the same change in the objective value as increasing  $x_1$  by  $c_2/c_1$ . Therefore, the potential improvements in the maximum waiting times should be computed for  $x_1 + c$  and  $x_2 + 1$ , respectively, where  $c = \lceil c_2/c_1 \rceil$  due to the integrality condition. The capacity variable with the largest potential improvement is incremented by one and the simulation model is run for another iteration. In each iteration, if  $(1 - \epsilon)\%$  of the patients' total waiting time is lower than  $\overline{W}$ , where  $\epsilon$  is the desired confidence level, e.g. 1%, then the heuristic stops. Otherwise, the process is repeated again. Algorithm 1 presents the pseudo-code of the SO heuristic. The heuristic can be applied to  $AE^{red}$  by just removing the parameters related to type 2 treatment as presented in Algorithm 2.

The computation time of the heuristic depends on the running time of the simulation model and therefore the levels of  $n$ ,  $j$  and  $T$ . As the number of runs and scenarios increases, the robustness of the solution obtained by the heuristic increases as well.

**Algorithm 1** SO Heuristic for  $AE$ 


---

Set  $\epsilon$ ,  $W^i$  to a very large number,  $i = 0$  and compute  $x_1^0 = \lambda/\mu$ ,  $x_2^0 = \min\left\{\frac{\lambda_{12}}{\mu_{12}}, \frac{\lambda_{22}}{\mu_{22}}\right\}$ , and  $W^0 = \max\left\{(\overline{W}_1(x_1^0) + \overline{W}_{21}(x_2^0)), (\overline{W}_1(x_1^0) + \overline{W}_{22}(x_2^0))\right\}$  using 9, 10, 11.

**while**  $W^i \geq \overline{W}$ , **do**

    Compute  $\Delta W_1 = \overline{W}_1(x_1^i) - \overline{W}_1(x_1^i + c)$  and

$\Delta W_2 = \max\left\{(\overline{W}_{12}(x_2^i) - \overline{W}_{12}(x_2^i + 1)), (\overline{W}_{22}(x_2^i) - \overline{W}_{22}(x_2^i + 1))\right\}$ .

**if**  $\Delta W_1 > \Delta W_2$  **then**

$x_1^{i+1} = x_1^i + 1$ ,  $x_2^{i+1} = x_2^i$ .

**else**

$x_2^{i+1} = x_2^i + 1$ ,  $x_1^{i+1} = x_1^i$ .

**end if**

$i := i + 1$ .

    Run simulation model for  $j$  runs and  $n$  scenarios with  $x_1^i$  and  $x_2^i$ . Set  $W^i$  to  $(1 - \epsilon)\%$  of the waiting times obtained by the simulation model.

**end while**

**return**  $x_1^i$  and  $x_2^i$ .

---

**Algorithm 2** SO Heuristic for  $AE^{red}$ 


---

Set  $\epsilon$ ,  $W^i$  to a very large number,  $i = 0$  and compute  $x_1^0 = \lambda/\mu$ ,  $x_2^0 = \left\{\frac{\lambda_{12}}{\mu_{12}}\right\}$ , and  $W^0 = (\overline{W}_1(x_1^0) + \overline{W}_{21}(x_2^0))$  using 9, 10, 11.

**while**  $W^i \geq \overline{W}$ , **do**

    Compute  $\Delta W_1 = \overline{W}_1(x_1^i) - \overline{W}_1(x_1^i + c)$  and

$\Delta W_2 = \overline{W}_{12}(x_2^i) - \overline{W}_{12}(x_2^i + 1)$ .

**if**  $\Delta W_1 > \Delta W_2$  **then**

$x_1^{i+1} = x_1^i + 1$ ,  $x_2^{i+1} = x_2^i$ .

**else**

$x_2^{i+1} = x_2^i + 1$ ,  $x_1^{i+1} = x_1^i$ .

**end if**

$i := i + 1$ .

    Run simulation model for  $j$  runs and  $n$  scenarios with  $x_1^i$  and  $x_2^i$ . Set  $W^i$  to  $(1 - \epsilon)\%$  of the waiting times obtained by the simulation model.

**end while**

**return**  $x_1^i$  and  $x_2^i$ .

---

## 5. COMPUTATIONAL EXPERIMENTS

The computational experiments aim to illustrate the performances of the approximation approach and the SO heuristic as well as the impact of several model parameters on the results. For this purpose, we design two sets of computational experiments. The first set of experiments compares the performances of the solutions computed by the approximate optimization models and the SO heuristic. The second set of the experiments investigates the impact of model parameters on the solutions obtained by the approximation approach. All computational experiments are carried out on a PC with Windows 10 Enterprise operating system, CPU 4GHz Intel Core i7 and 32GB of RAM.

## 5.1. INPUT DATA

As all major A&E's in the UK follow the same service process, we use the arrival data of University Hospitals Coventry & Warwickshire (UHCW) provided in the online resources of the NHS UK [44]. The average treatment and triage times

are obtained from Ahmed et al. [1]. The service times are assumed to follow an exponential distribution [53]. The effect of this assumption is investigated in the first set of computational experiments.

TABLE 1. Input data for model parameters used in the numerical experiments

Description of Parameters	Value/Range	Source of Data	Distribution
A&E arrival rates	0.57 patient/minute	[44], [1]	Exponential
Probability of discharge or diagnosis as type 1 patient after triage, respectively	0.1 and 0.61	[44]	Binomial
Mean triage service duration	15 minutes	[1], [53]	Exponential
Mean treatment duration for type 1 patients	90 minutes	[1], [53]	Exponential
Mean treatment duration for type 2 patients	25 minutes	[1], [53]	Exponential
Cost of doctors with respect to nurses ( $c$ )	5	[52], [47]	-

For the variability parameters, we first generate a dataset of arrival and service times by using the simulation model and the distribution information provided in Table 1. According to [7], the variability parameters ( $\Gamma^a$  and  $\Gamma^s$ ) are then set such that most of the uncertain parameters are covered. The time spent in the A&E ( $\bar{W}$ ) should be less than 4 hours for all patients.

## 5.2. COMPARISON OF SO HEURISTIC AND APPROXIMATE OPTIMIZATION MODEL

This section presents the results and the performances of two solution approaches for different (i) problem settings, (ii) conservativeness levels and (iii) service time distribution. The approach that solves the models,  $AE_{rob}$  and  $AE_{rob}^{red}$ , with a commercial solver (Gams/Bonmin) is referred as *Approximate Optimization (AO)* in the rest of this section.

**Impact of Problem Setting:** This section presents the results for the optimization models  $AE_{rob}$  and  $AE_{rob}^{red}$  solved by a commercial solver Gams, by using the solver Bonmin. Similarly, we solve the reduced and full models,  $AE$  and  $AE^{red}$ , with the SO heuristic that is implemented in Matlab with  $\epsilon = 0.0001$ . We set the planning horizon of the problem to 1500 minutes that is found to be large enough to observe the queue dynamics. The SO heuristic is run for 40, 30, 10, and 1 runs to understand the impact of the number of runs on the heuristic's performance.

First, we solve the problem where only type 1 patients are subject to 4 hours waiting time limit, i.e.  $AE^{red}$  and  $AE_{rob}^{red}$ . Table 2 shows the capacities found by AO and the SO heuristic for different number of runs and scenarios. The results of the SO heuristic with 40 runs is the same as those with 30 runs. The table also shows the computation times of these solution approaches in terms of seconds.

TABLE 2. Base results of AO and the SO Heuristic when only type 1 patients are subject to waiting time limit

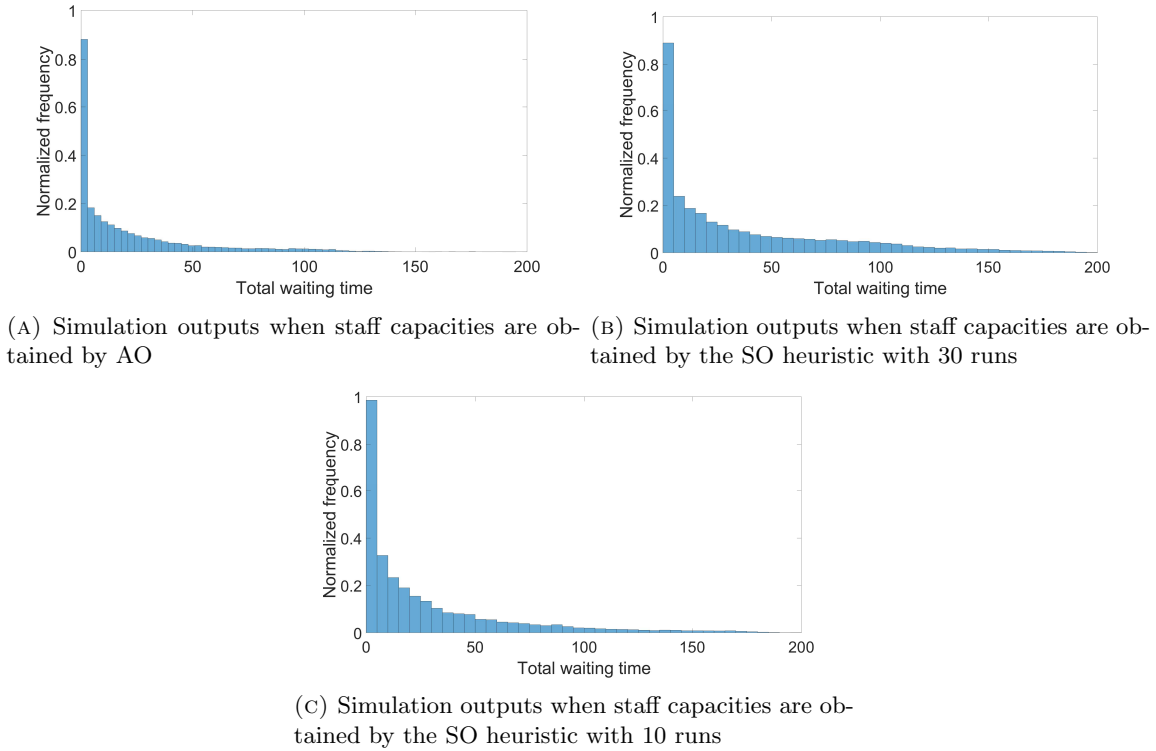
Approach	AO	SO Heuristic					
Number of runs ( $j$ )	-	<b>30</b>		<b>10</b>		<b>1</b>	
Number of scenarios ( $n$ )	-	1000	150	1000	150	1000	150
Capacities ( $x_1, x_2$ )	9, 31	9, 30	9, 29	9, 30	9, 28	9, 28	9, 28
Computation time (sec.)	2	3292	135	1161	198	52	4
Breaches	$10^{-5}\%$	$10^{-5}\%$	$10^{-5}\%$	$10^{-5}\%$	0.02%	0.02%	0.02%

The number of doctors found by two methods is slightly different; AO is more conservative to the uncertainties in the arrival and service times. The computation time of the SO heuristic is significantly larger than that of AO especially as the number of iterations increases. The conservativeness of the solutions found by the SO heuristic also increases with a higher number of iterations.

The NHS sources [46] indicate that there are around 58 full-time-equivalent doctors in UHCW in March 2018. Assuming the doctors make 2.5 shifts per day, this would be equivalent to 23 doctors. The NHS statistics show the percentage of patients treated within 4 hours in the A&E was 79.2% in that month [44]. Therefore, our results indicate that the performance can be improved by increasing the staff level from 23 to 31.

Next, we evaluate the performances of the solutions obtained by two approaches by giving these capacities to the simulation model as inputs. Figure 2 shows three frequency histograms for total waiting time in the A&E of (type 1) patients computed by the simulation model for 30 runs and 500 scenarios with the capacity levels found by AO and the SO heuristic for different number of runs.

FIGURE 2. Frequency histograms for total waiting time (of type 1 patients) in the A&E computed by the simulation model with the staff capacities obtained by different solution methods



The maximum waiting time in the A&E should not be larger than 135 minutes that is the difference between the limit on total time spent (240 minutes) and total average service time for triage and treatment (15 and 90 minutes, respectively). The graphs show that the frequency of patients waiting more than 135 minutes is significantly larger with the heuristic solutions, and the maximum waiting time computed by the Simulation model with AO solutions is 130 minutes while that

for the SO heuristic (with 30 iterations) is 150 minutes. This can be interpreted as the heuristic solutions perform worse than the AO solutions.

The results indicate that the proposed robust optimization based approximation is appropriate for the A&E capacity planning problem. They also suggest that AO performs better than the SO heuristic in terms of the computation time and the solution performance.

Second, we assume that the health authorities replace the waiting time target of 4 hours for both types of patients served in the A&E, i.e.  $AE$  and  $AE_{rob}$  are solved. Since the effect of number of runs is already shown, we set  $j = 10$  and  $n = 1000$  for the SO heuristic that produced the same solution with  $j = 30$  and  $n = 1000$ .

TABLE 3. Base results of AO and the SO heuristic when both types of patients are subject to 4 hours

Approach	AO	SO Heuristic	
Number of scenarios ( $n$ )	-	1000	100
Capacities ( $x_1, x_2$ )	12, 32	11, 32	11, 32
Computation time (sec.)	2	1003	210
Breaches	0.08%	0.08%	0.08 %

The computation times of both methods do not change significantly compared to the reduced problem setting. Note that type 2 patients have a lower priority and wait longer than type 1 patients. Therefore, the capacities are higher than those obtained for the reduced problem; they both increase by 2 units. We observe that the breaches have increased slightly for both methods compared to the previous problem setting. This may be due to increased complexity with the additional type 2 waiting time limit. These results indicate that the approximation of the waiting time for type 2 patients may not be as good as that for type 1 patients. In the rest of the experiments, we consider the reduced problem setting; only type 1 patients are subject to the 4-hour waiting time limit. Also, the SO heuristic is always run for  $j = 10$ ,  $n = 1000$  unless stated otherwise.

**Impact of Conservativeness Levels:** In this section, we investigate the impact of conservativeness levels on the solutions obtained by two approaches. As explained before, the variability parameters,  $\Gamma^a$  and  $\Gamma^s$ , define the conservativeness levels of the AO; a larger variability corresponds to a more robust model against the uncertainties in the arrival and service times. To investigate the effect of the conservativeness, we solve the optimization model for two more variability levels: the double and half of the base variability levels used in the previous set of experiments. Table 4 shows the optimum number of triage nurses and doctors in the doubled, base and halved variabilities.

TABLE 4. Capacities ( $x_1, x_2$ ) found by the AO in different variability levels ( $\Gamma^s$  and  $\Gamma^a$ )

Variabilities	Service Time	Arrival Time
Doubled	18, 37	9, 31
Base	9, 30	9, 31
Halved	8, 28	9, 30

As different from the service time variability, the arrival time variability does not affect the solutions significantly. This is probably due to a lower variance in the arrival times leading to a lower arrival time variability.



For the SO heuristic, the conservativeness level is defined via parameter  $\epsilon$ . Table 5 shows computation times and the solutions obtained by the SO heuristic in different ( $\epsilon$ ) levels. Note that the conservativeness level of the heuristic affects its computation time and the solutions significantly. Although a higher  $\epsilon$  results in a shorter computation time, the quality of the solution drops significantly.

TABLE 5. Computation time and capacities  $(x_1, x_2)$  found by the SO heuristic in different conservativeness levels ( $\epsilon$ )

$\epsilon$	<b>0.0001</b>	<b>0.001</b>	<b>0.05</b>	<b>0.1</b>	<b>0.2</b>
Capacities $(x_1, x_2)$	9, 30	9, 26	9, 24	9, 23	9, 22
Computation time (sec.)	1161	64.18	38.03	25.3	12.5
Breaches	0	0.0016 %	0.021 %	0.061 %	0.1294 %

Note that the number of doctors found for the conservativeness level of 0.2 is almost equal to the that in UHCW in March 2018, i.e. 23 doctors [46], in which the hospital had breaching rate of 0.2 [44].

**Impact of Service Time Distribution:** In this experiment, we investigate the impact of the distribution of treatment and triage durations on the results. For this purpose, the treatment and triage durations are assumed to be uniformly distributed, respectively, over  $[60, 120]$  and  $[10, 20]$  based on Ahmed et al. [1]. The solutions obtained by two approaches and the computation times are shown in Table 6. The rate of breaches is negligible for all solutions.

TABLE 6. Computation times and the solutions of AO and the SO heuristic when triage and treatment time distributions are uniform

<b>Approach</b>	<b>Aproximate Optimization</b>	<b>SO Heuristic</b>	
Number of scenarios ( $n$ )	-	1000	100
Capacities $(x_1, x_2)$	9 , 28	9, 28	9, 27
Computation time (sec.)	2	331	25

The computation times are not different to those with the exponential distribution assumption. However, the capacities obtained by the SO heuristic (with 1000 scenarios) and the AO are the same as different from the base case. This may indicate that the SO heuristic performs better when the service times follow a uniform distribution.

### 5.3. SENSITIVITY ANALYSIS

In this section, we test the impact of several model parameters on the solutions obtained by the AO.

**Impact of Waiting Time Limit:** As the NHS looks for alternative performance monitoring policies, we investigate the impact of different waiting time limits on the solutions obtained by the AO. Table 7 shows these capacities for different waiting time limits for type 1 ( $\overline{W}_1$ ) and type 2 ( $\overline{W}_2$ ) patients.

The results indicate that for a fixed type 2 waiting time limit, the solutions are not affected significantly when the waiting time limit for type 1 patients is above 4 hrs. Similarly, when the waiting time limit is above 5 hrs for type 2 patients, the solutions do not change.

**Impact of Patient Arrival Rates:** The NHS report [6] shows that the percentage of type 1 patients can vary in different hospitals. Therefore, this experiment investigates the effect of rates of type 1 patients and discharges after triage

TABLE 7. The solutions obtained by the AO for different waiting time limits

$\overline{W}_1$	$\overline{W}_2$		
	4 hrs	5 hrs	6 hrs
	$x_1, x_2$	$x_1, x_2$	$x_1, x_2$
<b>3 hrs</b>	12, 32	15, 31	15, 31
<b>4 hrs</b>	12, 32	10, 30	10, 30
<b>5 hrs</b>	12, 32	10, 30	9, 29

among all arrivals. Based on [6], we obtain the solutions by the AO for three rates for type 1 patients and two rates for discharge. Table 8 shows the solutions for different rates of discharge and type 1 patients. The base levels of these parameters are shown with \* in the table. We have not conducted the experiments for the unrealistic case where 0.75 of all arrivals are type 1 and 0.2 of all arrivals are discharged after triage.

TABLE 8. Computation time and results of optimization model and SO heuristic when both categories of patients are subject to waiting time limits

Ratio of type 1 patients after triage	Discharge rate after triage	
	0.1*	0.2
	$x_1, x_2$	$x_1, x_2$
<b>0.45</b>	16, 27	14, 27
<b>0.61*</b>	12, 32	12, 31
<b>0.75</b>	13, 34	-

The results indicate that when the rate of type 1 patients decreases, it is beneficial to increase the number of triage nurses instead of doctors. On the other hand, when type 1 patients increase, the numbers of both triage nurses and doctors should be increased.

## 6. CONCLUSIONS

A&E's are the first point of contact for urgent and complex cases. The performance targets of A&E's have long been to reduce the patient waiting times below 4 hours. As the hospitals have failed to satisfy this target, the government is planning to adopt alternative policies such as considering the waiting time targets for only serious cases. The staff planning in A&E's affect the performance levels critically. Due to the uncertainties involved in the A&E services, finding optimum capacities satisfying the performance targets is difficult with classical methods. This paper proposes to use a robust optimization based approximation for the computation of the worst-case waiting times in an A&E where patient are first triaged and then prioritized based on the urgency. We first model the problem with the approximation and then show that the model can be solved to optimality with the commercial solvers. We also develop an SO based heuristic where we use the approximation for the maximum waiting time in the search for a better feasible solution.

The computational experiments show that the AO outperforms the SO heuristic in terms of the computation time and the solution performance. The advantage of both the AO and the SO heuristic is their speed to provide an approximately good solution very quickly. As the problem complexity is increased with the waiting time limits for both patient types, the

performance of the solutions obtained by the approximation drops when both patients are subject to waiting time limit. The experiments also indicate that the approximation approach still works well for different distribution assumptions for the treatment and triage durations. Another observation drawn from the experiment results is the non-linear effect of the waiting time limits on the solutions. The future studies may investigate the suitability of the approximation method for more complex A&E operations including more than two prioritization categories, the diagnostic test queues, etc.

## REFERENCES

- [1] Ahmed, M.A., Alkhamis, T.M.: Simulation optimization for an emergency department healthcare unit in Kuwait. *Eur. J. Oper. Res.* **198**(3), 936–942 (2009)
- [2] Alfonso, E., Xie, X., Augusto, V., Garraud, O.: Modelling and simulation of blood collection systems: improvement of human resources allocation for better cost-effectiveness and reduction of candidate donor abandonment. *Vox Sang.* **104**(3), 225–233 (2013)
- [3] Allen, A.O.: Probability, statistics, and queueing theory. Academic press (2014)
- [4] Alrefaei, M.H., Diabat, A.: Modelling and optimization of outpatient appointment scheduling. *RAIRO-Oper. Res.* **49**(3), 435–450 (2015). URL <https://doi.org/10.1051/ro/2014041>
- [5] Asaduzzaman, M., Chausalet, T.J., Robertson, N.J.: A loss network model with overflow for capacity planning of a neonatal unit. *Ann. Oper. Res.* **178**(1), 67–76 (2010)
- [6] Audit General for Scotland: Emergency Departments. Tech. rep., Audit General for Scotland (2010). URL [http://www.audit-scotland.gov.uk/uploads/docs/report/2010/nr{\\\_}100812{\\\_}emergency{\\\_}departments.pdf](http://www.audit-scotland.gov.uk/uploads/docs/report/2010/nr{\_}100812{\_}emergency{\_}departments.pdf)
- [7] Bandi, C., Bertsimas, D.: Tractable stochastic analysis in high dimensions via robust optimization. *Math. Program.* **134**, 23–70 (2012). DOI 10.1007/s10107-012-0567-2
- [8] Bonami, P., Kilinç, M., Linderoth, J.: Algorithms and software for convex mixed integer nonlinear programs. In: *Mix. integer nonlinear Program.*, pp. 1–39. Springer, New York (2012)
- [9] Boyd, S., Vandenberghe, L.: *Convex Optimization*, vol. 25. Cambridge University Press (2004). DOI 10.1017/CBO9780511804441. URL <http://ebooks.cambridge.org/ref/id/CBO9780511804441{\%}5Cnhttp://www.informaworld.com/openurl?genre=article{\&}doi=10.1080/10556781003625177{\&}magic=crossref>
- [10] Bretthauer, K.M., Heese, H.S., Pun, H., Coe, E.: Blocking in healthcare operations: A new heuristic and an application. *Prod. Oper. Manag.* **20**(3), 375–391 (2011)
- [11] Castillo, I., Ingolfsson, A., Sim, T.: Social optimal location of facilities with fixed servers, stochastic demand, and congestion. *Prod. Oper. Manag.* **18**(6), 721–736 (2009)
- [12] Chen, T.L., Wang, C.C.: Multi-objective simulation optimization for medical capacity allocation in emergency department. *J. Simul.* **10**(1), 50–68 (2016)
- [13] Cochran, J.K., Roche, K.T.: A multi-class queueing network analysis methodology for improving hospital emergency department performance. *Comput. Oper. Res.* **36**(5), 1497–1512 (2009)

- [14] Computational Infrastructure for Operations Research: Bonmin (2018). URL <https://www.coin-or.org/Bonmin/index.html>
- [15] Costa, A.X., Ridley, S.A., Shahani, A.K., Harper, P.R., De Senna, V., Nielsen, M.S.: Mathematical modelling and simulation for planning critical care capacity. *Anaesthesia* **58**(4), 320–327 (2003)
- [16] Creemers, S., Lambrecht, M.: An advanced queueing model to analyze appointment-driven service systems. *Comput. Oper. Res.* **36**(10), 2773–2785 (2009)
- [17] De Angelis, V., Felici, G., Impelluso, P.: Integrating simulation and optimisation in health care centre management. *Eur. J. Oper. Res.* **150**(1), 101–114 (2003)
- [18] Eskandari, H., Riyahifard, M., Khosravi, S., Geiger, C.D.: Improving the emergency department performance using simulation and MCDM methods. In: *Proc. winter Simul. Conf.*, pp. 1211–1222. Winter Simulation Conference (2011)
- [19] Flessa, S.: Where efficiency saves lives: A linear programme for the optimal allocation of health care resources in developing countries. *Health Care Manag. Sci.* **3**(3), 249–267 (2000)
- [20] Fletcher, A., Halsall, D., Huxham, S., Worthington, D.: The DH Accident and Emergency Department model: a national generic model used locally. *J. Oper. Res. Soc.* **58**(12), 1554–1562 (2007). DOI 10.1057/palgrave.jors.2602344. URL <https://doi.org/10.1057/palgrave.jors.2602344>
- [21] Fomundam, S., Herrmann, J.W.: A survey of queueing theory applications in healthcare (2007)
- [22] Fruggiero, F., Lambiase, A., Fallon, D.: Computer simulation and swarm intelligence organisation into an emergency department: a balancing approach across ant colony optimisation. *Int. J. Serv. Oper. Informatics* **3**(2), 142–161 (2008)
- [23] Ghanes, K., Wargon, M., Jouini, O., Jemai, Z., Diakogiannis, A., Hellmann, R., Thomas, V., Koole, G.: Simulation-based optimization of staffing levels in an emergency department. *Simulation* **91**(10), 942–953 (2015)
- [24] Gourgand, M., Grangeon, N., Klement, N.: Activities planning and resources assignment on distinct places: a mathematical model. *RAIRO-Oper. Res.* **49**(1), 79–98 (2015). URL <https://doi.org/10.1051/ro/2014028>
- [25] Govind, R., Chatterjee, R., Mittal, V.: Timely access to health care: Customer-focused resource allocation in a hospital network. *Int. J. Res. Mark.* **25**(4), 294–300 (2008)
- [26] Green, L.V., Kolesar, P.J., Whitt, W.: Coping with Time-Varying Demand When Setting Staffing Requirements for a Service System. *Prod. Oper. Manag.* **16**(1), 13–39 (2007)
- [27] Gul, S., Denton, B.T., Fowler, J.W.: A Progressive Hedging Approach for Surgery Planning Under Uncertainty. *INFORMS J. Comput.* **27**(4), 755–772 (2015)
- [28] Güneş, E.D., Yaman, H.: Health network mergers and hospital re-planning. *J. Oper. Res. Soc.* **61**(2), 275–283 (2010)
- [29] Gupta, V., Osogami, T.: On Markov–Krein characterization of the mean waiting time in M/G/K and other queueing systems. *Queueing Syst.* **68**(3–4), 339–352 (2011)
- [30] Harper, P.R., Powell, N.H., Williams, J.E.: Modelling the size and skill-mix of hospital nursing teams. *J. Oper. Res. Soc.* **61**(5), 768–779 (2010)

- [31] Hu, X., Barnes, S., Golden, B.: Applying queueing theory to the study of emergency department operations: a survey and a discussion of comparable simulation studies. *Int. Trans. Oper. Res.* **25**(1), 7–49 (2018)
- [32] Hulshof, P.J.H., Mes, M.R.K., Boucherie, R.J., Hans, E.W.: Tactical planning in healthcare using approximate dynamic programming. Tech. rep., University of Twente (2013)
- [33] Ibrahim, I.M., Liong, C.Y., Bakar, S.A., Ahmad, N., Najmuddin, A.F.: Estimating Optimal Resource Capacities in Emergency Department. *Indian J. Public Heal. Res. Dev.* **9**(11) (2018)
- [34] Izady, N., Worthington, D.: Setting staffing requirements for time dependent queueing networks: The case of accident and emergency departments. *Eur. J. Oper. Res.* **219**(3), 531–540 (2012)
- [35] Köllerström, J.: Heavy traffic theory for queues with several servers. I. *J. Appl. Probab.* **11**(3), 544–552 (1974)
- [36] Lakshmi, C., Iyer, S.A.: Application of queueing theory in health care: A literature review. *Oper. Res. Heal. Care* **2**(1), 25–39 (2013)
- [37] Lunenfeld, B., Stratton, P.: The clinical consequences of an ageing world and preventive strategies. *Best Pract. Res. Clin. Obstet. Gynaecol.* **27**(5), 643–659 (2013)
- [38] Mahar, S., Bretthauer, K.M., Salzarulo, P.A.: Locating specialized service capacity in a multi-hospital network. *Eur. J. Oper. Res.* **212**(3), 596–605 (2011)
- [39] Mayhew, L., Smith, D.: Using queueing theory to analyse the government’s 4-h completion time target in accident and emergency departments. *Health Care Manag. Sci.* **11**(1), 11–21 (2008)
- [40] Mingzhu, Z., Ershi, Q.: A Multi-Type Queuing Network Analysis Method for Controlling Server Number in the Outpatient. *Open Autom. Control Syst. J.* **8**(1) (2016)
- [41] Mohiuddin, S., Busby, J., Savović, J., Richards, A., Northstone, K., Hollingworth, W., Donovan, J.L., Vasilakis, C.: Patient flow within UK emergency departments: a systematic review of the use of computer simulation modelling methods. *BMJ Open* **7**(5), e015007 (2017)
- [42] Mortimore, A., Cooper, S.: The "4-hour target": emergency nurses’ views. *Emerg. Med. J.* **24**(6), 402–404 (2007). DOI 10.1136/emj.2006.044933. URL <https://www.ncbi.nlm.nih.gov/pubmed/17513535><https://www.ncbi.nlm.nih.gov/pmc/PMC2658273/>
- [43] Munir, W.: Critical analysis of the 4-hour A&E policy’s impact on elderly patients. *Br. J. Nurs.* **17**(18) (2008)
- [44] NHS England: AE Waiting times and activity. Tech. rep., NHS England, London (2018). URL <https://www.england.nhs.uk/statistics/statistical-work-areas/ae-waiting-times-and-activity/ae-attendances-and-emergency-admissions-2018-19/>
- [45] NHS England: Commissioning Committee Report to Board. Tech. rep., NHS England (2018). URL <https://www.england.nhs.uk/wp-content/uploads/2018/11/13i-pb-28-11-18-commissioning-committee-report-to-board-24-october.pdf>
- [46] NHS UK: Hospital Accident and Emergency Activity, 2017-18, Provider Level Analysis. Tech. rep., National Health Services UK (2018). URL <https://digital.nhs.uk/data-and-information/publications/statistical/>

hospital-accident--emergency-activity/2017-18

- [47] NHS UK: Pay for doctors (2019). URL <https://www.healthcareers.nhs.uk/explore-roles/doctors/pay-doctors>
- [48] NHS UK: Urgent and Emergency Care: When to go to A&E (2019). URL <https://www.nhs.uk/using-the-nhs/nhs-services/urgent-and-emergency-care/when-to-go-to-ae/>
- [49] OptTek Systems: OptQuest (2019). URL <https://www.opttek.com/products/optquest/>
- [50] Pehlivan, C., Augusto, V., Xie, X., Crenn-Hebert, C.: Multi-period capacity planning for maternity facilities in a perinatal network: A queuing and optimization approach. In: Autom. Sci. Eng. (CASE), 2012 IEEE Int. Conf., pp. 137–142. IEEE (2012)
- [51] Rico, F., Salari, E., Centeno, G.: Emergency departments nurse allocation to face a pandemic influenza outbreak. In: Proc. 39th Conf. Winter Simul., pp. 1292–1298. IEEE Press (2007)
- [52] Royal College of Nursing: NHS pay scales 2017-18 (2019). URL <https://www.rcn.org.uk/employment-and-pay/nhs-pay-scales-2017-18>
- [53] Saghaian, S., Austin, G., Traub, S.J.: Operations research/management contributions to emergency department patient flow optimization: Review and research prospects. IIE Trans. Healthc. Syst. Eng. **5**(2), 101–123 (2015)
- [54] Santibáñez, P., Bekiou, G., Yip, K.: Fraser Health uses mathematical programming to plan its inpatient hospital network. Interfaces (Providence). **39**(3), 196–208 (2009)
- [55] Stummer, C., Doerner, K., Focke, A., Heidenberger, K.: Determining location and size of medical departments in a hospital network: a multiobjective decision support approach. Health Care Manag. Sci. **7**(1), 63–71 (2004)
- [56] Syam, S.S., Côté, M.J.: A location allocation model for service providers with application to not-for-profit health care organizations. Omega **38**(3), 157–166 (2010)
- [57] The Telegraph: One in three hospital patients admitted as an emergency to be sent home without an overnight stay (2019). URL <https://www.telegraph.co.uk/news/2019/03/08/one-three-hospital-patients-admitted-emergency-sent-home-without/>
- [58] Tijms, H.C., Van Hoorn, M.H., Federgruen, A.: Approximations for the steady-state probabilities in the M/G/c queue. Adv. Appl. Probab. pp. 186–206 (1981)
- [59] Weng, S.J., Cheng, B.C., Kwong, S.T., Wang, L.M., Chang, C.Y.: Simulation optimization for emergency department resources allocation. In: Proc. 2011 Winter Simul. Conf. (WSC),, pp. 1231–1238. IEEE (2011)
- [60] Wiler, J.L., Griffey, R.T., Olsen, T.: Review of modeling approaches for emergency department patient flow and crowding research. Acad. Emerg. Med. **18**(12), 1371–1379 (2011)