# Modelling and Solving Healthcare Decision Making Problems under Uncertainty

by

**Elvan Gökalp**

A thesis submitted in partial fulfilment of the requirements for the

Degree of Doctor of Philosophy

Operational Research and Management Sciences

Warwick Business School

The University of Warwick

September 2017

# Contents

# List of Figures

# List of Tables

# ACKNOWLEDGEMENTS

# DECLARATIONS

I declare that this thesis is my own work. This thesis has not been submitted in any other university and not published in any scientific journal.

# ABSTRACT

The efficient management of healthcare services is a great challenge for healthcare managers because of ageing populations, rising healthcare costs, and complex operation and service delivery systems. The challenge is intensified due to the fact that healthcare systems involve various uncertainties. Operations Research (OR) can be used to model and solve several healthcare decision making problems at strategic, tactical and also operational levels. Among different stages of healthcare decision making, resoure allocation and capacity planning play an important role for the overall performance of the complex systems. This thesis aims to develop modelling and solution tools to support healthcare decision making process within dynamic and stochastic systems. In particular, we are concerned with stochastic optimization problems, namely i) capacity planning in a stem-cell donation network, ii) resource allocation in a healthcare outsourcing network and iii) real-time surgery planning. The patient waiting times and operational costs are considered as the main performance indicators in these healthcare settings. The uncertainties arising in patient arrivals and service durations are integrated into the decision making as the most significant factors affecting the overall performance of the underlying healthcare systems. We use stochastic programming, a collection of OR tools for decision-making under uncertainty, to obtain robust solutions against these uncertainties. Due to complexities of the underlying stochastic optimization models such as large real-life problem instances and non-convexity, these models cannot be solved efficiently by exact methods within reasonable computation time. Thus, we employ approximate solution approaches to obtain feasible decisions close to the optimum. The computational experiments are designed to illustrate the performance of the proposed approximate methods. Moreover, we analyze the numerical results to provide some managerial insights to aid the decision-making processes. The numerical results show the benefits of integrating the uncertainty into decision making process and the impact of various factors in the overall performance of the healthcare systems.

# Chapter 1

# Introduction

Healthcare industry is one of the largest and crucial sectors affecting millions of lives worldwide. Recently, the industry has been facing significant challenges due to several sociological and technological changes. First of all, the increasing amount of publicly available data has resulted in a higher demand for better quality of services. Secondly, the incidence rates of long-term conditions such as hypertension or diabetes have increased because of the modern lifestyle. Lastly, one of the biggest challenges for the industry is the aging phenomenon i.e. the increasing rate of over-aged population. The emergence of these challenges put healthcare managers under a serious pressure to improve the efficiencies of their services.

Healthcare management is a complex task due to several distinguishing features of healthcare services. First of all, they serve a large number of patients and engage with multiple stakeholders, such as hospital managers, doctors, and nurses. Often in time, these stakeholders have conflicting objectives like maximising profit or minimising waiting time. From a managerial point of view, these performance measures are mainly influenced by several tactical and strategic decisions regarding the capacity planning or resource allocation. These

decisions should take into account inflexible and expensive resources as well as operational complexities involving different specialities and resources at the same time. Besides, healthcare services are subject to several uncertainties such as patient arrivals, service durations, treatment outcomes, test results, or disease progression.

Mathematical modelling can be used as a tool to simplify the complex systems and analyze them in a more effective manner. Operations Research (OR) provides useful modelling tools for healthcare management due to its success in handling large and complex systems (Brandeau et al., 2004). Healthcare modelling is specifically concerned with the design of healthcare delivery systems to achieve cost-effective quality of services. In particular, the decision-making in healthcare modelling may be strategic/tactical such as capacity planning and resource allocation or operational such as patient and staff scheduling. Capacity planning, as one of the main interest areas in OR, deals with an effective use of available resources to meet the changing demand for products or services. An effective capacity planning model requires to take optimal decisions to minimize operational costs while satisfying the demand even at emergency situations. In particular, the capacity planning for healthcare facilities such as intensive care units (Gallivan et al., 2002; Harper et al., 2010), inpatient clinics (Gnanlet and Gilland, 2009; Creemers and Lambrecht, 2009), and hospitals (Utley et al., 2003) is crucial to utilise resources such as nurses, beds and operating rooms. Similarly, healthcare resource allocation problems aim to allocate a given set of resources among the operational entities such that the overall service performance is improved.

Healthcare modelling should also consider the inherent uncertainties, as they may have a significant impact on the solution and the quality of the service provided to patients. The uncertainties affecting healthcare processes may

be grouped under two categories: medical and managerial. Medical uncertainties can be counted as treatment outcomes, results of medical tests, disease progression, etc. These uncertainties are usually independent from a specific hospital, region or management which implies a higher chance of finding relevant data for the analysis. Managerial uncertainties can be listed as the variations in the demand and service times, availability of resources (especially the medical staff), business environment, and the emergence of new technologies. Unlike the medical ones, these uncertainties may be specific to the hospital or the country under consideration and may require different modelling approaches.

Two critical managerial uncertainties affecting healthcare operations are the variations in patient arrivals and service durations. In most of the healthcare services, the overall patient demand is not known with certainty. When a healthcare delivery system is not designed according to this variation, the resulting chaotic environment puts the lives of patients in danger. Appointment systems aim to reduce the impact of this variation. However, even with an appointment system, the arrival of an emergent patient is unavoidable. These patients should either be diverted or inserted into the existing list of admitted patients which implies extra waiting times for the existing patients. Thus, both the strategic and operational healthcare planning should take into account the uncertainty in patient demand. The other critical uncertainty, the variations in service durations, can be managed by creating robust schedules with plenty of buffer times. However, this may result in an inefficient usage of resources which are very limited and costly. Therefore, the trade-off between the conservative (robust) and cost-effective approaches should be balanced very carefully in healthcare decision-making.

Incorporating the uncertainty into the modelling of a problem is expected to increase the robustness of the results. OR provides different modelling ap-

proaches depending on the nature of the decision-making problem under consideration. When the probability distributions of the uncertain data are known or can be estimated, stochastic programming can be used to model the problem. Stochastic programming is a collection of the OR tools used for the decision-making problems under uncertainty. A stochastic programming model aims to find the optimum policy that is feasible for all (or almost all) possible realizations of the uncertainty. In other words, it is concerned with the expected performance of a system involving uncertainties. Stochastic programming offers several modelling tools suitable for different types of uncertainties and decision maker attitudes towards these uncertainties. For example, uncertainties may follow a specific distribution or attain no information at all. Some uncertainties may be resolved after an initial set of decisions is taken. Besides, decision-makers may be risk-averse, very cautious against uncertainties, or risk-seeker, willing to take risks. Next section provides more detailed information about different stochastic programming tools suitable for various kinds of uncertainties and risk attitudes.

This thesis focuses on the capacity planning and resource allocation problems arising in different healthcare management practices. Specifically, we model and solve three healthcare decision making problems under uncertainty that can be listed as:

- Capacity planning for a network of stem-cell donation centres,

- Resource allocation for a healthcare network with outsourcing,

- Real-time surgery management in a surgery suite.

The chapters incorporate several common features. First of all, the uncertainties in patient arrivals and service times are considered in all chapters. Secondly, the minimization of the patient waiting time is the main concern in all models developed in the chapters. Lastly, the underlying uncertainties in the

problems are incorporated into the modelling by using several stochastic programming tools.

## 1.1   Thesis Outline

This thesis is composed of five chapters. Chapter 1 provides a brief introduction, an overview of the thesis and a review of the OR methodologies for decision-making under uncertainty. Chapter 2 focuses on the capacity planning for a network of stem-cell donation centres. The chapter starts with providing some background information related to stem-cell donation centres. Then, the underlying capacity planning problem is explained in more detail and a mathematical framework is proposed. Afterwards, we present a scenario-based stochastic programming model where the maximum patient waiting time is approximated by a robust optimization based approach. Finally, we design several computational experiments to investigate the model sensitivity and the impact of different network structures on the overall service performance.

In Chapter 3, we focus on a resource allocation problem in a healthcare network with outsourcing. First, we provide a literature review and then describe the underlying problem in detail. We develop a non-linear integer programming model by incorporating a robust queuing approach. The structural properties of the model are presented and an alternating optimization based heuristic is proposed to solve the model. The chapter concludes with the computational experiments that compare the performances of the proposed heuristic and the available commercial solvers and investigate the effect of the model parameters on the overall service measures.

Chapter 4 presents a stochastic dynamic programming approach for the real-time management of a surgery schedule. First, we provide a review of the

related literature and describe the problem in detail. Then, a stochastic dynamic programming model is presented. Due to the large problem size, the real-sized instances are solved with an approximate dynamic programming (ADP) algorithm. The computational experiments analyze the impact of the model parameters on the algorithm's performance and evaluate it with respect to an exact method and a myopic heuristic. Finally, we compare different elective surgery scheduling strategies in terms of the overall cost. Chapter 5 concludes the thesis by summarizing the research, the key findings and the direction of the future research.

## 1.2   Contributions of Thesis

This thesis aims to contribute to the literature by:

- Modelling three healthcare decision-making problems under uncertainty,

- Using several stochastic programming tools to deal with different types of uncertainties emerging within these problems,

- Analyzing the structures of the resulting models and developing and implementing appropriate solution methods,

- Designing and analyzing several computational experiments to investigate the performance of the solution methods and produce a set of useful managerial insights.

Each chapter provides more detailed information regarding the contributions of the respective study. In the next section, we present several modelling and solution approaches for decision-making under uncertainty.

## 1.3 Decision-making under Uncertainty: A Review

Healthcare management problems are affected by several uncertainties that should be taken into account to obtain robust solutions. These uncertainties may be due to the measurement errors or simply because the relevant data are not realized yet. Traditional (deterministic) optimization models assume that the input data of a problem are known beforehand. However, input data are usually neither available nor fixed. This section provides a review of various methods used to model and solve decision-making problems under uncertainty.

The most intuitive approach to handle the uncertainties in a decision-making problem is to replace them with their average values or point-wise estimates. This method reduces the stochastic problem into a deterministic one. But, Ben-Tal and Nemirovski (1999, 2000) have shown that the solutions obtained by using these estimates may become infeasible even with the slight changes in the levels of the uncertain parameters. Thus, the solution found with the expected values may not be feasible if the data change or realize differently from the expectation. This is usually not acceptable to decision makers who need solutions suitable for most of the future realizations. Therefore, as many realizations as possible should be included in the model. However, this may result in very large problems that are computationally expensive.

Two main OR approaches for decision-making under uncertainty are stochastic programming and robust optimization. The following sections provide some background information related to these approaches.

## 1.3.1 Stochastic Programming

Stochastic programming is a wide framework for modelling and solving optimization problems in the presence of uncertainty. The framework is mainly introduced in 1955 by Dantzig using the fact that uncertain data can be described by probability distributions. Readers are directed to Birge and Louveaux (2011) for a comprehensive review of stochastic programming. In this section, we present different modelling and solution approaches under stochastic programming.

**Modelling Approaches**

This section provides an overview of stochastic programming modelling approaches by emphasizing several problem features such as uncertainties, number of objective functions, and convexity. First, let's introduce the basic concepts of stochastic programming. A stochastic programming model can be defined as,

$$\max_{\mathbf{x} \in \mathbf{X}} E[f(\mathbf{x}, \boldsymbol{\xi})], \tag{1.1}$$

where $\mathbf{x} \in \mathbf{X}$ and $\boldsymbol{\xi} \in \boldsymbol{\Xi}$ denote the vectors of decision variables and random variables belonging to a feasible set $\mathbf{X}$ and a probability space $\boldsymbol{\Xi}$, respectively. The objective of the model is to maximize the expected performance represented in the form of function, $f(\mathbf{x}, \boldsymbol{\xi})$, under the presence of uncertainties $\boldsymbol{\xi}$. The stochastic programming framework assumes that the distributions of the random variables are known or can be estimated from the historical data.

**Multi-stage Modelling:** The uncertainties in a stochastic problem may be realized in different points of the planning horizon. Some decisions have to be taken before the uncertainties are realized while the others may be taken afterwards. The timing of a set of decisions is called a stage. Between each stage, some relevant information is unfolded i.e. uncertainties are realized.

A single-stage stochastic programming model requires all decisions to be taken before the uncertainties are realized. If some decisions may be taken afterwards, then the problem can be modelled as a multi-stage stochastic programming model. A two-stage stochastic programming problem can be written as,

$$\max_{\mathbf{x} \in \mathbf{X}} \quad f(\mathbf{x}) + E[Q(\mathbf{x}, \boldsymbol{\xi})], \tag{1.2}$$

where $Q(\mathbf{x}, \boldsymbol{\xi})$ is the optimal value of the second stage problem,

$$\max_{\mathbf{y} \in \mathbf{Y}} \quad \{q(\mathbf{y}, \boldsymbol{\xi}) \mid T(\boldsymbol{\xi})\mathbf{x} + W(\boldsymbol{\xi})\mathbf{y} = h(\boldsymbol{\xi})\},$$

and $\mathbf{x} \in \mathbf{X}$, $\mathbf{y} \in \mathbf{Y}$, $T$, $W$ and $h$ represent the first and second stage variables, where $\mathbf{Y}$ is the feasible set $y$ belongs to, and the elements of the second-stage problem, respectively. In a multi-stage model, the decisions taken in the second or later stages are called recourse actions. In other words, each possible realization of the uncertainty is associated with a recourse action.

**Chance-constrained Formulations:** In some problems, constraint violation cannot be avoided due to the inherent uncertainties. For these cases, the probability of violating a constraint can be bounded in a chance-constraint, which can be formulated as,

$$\Pr\left(\phi(\mathbf{x}, \boldsymbol{\xi}) \geq 0\right) \leq p,$$

where $\phi(\mathbf{x}, \boldsymbol{\xi}) \geq 0$, $\Pr(.)$ and $p \in [0, 1]$ represent a finite system of inequalities, the probability function and the threshold probability level, respectively. A chance-constraint can be for a single constraint only or multiple constraints at the same time in which the dependencies between the random

variables should also be considered.

**Scenario-based Modelling:** In a stochastic programming model, probability distribution functions can be directly used as in the form of chance-constrained formulations. Alternatively, they can be discretized in the form of scenarios; the random vector $\boldsymbol{\xi}$ can be replaced with its $K$ possible realizations (scenarios), $\xi_1, \cdots, \xi_K$, with the respective probabilities of occurrence represented as $p_1, \cdots, p_K$. Then, model (1.1) can be written as,

$$\max_{\mathbf{x} \in \mathbf{X}} \quad \sum_{k=1}^{K} p_k f(\mathbf{x}, \xi_k). \tag{1.3}$$

With this discretization, the stochastic programming model (1.1) reduces to a deterministic equivalent (1.3). As the number of scenarios, $K$, increases, the solution of (1.3) is expected to approach to the (exact) optimum of model (1.1).

A significant challenge in scenario-based modelling is to generate the scenarios in such a way that the uncertainty representation is rigorous enough. Several scenario generation methods have been proposed in the literature (see Kaut and Wallace (2007) for more information). The choice of an appropriate scenario generation method depends on the problem features such as the number of stages or the nature of the information revealed between the stages. The most popular scenario generation method is Monte-Carlo simulation (Chen, 2015) in which, first, a random and independent sequence of numbers, $U_1, \cdots, U_K$ is generated from a uniform distribution over [0,1]. Then, by using an appropriate transformation, these random numbers are converted into a sample of $\boldsymbol{\xi}$: $\boldsymbol{\xi}' = \{\xi_1, \cdots, \xi_K\}$. In other words, the sequence $\omega = \{U_1, \cdots, U_K\}$ is an element of the probability space, while the generated sample $\boldsymbol{\xi}'$ is a function of $\omega$. Given this sample, the expectation

in (1.2), $E[Q(\mathbf{x}, \boldsymbol{\xi})]$, can be approximated as,

$$E[Q(\mathbf{x}, \boldsymbol{\xi})] = \sum_{k=1}^{K} Q(\mathbf{x}, \xi_k)/K. \tag{1.4}$$

This type of uncertainty modelling is also called Sample Average Approximation (SAA). More information regarding SAA can be found in Shapiro (2013).

**Objective:** An important feature of a stochastic programming model is the number of objective functions. When a model contains more than one objective function, $f_i(\mathbf{x}, \boldsymbol{\xi})$ for $i = 1, \cdots, m$, they can be combined by multiplying each objective function with a weight $w_i$, and summing them up: $h(\mathbf{x}, \boldsymbol{\xi}) = \sum_{i=1}^{m} w_i f_i(\mathbf{x}, \boldsymbol{\xi})$. The weights can be elicited from the decision-makers by using the weight elicitation techniques (see Riabacke et al. (2012) for more information). Alternatively, the problem may be formulated as a multi-objective stochastic programming model which can be stated as,

$$\max_{\mathbf{x} \in \mathbf{X}} \quad h(\mathbf{x}, \boldsymbol{\xi}) = \big(f_1(\mathbf{x}, \boldsymbol{\xi}), f_2(\mathbf{x}, \boldsymbol{\xi}), \cdots, f_m(\mathbf{x}, \boldsymbol{\xi})\big),$$

and can be solved with one of the multi-objective solution methods such as evolutionary algorithms (Deb, 2001). These methods try to find non-dominated solutions. A solution is called as non-dominated if none of the objective functions can be improved in value without worsening some of the other objective values. Readers are referred to Stancu-Minasian (1984) and Ben-Abdelaziz (2012) for detailed information regarding the stochastic programming with multiple objectives.

**Convexity:** Other than the modelling of objective functions, the appropriate solution method for a stochastic programming model depends on whether it is

convex or not. Convex models can be solved with exact methods while non-convex models usually need to be solved with approximate or heuristic methods. The convexity of a stochastic programming model is mainly influenced by the representation of the uncertainty. For example, chance-constrained formulations are usually non-convex due to the underlying probability functions. In multi-stage models, non-convexity may especially arise in recourse functions which can be integer, non-convex and discontinuous (Sahinidis, 2004). Next section provides the details of the solution approaches used for stochastic programming models.

## Solution Approaches

This section reviews possible solution methods for stochastic programming models. We categorize the solution methods into two: analytical, which provide the global optimum, and non-analytical (approximate or heuristic), preferred when the analytical methods are not applicable or computationally expensive.

**Analytical Approaches:** Depending on the structure of a model, there are various analytical solution methodologies available in the literature. For example, if the objective function of a model is differentiable and there is no constraint, then the most intuitive solution is to take the derivative of the objective function. However, usually, stochastic programming models have constraints. In such a case, when the objective function and the constraints are linear, the model can be solved with a linear programming method such as simplex. Yet, the convexity is lost when a stochastic programming model contains integer elements as in most of the scenario-based formulations. These models can be solved with branch-and-bound, an exact method used to solve integer linear programming models. In this method, the candidate solutions are systematically enumerated by using a rooted tree with the full

set at the root. Then, the branches of this tree are explored sequentially. Before going into the other nodes in a branch, the root solution is compared with the estimated upper and lower bounds of the optimal solution and discarded if it does not provide a better solution. The main drawback of branch-and-bound is the need for a large memory to keep track of the solution tree.

Another exact method that can be used for stochastic programming models with a discrete solution space is the exhaustive search: evaluating and comparing all possible solutions. However, due to the computational concerns, this method is only preferred when the feasible solution set is small.

**Approximate Approaches and Heuristics:** If a problem is analytically intractable or computationally expensive to solve, then approximate or heuristic methods can be used to obtain a solution. Approximate solution methods find a solution close enough to the optimum within a reasonable time. On the other hand, heuristics provide any satisfactory solution within a short time. Heuristics can be preferred over approximate methods when obtaining a solution quickly is more important than the quality of that solution.

Stochastic programming models are usually hard to solve and require approximate or heuristic approaches (Stougie and van Der Vlerk, 2003). An overview of approximation algorithms for stochastic programming models and the analysis of their performance are given in Stougie and van Der Vlerk (2003). A popular approximate method used to solve stochastic programming models is Langrangian relaxation. This method involves adding penalty costs to the objective function due to the violations in the inequality constraints. Another approximate solution method especially used for two-stage stochastic programming models is L-shaped decomposition algo-

rithm. Instead of solving the original problem, this method decomposes it into smaller problems and solves them sequentially. There are several variants of L-shaped decomposition such as regularized L-shaped decomposition. The main framework can also be extended to multi-stage problems (Birge, 1985) which result in the algorithms like nested Benders decomposition, stochastic dual dynamic programming, etc. The decomposition approach can also be combined with branch-and-bound to generate an efficient exact algorithm for two-stage stochastic programming models (Ahmed et al., 2004). When a stochastic problem is too complex to build an optimization model, simulation-optimization can be used to obtain a solution. In this method, a simulation model is used as a black-box to map the decision variables to an estimate of the performance measure (April et al., 2003).

For chance-constrained formulations, an approximate solution method is p-efficient point-based algorithm that enumerates p-efficient points of the joint probability function in the chance-constraint. A point $v \in R^n$ is called a p-efficient point of the probability function $F$, if $F(v) \geq p$ and there is no $y < v$ such that $F(y) \geq p$, where $y \in \mathbb{R}^n$. (Lejeune and Noyan, 2010).

If the main concern is obtaining any solution but not nexessarily the best one, then a stochastic programming model can be solved with a local search algorithm that moves to a (better) neighbour solution iteratively. The algorithm stops when the time limit is achieved or a deemed optimal solution is found. Another possible heuristic for stochastic programming models is genetic algorithm (Ma and Zhang, 2002) which iteratively modifies a population of candidate solutions based on the natural selection idea. Multi-objective tabu search, a heuristic especially used for multi-objective problems, improves an initial solution by searching for optimal solutions

using parallel agents (Erdogan et al., 2010). Each agent searches for the non-dominated solutions and shares the information with other agents to get a better search performance. Another heuristic, progressive hedging, is based on the scenario aggregation idea and used for multi-stage stochastic programming models with integer variables (Wallace and Helgason, 1991).

**Markov Decision Process**

Within stochastic programming, a special area of study, Markov Decision Process (MDP), focuses on the modelling of stochastic problems with dynamic decision-making. Specifically, MDP is a framework for modelling Markovian processes where the outcomes are affected by decisions and uncertainty at the same time. In an MDP formulation, the process is defined by its selected features, called states, that capture all information required to make decisions. Each possible state $s$ and action $a$ are assumed to belong to finite sets $S$ and $A$, respectively. The planning horizon can be finite or infinite and the decisions (actions) are taken at discrete time points (epochs) represented with $t \in \{1, \cdots, T\}$ ($T = \infty$ for the infinite case). When the decisions can be taken at any time point, the problem becomes a continuous-time MDP. At each decision epoch $t$, the system randomly moves from state $s$ to another state $s'$ partially affected by the selected action $a$ and the transition probabilities represented by $P_a^t(s, s')$. The reward obtained as a result of this transition is denoted by $R_a^t(s, s')$. The Markovian property implies that the sets of available actions, rewards and transition probabilities at epoch $t$ only depend on the current state and action at this epoch, not the past ones. The objective of an MDP formulation is to find the optimum policy, the optimum action for each state and epoch maximizing total expected reward. If the states cannot be observed with certainty, the problem can be modelled as a partially observable Markov decision process (POMDP) (Dutech and Scherrer, 2013).

MDP formulations can be solved with dynamic programming methods such as value or policy iteration. The value iteration method, known as backward induction as well, iteratively calculates the optimum value of being in state $s \in S$, $V(s)$, by using an optimality equation. For an infinite horizon problem, the optimality equation can be written as

$$V_{(i+1)}(s) := \max_{a \in A} \left\{ \sum_{s' \in S'} P_a(s, s')\big(R_a(s, s') + \gamma V_i(s')\big) \right\},$$

where $i$, $\gamma$, and $S'$ represent the iteration counter, the discount factor and the set of all possible next states, respectively. In a finite horizon problem, the value iteration method requires to compute the values of all possible states at the end of the planning horizon, $V_T(s_T)$ for all $s_T \in S$. Then, the state values in the previous time periods are calculated iteratively by moving backwards in time based on the optimality equation:

$$V_t(s_t) := \max_{a_t \in A} \left\{ \sum_{s'_{t+1} \in S'} P_{a_t}(s_t, s'_{t+1})\big(R_{a_t}(s_t, s'_{t+1}) + V_{t+1}(s'_{t+1})\big) \right\}, \quad t = 1, \cdots, T-1.$$

In the policy iteration, instead of iterating the value function, the optimum policy $\pi(s)$ is iteratively computed by using,

$$\pi(s) := \arg\max_{a \in A} \left\{ \sum_{s'} P_a(s, s')\big(R_a(s, s') + \gamma V(s')\big) \right\},$$

for the infinite horizon case, and

$$\pi(s_t) := \arg\max_{a_t \in A_t} \left\{ \sum_{s'_{t+1} \in S'} P_{a_t}(s_t, s'_{t+1})\big(R_{a_t}(s_t, s'_{t+1}) + V_{t+1}(s'_{t+1})\big) \right\}, \quad t = 1, \cdots, T-1,$$

for the finite horizon case. Then, the value function is calculated by using the optimum policy and the optimality equation. Each iteration of policy and value

iteration methods takes $O(card(S^3))$ and $O(card(S)card(A))$ times, respectively (Sun et al., 2013). Thus, when the state space is large and action set is relatively small in an MDP formulation, value iteration should be preferred over policy iteration.

For real-sized instances, MDP formulations can reach very large sizes quickly. This phenomenon is known as 'curse of dimensionality' in the literature. Since the exact methods suffer from curse of dimensionality, approximation techniques have been an active research area within the MDP community. Neuro-dynamic programming is an approximation technique for MDPs that combine tools from reinforcement learning to approximate the value functions (Bertsekas and Tsitsiklis, 1995). The most popular one of the approximation techniques for MDPs is ADP, a solution framework in which the value function is approximated by using linearization or simulation-based methods. For more detailed information related to ADP, readers are referred to Powell (2007).

ADP methods may be categorised under two main streams: linear programming and simulation-based methods. The first stream is useful when the expectation in optimality equations can be computed exactly. In these methods, the dynamic formulation is first converted to an equivalent linear programming model and then solved with well-established linear programming solution methods. However, the resulting model usually contains a very large number of constraints (Haugh and Kogan, 2007) that can only be solved with the reduction techniques, such as constraint sampling or column generation.

The simulation-based methods may also be divided into two categories: Q-learning algorithms and value/policy evaluation (Powell, 2007). The first one is based on estimating and updating the value function for each state-action pair, whereas the value and policy evaluation algorithms compute the value function approximation of each state and a single policy at each iteration, respectively.

Within both streams, direct methods use simulation to calculate the value function estimates of the sampled states and fit an approximation structure to these samples. The value function estimates may be stored in a table format, known as the lookup table. As the number of iterations increases, the previous estimates for state values are averaged with the new values. The disadvantage of a direct method is that it requires a large memory to store the lookup table. The indirect methods use a linear combination of basis function approximations to solve the optimality equations. In these methods, the approximate state value is obtained by weighting and summing the basis functions which represent selected features of the state variable. The basis function approximation allows to estimate the state values that are not visited by the lookup table method. The selection of these features is a state-of-art and depends on the problem structure.

There are two methods within the basis function approximation: on-policy and off-policy. The first method initially finds the state values visited in a sample path as in the lookup table based ADP algorithm. Then, it applies regression methods on these approximate state values to find the best weight levels for the basis functions. Finally, these weights are used to find the approximate state values in the next iteration, and the greedy policy accordingly. In other words, it approximates the state values *within* the policy selection. The off-policy basis function method applies the regression after the approximate state values and the policy are computed by the lookup table aggregation. In both methods, as the number of iterations increases, it is expected that the weights will converge through their true values. The direct and indirect methods may differ in the convergence speed depending on the problem, whereas, both may suffer from long simulation runs.

A simulation-based ADP algorithm may be single- or double-pass. In a double-pass ADP algorithm, first, a trajectory of states and outcomes are gen-

erated with the help of simulation and the initial approximate state values and the greedy actions are obtained. Then, these approximate values are updated with a backward pass i.e. by using the information from the future steps of the same trajectory. A double-pass algorithm is especially preferred when the value function differs for some periods of the planning horizon. In such a case, an action at a period may have an effect on the costs incurring in the future periods.

## 1.3.2 Queuing Theory

A specific area of study focusing on service systems is known as queuing theory. This section presents a brief overview of queuing theory, a modelling approach for service systems involving queue(s). A queuing system can be defined by customers, server(s), input process, service mechanism, system capacity and queue discipline. In this framework, the customer is the entity demanding the service while the server is the entity providing the service. The input process describes the arrival pattern of the customers, usually in terms of the distribution of the random variables. The service mechanism consists of the number of servers, the service time, and the form of providing service (batch or single). System capacity denotes how many customers can be present in the system at any time. Finally, queuing discipline explains all other factors related to the order of queue selection such as how servers accept the next customer to be served. The most popular queuing disciplines can be listed as first-come first-served, last-come first-served, and random selection for service.

All the information related to a queuing system can be represented with a notational taxonomy developed by Kendall (1953). With this method, a queue can be represented as $\alpha/\sigma/m/\beta/N/Q$. The first and second symbols in this notation describe the distributions of the input process and service time, respectively, and can take symbols like $M$ (exponential), $G$ (general), $D$ (deterministic), or

$E_k$ (Erlang). The third symbol, $m$, shows the number of servers. For example, M/M/1 represents a Poisson arrival process, exponential service times and a single server queuing system. As there are more factors involved, they are added sequentially to the whole notation. On some cases, each server may have its own queue which is known as a parallel queue system. Alternatively, there can be a network of queues in which customers move between different servers in a sequential way. For example, surgical process can be represented as a queuing network that requires an operating room, then a recovery bed, and finally an intensive-care unit bed.

The performance of a queuing system can be measured by average waiting time, number of people waiting in the queue, server utilization, etc. The analysis of a queuing system aims to provide a mathematical representation of these performance metrics. The analysis starts by assuming that the statistical equilibrium exists; the queuing system reaches to an equilibrium state in the long-run (Bhat, 2015). Let's consider the state transition probability of a system, represented with $\{Q(t), t \geq 0\}$ at time period $t$,

$$P_{ij}(s,t) = P[Q(t) = j | Q(s) = i], \quad s < t,$$

where the system is at state $j$ at time $t$ conditional on its state $i$ at time $s$. If the system attains a statistical equilibrium, then,

$$\lim_{t \to \infty} P_{ij}(s,t) = p_j,$$

which is independent from time $t$ and state $i$ (Bhat, 2015).

When the interarrival times in a queue follow exponential distribution, it is possible to obtain closed-form formulations of the performance metrics. However, if the interarrival process is not exponential, it is usually not possible to derive these closed-form formulations (Bandi and Bertsimas, 2012). These cases can be

analyzed by using approximate approaches (Whitt, 1993; Kimura, 1983).

Queuing theory is usually simple, provide generic results and require less data compared to other available methods for the analysis of queues such as simulation modelling (Fomundam and Hermann, 2007). Due to its success in analyzing complex service systems, queuing theory has been widely applied to healthcare management problems (Green, 2006). Queuing theory can be used to understand a healthcare system better, to figure out the reasons of an undesired performance or to make recommendations to improve it. For example, it can be used to find the optimum number of servers (capacity) to achieve a better performance in terms of the average waiting time or utilization. Queuing models can also be used to test different managerial strategies such as customer priority schemes or the degree of flexibility to be used in resource planning (Green, 2006). For a detailed review of queuing models in healthcare, readers are referred to Green (2006) and Lakshmi and Iyer (2013).

### 1.3.3 Robust Optimization

Stochastic programming assumes that uncertain parameters follow certain probability distributions. However, in some cases, this is not applicable due to the lack of data or simply because the data do not fit into any known probability distribution. For these cases, robust optimization (RO) offers a rigorous framework. This section provides an overview of RO.

RO handles optimization problems with certain degree of robustness against uncertainty that can be represented in deterministic and set-based forms. It provides a guaranteed performance even in the worst-case scenario. In other words, the general purpose of RO is to find a solution that is feasible for any realization of the uncertainty in an uncertainty set. Thus, it is usually preferred when the solutions are highly sensitive to the perturbations in the data or the worst-case

scenario cannot be afforded. Although it is a relatively young field, mostly flourished in the last 20 years, its computational tractability and suitability for many stochastic problems have generated a considerable RO literature (Bertsimas et al., 2011). RO differs from stochastic programming mainly because it does not require any information regarding the probability distributions of uncertain parameters. Also, it is significantly different from the sensitivity analysis since the solution of an RO formulation is feasible regardless of the data. Let's consider a generic optimization problem under uncertainty,

$$\max_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x}, \boldsymbol{\xi}). \tag{1.5}$$

The robust counterpart of (1.5) can be written as:

$$\max_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x}, \boldsymbol{\xi}), \forall \boldsymbol{\xi} \in U, \tag{1.6}$$

where $U$ is an uncertainty set that the random parameters, $\boldsymbol{\xi}$, can take any value from. In general, an uncertainty set specifies a set of values that the uncertain parameters can take. In other words, by optimizing over an uncertainty set, the original problem (1.5) is converted to its robust counterpart (1.6).

An uncertainty set can be defined as discrete or continuous. A discrete uncertainty set contains discrete values representing possible realizations of uncertainty. On the other hand, a continuous uncertainty set contains an infinite number of possible realizations. Other than these, uncertainty sets can be classified according to their shapes or structures. Readers are referred to Bertsimas et al. (2011) for a review of uncertainty set structures.

The most common set structures studied in the literature are ellipsoidal, polyhedral, and cardinality-constrained (Bertsimas et al., 2011). For the ease of understanding, we will explain different uncertainty set structures through a robust linear problem i.e. the robust counterpart of a linear optimization problem

that can be written as

$$\min_{\mathbf{x} \in \mathbf{X}} \quad c^T \mathbf{x},$$

$$\text{subject to} \quad a_i \mathbf{x} \leq b_i, \ a_i \in U, \ i = 1, \cdots, m,$$

where $a_i$ for $i = 1, \cdots, m$ are the uncertain parameters belonging to the uncertainty set $U$. An ellipsoidal uncertainty set can be defined as

$$U = \{(a_1, \cdots, a_m) : a_i = a_i^0 + \Delta_i \xi_i, \ \ i = 1, \cdots, m, ||\xi||_2 \leq \rho\},$$

where $a_i^0$ denotes the nominal value of the $i^{th}$ uncertain parameter and $\rho$ is a parameter controlling the size of the set defined by modeller. The robust linear problem with an ellipsoidal uncertainty set can be written as a second-order cone problem (Bertsimas et al., 2011):

$$\min_{\mathbf{x} \in \mathbf{X}} \quad c^T \mathbf{x},$$

$$\text{subject to} \quad a_i^0 \mathbf{x} \leq b_i - \rho ||\Delta_i \mathbf{x}||_2, \ i = 1, \cdots, m,$$

that can be solved with exact methods.

A polyhedral uncertainty set, a special case of an ellipsoidal set (Ben-Tal and Nemirovski, 1999), can be defined as

$$U = \{(a_1, \cdots, a_m) : a_i = a_i^0 + \Delta_i \xi_i, \ \ i = 1, \cdots, m, \ \ D\xi + q \geq 0\},$$

where $D$ and $q$ are the parameters defined by modeller. When a robust linear problem has a polyhedral uncertainty set, it can be written as a linear optimization problem that can be solved to optimality with many commercial solvers (Bertsimas et al., 2011). As an alternative, a cardinality-constrained set limits the number of parameters that are allowed to deviate from their nominal values. With this type of uncertainty set, modellers can control the trade-off between the conservativeness and the optimality of the solution.

The shape of the uncertainty set affects the tractability of an RO problem significantly (Bertsimas et al., 2011). Let's define the feasible (solution) set of an RO problem as,

$$X(U) = \{\mathbf{x}|g(\mathbf{x}, \boldsymbol{\xi}) \leq 0, \quad \forall \boldsymbol{\xi} \in U\}.$$

Then, the tractability usually refers to $X(U)$ being convex in $\mathbf{x}$ (Bertsimas et al., 2011). However, the robust counterparts of convex optimization problems are mostly intractable (Bertsimas et al., 2011).

Recently, several important developments have occurred in RO; Gabrel et al. (2014) present a detailed review of these advances since 2007. Two main streams with recent advances are the use of risk theory to describe the uncertainty sets and multi-stage RO models. In the first stream, uncertain parameters are assumed to follow unknown probability distributions. This new approach is also called distributionally robust optimization (Delage and Ye, 2010; Goh and Sim, 2010) as an attempt to close the gap between stochastic programming and RO. The second stream, multi-stage RO, is also known as dynamic RO in which recourse decisions are included into a robust formulation in a tractable fashion. Dynamic RO aims to improve the long-term tradition of static RO modelling; for a detailed review, readers are referred to Duzgun and Thiele (2010).

The next chapter presents the capacity planning study for a stem-cell donation network under uncertainty.

# Chapter 2

# Capacity Planning for Network of Stem-cell Donation Centres

## 2.1 Introduction

A healthcare system providing medical services involves multiple stakeholders such as doctors, managers, and public policy makers. These stakeholders have to deal with different operational complexities and various uncertainties inherent in the systems to decrease cost and to increase patient satisfaction. The operational complexities arise from the fact that healthcare procedures usually consist of multiple steps involving different specialities and specific resources at the same time. These steps incorporate several uncertainties such as treatment outcomes, test results, and disease progression. Due to the increasing pressure to minimize operational costs and high demand for an improved service, an effective capacity planning becomes crucial for the healthcare management.

Healthcare is one of the largest sectors in the world, accounting around 10% of global GDP (Deloitte, 2016). Transplantation services, including organs, tissues or cells, take an important place within the healthcare sector. In partic-

ular, kidney transplantation comprises around 70% of all transplants in the UK (ODT, 2016). The most common cell transplantation is that of stem-cells defined as infusion or injection of healthy stem-cells to replace the diseased or damaged ones (Fruchtman, 2003). The stem-cell transplant is crucial for several illnesses: leukaemia, anaemia, various blood diseases and immune system problems. Gratwohl et al. (2013) reported that around 60,000 stem-cell transplantations occur annually worldwide.

The main distinguishing feature of stem-cell donation is that it does not cause any harm to the donor. For this reason, a patient requiring a stem-cell transplant has the option of conducting a search for matching cells from unrelated living donors. Due to this feature, it is possible to develop a large donor database for stem-cell donation. On the other hand, in the organ transplants, the donor's health may be affected from the operation (for kidney transplants) or it is not even possible to have transplant from a living donor (e.g. for heart, lung, pancreas, etc.). Thus, the donor for an organ transplantation is usually a relative of the patient. Otherwise, the patient has to wait for a deceased (recently died) donor.

Stem-cell donation centres are public facilities serving to patients in need of stem-cell transplants that require complex search procedures and the advanced blood-gene tests to reveal any matching between donors and patients. These advanced blood-gene tests may be done in the laboratories belonging to the state hospitals such as in Turkey or private transplantation centres as in the UK. In this chapter, we consider the former case, i.e. when the centre and laboratories belong to the same body. Even though the search process is expensive, ranging between $25,000 to $150,000 (Lee et al., 2000), a transplantation is the last chance for the survival of some patients. Also, international sources cost nearly 10 times more than the national sources due to the special carriage requirement for stem-cells.

In general, life expectancy of patients waiting for a stem-cell transplant

is short (Odejide, 2014); in other words, the probability of patient death during the search is high. Therefore, the processing time to search for the best match is crucial for the survival of these patients. For a donation centre with over-utilized work capacity, the search process generally takes longer and this of course increases the number of deaths. For example, Anthony Nolan, a stem-cell donation register in the UK, reported that they were able to supply suitable donors for only half of the patients needing transplantation in 2014 (Antony Nolan, 2016). As another example, in Turkey, around 70 out of 1000 patients die within a year due to the lack of suitable stem-cell donations (Beksac, 2014).

All stem-cell donation centres operating within a country are managed by an authority such as a foundation or the state which is also responsible for the distribution of the budget among them. The centres usually operate independently, but have access to the same national donor database. The main reason behind the incoordination is the geographical distance between centres due to their scarcity and high establishment costs; for example, there are only two centres in Turkey (Beksac, 2014). Each centre typically aims to increase the number of successful searches as well as utilizing the resources in a cost-effective way. In order to decrease the number of deaths, capacities of centres could be enlarged or new facilities could be opened. However, high operational and infrastructure costs limit the capacity expansion opportunities. Thus, the central authorities need to plan the capacities of donation centres to maximize the overall performance considering the inflexibility of resources. In this chapter, we are concerned with a stochastic capacity planning problem that takes into account search operations of stem-cell donation centres.

In the literature, there exist several statistical studies that aim to determine the optimum donor database levels for the stem-cell donation centres in terms of cost-effectiveness (see for instance, Hurley et al., 2003; Muller et al., 2003;

Kollman et al., 2004). The authors calculate the probability of finding a suitable donor given a donor database level and gene compositions. Then, they consider the average cost of adding a donor to the database. It is found that adding more donors provides a decreasing utility after a certain database level. Thus, they identify the best donor database level in terms of the cost-effectiveness. However, according to our best knowledge, the stochastic capacity planning problem for a network of stem-cell donation centres has not been studied yet.

The closest service system to stem-cell donation centres is blood banks. The operational and strategic aspects of blood banks are handled with various optimization and simulation based techniques; for instance, see Alfonso et al., (2013) and Gunpinar (2013). However, those models cannot be directly applied for management of stem-cell donation centres due to the several distinguishing differences between the operations of these service providers. Blood banks collect and store bloods before releasing them to hospitals. Blood is a perishable product with limited storage period in blood banks whereas blood collected in stem-cell donation centres is discarded once the gene information of the prospective donor is revealed. Also, blood banks can work as mobile teams and most importantly do not require expensive advanced blood testing machines unlike stem-cell donation centres.

In this chapter, we develop a stochastic capacity planning model for a network of stem-cell donation centres. In particular, we are concerned with allocation of the capacity budget among stem-cell donation centres such that their performance in finding suitable donors of stem-cells under uncertainty is maximized. The contributions of this chapter are twofold in terms of modelling and solution approaches.

- We model complex search operations to be integrated into the capacity decision-making problem using stochastic programming. We investigate the

28

impact of uncertainty on search operations as well as strategic capacity decisions. In particular, uncertainties arising in patients' lifetime and test results are represented by a probability distribution and a scenario-based approach, respectively. The service capacity of these centres cannot be easily adapted to the variation in number of patient arrivals. Therefore, patient arrivals, that affect the average service time of the centre, and consequently the number of successful searches are also assumed to be uncertain. For the advanced blood testing system searching for the best match, we use a first-come first-served (FCFS) queue assuming that arrival and service times follow general distributions.

- In order to obtain robust solutions against the uncertainty in patient arrival and blood test duration, we consider an approximate upper-bound of the waiting time in blood testing from robust queueing. The resulting problem formulation leads to a non-linear integer programming model that is computationally difficult to solve. We reformulate the model into an integer linear programming model under certain distribution assumptions. Finally, we design a series of computational experiments to illustrate the performance of the developed model. The numerical results obtained by the in-sample and out-of-sample experiments show that the proposed method provides a good approximation to the waiting time of blood samples. Moreover, they imply that the cost-effectiveness of the network improves as the number of centres increases.

The rest of this chapter is organized as follows. In Section 2.2, the literature on the capacity planning problems within healthcare is reviewed. Section 2.3 describes the patient and donor side operations of a stem-cell donation centre. Section 2.4 focuses on formulation of the stochastic capacity planning model. In Section 2.5, we introduce a scenario-based capacity planning model for a network

of stem-cell donation centres. The numerical results are displayed in Section 2.6.

## 2.2   Stochastic Capacity Planning in Healthcare

Capacity planning, as one of the main problems in OR, deals with an effective utilization of available resources to meet changing demand for products or services. The capacity planning problems have been widely studied in the literature for various service providers such as banks, hotels, and hospitals as well as for production and supply chain management purposes.

The healthcare capacity planning problems involve various uncertainties like demand, staff availability and medical results. In particular, variabilities in different factors such as patient arrivals and service processes may result in excessive waiting times and poor utilization of facility resources (Salzarulo et al., 2011). In order to handle uncertainties, there exists different approaches. For example, the sensitivity analysis is applied only after a solution is obtained as a post-optimization tool. Unlike the sensitivity analysis, the expected value of underlying random factor can be used to find the solution of an optimization problem under uncertainty. Note that this solution is the optimum only for a single realization of the uncertainty (corresponding to the expected value), possibly giving undesirable results for other realizations. On the other hand, the scenario-based stochastic programming approach considers an adequate range of possible realizations as well as probabilities of their occurrences and optimizes the expected performance of the system in view of a finite number of discrete scenarios. The scenario-based uncertainty modelling approach provides a flexible way of defining the decision process where each scenario represents a possible realization of the uncertainty associated with the occurrence probability. In this section, we briefly review the capacity planning models within the healthcare service management

under uncertainty.

Queuing theory, as a modelling approach to obtain performance measures in a queuing system, has been widely applied for the capacity planning of healthcare services; a related review can be found in Fomundam and Hermann (2007). For example, the built-in queuing formulas can be used to find the number of servers (capacity) required to achieve a certain degree of performance as in Creemers and Lambrecht (2009). Hulshof et al. (2013) also use queuing theory to model the elective patient admission and intermediate term resource allocation for hospitals with uncertain treatment paths and number of arrivals. They consider different queues for different types of services with time-dependent resource-capacity levels. The objective is to obtain the optimum number of patients to be served at each time period. Similarly, Cochran and Roche (2009) apply queuing theory to test various capacity design alternatives to be used in real time when the capacity cannot meet the demand. Bretthauer et al. (2011) consider the capacity planning problem for healthcare operations with blocking between different units. Similarly, Castillo et al. (2009) study the optimal capacity and location of healthcare facilities modelled as queues with exponential service times and Poisson arrivals. By considering the time-varying demand in hospitals, Green et al. (2007) analyze the staffing requirement in hospitals based on the queuing analysis. The main drawback of queuing models comes from their intractability due to non-linear 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. (2010) introduces a discrete-event simulation model to analyze the operations management of an intensive care unit and uses 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. DeAngelis et al. (2003) consider simulation optimization to determine the capacity of a transfusion centre under multiple objectives: cost minimization to achieve a fixed waiting time and minimization of the waiting time under a limited budget. The queuing system in centres is modelled with a discrete event simulation. The objective functions are approximated by fitting functions to the data generated by the simulation model. The processes in a blood collection unit are modelled by using a simulation-based approach by Alfonso et al. (2013). 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.

The optimization models in healthcare capacity planning mostly focus on single hospital or department. However, the interconnection between departments and hospitals has a significant effect on their performance. There are several studies considering this interconnection for the capacity planning under uncertainty for a network of hospitals. Flessa (2000) develops a model to allocate resources in the preventive and curative services in hospitals and dispensers. The author considers different types of diseases and assumes fixed service units required for each disease type in different institutions. The optimization model distributes a fixed budget among different institutions based on the expected patient arrivals. Stummer et al. (2004), Govind et al. (2008), Santibanez et al. (2009) and Gunes et al. (2010) focus on the location and number of beds in hospitals within a network to minimize operation cost and maximize patient utility. They consider the patient flows either at the unit level or regional level to find the optimum bed/staff capacities. However, they do not model the operational details, but rather assume that patients stay for a fixed period of time. Unlike these authors, Mahar et al. (2011) study the location of the specialized services such as imaging

or neonatal intensive care. Their model identifies which hospitals in a network should have the specialized care services.

The hospital network capacity planning models developed in the literature are not directly applicable to modelling of search operations within stem-cell donation centres. The stem-cell donation centres have distinctive and complex operations, making the problem novel in this sense. Besides, the capacity planning model introduced in this chapter incorporates the queuing theory that has not been widely studied in the network capacity planning under uncertainty. The most relevant papers using queuing theory for healthcare network capacity planning are Pehlivan et al. (2012) and Asaduzzaman et al. (2010). Pehlivan et al. (2012) develop a mixed-integer 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. On the other hand, Asaduzzaman et al. (2010) develop a queuing model to find the optimum capacities of neonatal centres to minimize refusal and overflow probabilities. Unlike these authors, we employ a novel robust approach to derive the maximum waiting time in a queuing system. The resulting non-linear integer formulation is then approximated as a linear integer model that can be solved by exact methods.

## 2.3 Operations of a Stem-cell Donation Centre

We consider a network of stem-cell donation centres at each of which the same kind of search operations takes place. Stem-cell donation centres located in different areas operate independently to find suitable donors for the patients applying to the centre. They are financially coordinated by a central authority who may also set up national targets for centres to achieve. Before introducing the formulation of the capacity planning problem of the central authority, we summarize search

operations of a stem-cell donation centre in this section.

The search operations of a centre can be classified into two groups related to donor and patient sides. The patient side operations mainly consist of searching suitable donors for the patients who need a transplantation. In the donor side, the centre accepts and tests the bloods of prospective donors and updates the donor database. Figure 2.1 depicts various activities taking place at both patient and donor sides in a stem-cell donation centre. The search operations taking place at patient and donor sides use the same (national) donor database.



Figure 2.1: Main operations in a stem-cell donation centre

The patient-side operations are based on the search for the best match between a patient and donors' blood genes required for a transplantation. The possibility of finding a perfect match between family members is around 30% (Antony Nolan, 2016). The patients who cannot obtain a suitable donor from their family members require cells from suitable non-related donors. The search operations take place sequentially at three different levels. At the first stage, an online database search is performed as soon as a patient is admitted to the centre. The initial database search is relatively simple and can be completed in a short period. Suitable donors possessing the same blood characteristics with

the patient's blood gene structure are identified at the end of the online search. The number of eligible donors depends on the database level at the time of the initial search. As expected, a larger database increases the chance of finding more suitable donors.

At the second stage of the patient side operations, the suitable donors are contacted to provide a new blood sample. The donor can either go to the center to supply blood sample or send his/her own blood sample to the center. The blood samples collected from the suitable donors are stored temporarily and tested by a first-come first-served basis using advanced equipment in the blood-test laboratories working in collaboration with the centre. Further blood tests are conducted to find the best match for the patient among the eligible donors (identified at the first level). Total number of tests to be conducted for each patient is limited and can be determined by the central authority as a national policy. The duration of an advanced test is subject to a small variation. The patient remains in the process until all blood samples are processed. When a suitable donor with a perfect matching is found, the donation search for the patient is terminated and then a transplantation can take place. If there is no match between a patient and suitable donors, then an international search may be initiated as the final stage of the search process. If the international search becomes unsuccessful, then the search is completed without a transplantation.

The donor and patient side operations are performed independently. In the donor side operations, a number of donors arrives to the centre each day and provides blood samples at any convenient time. Then, some preliminary tests are conducted on the blood samples in a laboratory. These tests do not require any special equipment and are not as advanced as the ones applied for the donor bloods on the patient side operations; only preliminary information is gathered at this stage. After completing the tests, the bloods can be either stored (if there is

enough space) or discarded. Note that the characteristic of the blood revealed by the test is important rather than the blood itself. The unique information of each blood sample is recorded in the online donor database. When a perfect matching between the blood characteristic of a donor and a patient is found, the donor is called back to supply the actual donation of stem-cells.

The actual donation and transplantation take place in a hospital rather than the donation centre. The donor side operations only affect the performance of the stem-cell donation centre through the database level. On the other hand, the database level is mainly determined by the donation willingness in the country. Therefore, the donor-side operations are not taken into account for modelling search operations of stem-cell donation centres. Besides, these operations are simple and do not require strategic decision making.

The existence of sufficient capacity for the advanced blood tests plays an important role on the success probability of having a transplantation. On the other hand, a large capacity (at low demand season) may unnecessarily increase the operational cost. The search process involves real-time complex operations (as described above) and various uncertainties arising at different levels of the donation search processes. Thus, it is crucial to determine the optimum service capacity of each stem-cell donation centre by taking various uncertainties into account.

## 2.4 Formulation of the Stochastic Capacity Planning Problem

Stem-cell donation centres within a country usually do not interact with each other, but are controlled and financially supported by a central authority. Although the location, size and capacity of centres differ from each other, they

perform the same kind of search operations and contribute to the government's national targets. In this section, we introduce a stochastic capacity planning model for a network of stem-cell donation centres to optimize the overall performance of the network. Before introducing the capacity planning model, we first describe the model assumptions and the underlying uncertainties.

**Assumptions:** We make the following assumptions for the model development.

- The service capacities of the centres are assumed to remain the same during the planning horizon. We consider a first-come first-served queue with random service time in order to model the operations of advanced blood test.

- In general, an international search starts only after the results of all advanced tests are revealed. However, for cases where the medical situation of a patient is very critical, an international search may start as soon as the patient is admitted. Although these special cases are not taken into account for the sake of simplicity, the model introduced in this chapter can be easily modified to incorporate the medical condition of a patient.

- The international search is an independent process; therefore, a local authority, patient or any other external factors cannot influence its duration. In addition, the advanced tests of blood samples to be collected from suitable international donors, are usually conducted at their own centres. It is worthwhile to mention that the advanced-blood tests of the international donors might be done at the stem-cell donation center where the patient is registered. But for the sake of simplicity, these cases are omitted in the problem formulation.

- We also assume that the patient leaves the system at the end of unsuccessful national and international search operations. However, in practice, if the

patient is in severe medical conditions, a new search may be started for him/her. These cases are not taken into account in the current model.

**Uncertainty**: The search operations within a stem-cell donation center involve various exogenous and endogenous factors that directly affect the overall goals and capacity planning strategies of the center. We can classify these factors according to stages of patient arrivals and blood samples as well as the search operations at the national and international levels. They are described in more detail below:

- The arrival time of a patient is not known in advance and the total number of patients that have been waiting for a suitable donor varies over time. Similarly, number of suitable donors to be tested for each patient is uncertain.

- The completion time of a donor search is crucial on the success of donor search operations. It is basically determined as the sum of the durations taken for the national and international search operations.

    - Search duration at the national level depends on the travelling time of blood donors and the waiting time of blood samples in the system for the advanced test. The donor travelling time depends on an individual's behaviour and personal preferences; therefore, it is not known in advance by the center. In addition, the waiting time of blood samples in the system is also uncertain due to the variations in patient arrivals and testing duration.

    - Duration of the international search is affected by various factors such as the capacity and the demand at the international centres and also the frequency of patients' gene structure.

Thus, we model the donor travel time, the blood test waiting time and the international search duration as uncertain parameters.

- The results of matching tests using the national and international sources may be influenced by various endogenous (frequency of patient's gene structure) and exogenous (capacity of the donor pool) uncertainties. Note that the blood gene structure of a patient is constructed by millions of different combinations. Therefore, the test results are also assumed to be uncertain in the formulation of the capacity planning problem.

- Apart from these factors, the health condition of a patient independently influences the success of the donor search for a possible transplantation. Most of the patients seeking suitable donors have a critical health condition. Thus, the patient's lifetime is considered as an uncertain factor to be roughly predicted by the doctors. As explained further in the next section, we assume that the patient lifetime is assumed to follow a known distribution.

**Problem Formulation**: We consider a network consisting of $J$ stem-cell donation centres labelled as $j = 1, \cdots, J$. A central authority is responsible to allocate the budget capacity $B$ among centres. Consider a planning horizon $T$, that is discretized by time periods $t = 1, \cdots, T$. In practice, $T$ may represent a year while each time period corresponds to one week. Throughout the chapter, uncertain parameters are indicated by a tilde, $\tilde{*}$.

Let $\tilde{I}_j$ denote the number of patients (labelled as $i = 1, \cdots, \tilde{I}_j$) arrive to centre $j$ during the planning horizon. There can be a single or a batch arrival of patients at any time period. For each patient $i$, suppose that $\tilde{p}_{ij}$ number of suitable donors are found from the online search. The candidate donors are then invited to supply another blood sample for further testing. For each patient $i$, we introduce indices $k \in \{1, \cdots, \tilde{p}_{ij}\}$ to label the blood samples that are received from $\tilde{p}_{ij}$ donors.

Let $x_j$ be a discrete decision variable representing the capacity of centre $j$

for the advanced testing during the planning horizon. The waiting time of blood samples in the service system, $\widetilde{W}_{ijk}(x_j)$, is defined as the duration between the arrival of blood sample $k$ of patient $i$ to the blood-testing queue and the start of its test in centre $j$. Note that the time taken from the arrival of a blood sample till the test completion directly depends on the capacity of the centre. Specifically, each machine-staff pair in the advanced blood testing is identified as one unit of capacity.

Let $\tilde{t}_{ijk}$ denote time taken between the arrival of patient $i$ to centre $j$ and collection of the $k$-th blood sample for the patient. We introduce $\tilde{o}_{ijk}$ to represent the duration of the advanced blood test for blood sample $k$ of patient $i$ in centre $j$. Let $\tilde{u}_{ijk}(x_j)$ define the duration between the arrival of patient $i$ to centre $j$ and the test completion time of its $k$-th blood sample. We can compute $\tilde{u}_{ijk}(x_j)$ for $i = 1, \cdots, \tilde{I}_j$, $j = 1, \cdots, J$ and $k = 1, \cdots, \tilde{p}_{ij}$ as an accumulated outcome of uncertain waiting time and arrival time of the blood sample as follows:

$$\tilde{u}_{ijk}(x_j) = \tilde{t}_{ijk} + \tilde{o}_{ijk} + \widetilde{W}_{ijk}(x_j). \tag{2.1}$$

The search process is not only affected by the time to obtain the test results, but also the medical outputs. Thus, we need to take into account the search results of each patient arriving to the centre. Let $\tilde{r}_{ij}$ and $\tilde{z}_{ij}$ represent the search results obtained by the national and international sources, respectively, for patient $i$ admitted to centre $j$. If at least one blood test result is positive, then $\tilde{r}_{ij}$ takes 1. If the results of all blood tests are negative, then $\tilde{r}_{ij}$ is assigned to 0. Similarly, if the search using international sources for patient $i$ is successful, then $\tilde{z}_{ij}$ takes 1; otherwise, it takes 0. If a search at the international level is never initiated, then it is fixed at zero ($\tilde{z}_{ij} = 0$). According to the revealed results of the tests taken at the national level, either a transplantation takes place or the search

for a suitable donor is carried out at the international level, assuming that the patient is still alive. The duration of an international search is represented with $\tilde{v}_{ij}$ if activated. Thus, for patient $i$ in the system, one of the following possible outcomes is realized:

- At least one positive result is obtained from the national sources ($\tilde{r}_{ij} = 1$ and $\tilde{z}_{ij} = 0$).

- No positive result is achieved from the searches using the national and international sources ($\tilde{r}_{ij} = 0$ and $\tilde{z}_{ij} = 0$).

- No positive result is obtained from the national sources, but a positive result is acquired from the international search ($\tilde{r}_{ij} = 0$ and $\tilde{z}_{ij} = 1$). Notice that the case $\tilde{z}_{ij} = \tilde{r}_{ij} = 1$ never occurs since a search at the international level for patient $i$ is initiated only after no suitable donor is identified from the national sources.

Let $\tilde{d}_{ij}(x_j)$ define the search completion time taken from the admission of patient $i$ to centre $j$ until the end of all searches for a suitable donation to be completed. The search completion time depends on the number of suitable donors as well as the search outcomes at the national and international levels. A search at the national level is terminated only when the advanced tests for all suitable donors are completed. Recall that a transplantation can be conducted only when the perfect match from suitable donors is found (Antony Nolan, 2016). Thus, if at least one positive outcome from the national search ($\tilde{r}_{ij} = 1$) is achieved, then the search completion time for patient $i$ is determined as the maximum of whole test completion times as $\tilde{d}_{ij}(x_j) = \max\limits_{k=1,\cdots,\tilde{p}_{ij}} \{\tilde{u}_{ijk}(x_j)\}$. The reason of using a 'maximum' function is that the centre waits until the last test to be able to have backup donors and the best possible match.

On the other hand, if no perfect match is found from the national sources ($\tilde{r}_{ij} = 0$), but there is at least one positive result obtained from the international

sources ($\tilde{z}_{ij} = 1$), then the search completion time is computed as sum of the completion times of the national and international searches:

$$\tilde{d}_{ij}(x_j) = \tilde{v}_{ij} + \max_{k=1,\cdots,\tilde{p}_{ij}} \{\tilde{u}_{ijk}(x_j)\}.$$

Finally, if no search at the national and international levels is successful ($\tilde{r}_{ij} = \tilde{z}_{ij} = 0$), then the search completion time is assigned to a big number, $\tilde{d}_{ij}(x_j) = M$, to imply that the patient remains in the system as long as being alive. Three cases showing computation of $\tilde{d}_{ij}(x_j)$ can be summarised in a compact form as follows:

$$\tilde{d}_{ij}(x_j) = \begin{cases} \max_{k=1,\cdots,\tilde{p}_{ij}} \{\tilde{u}_{ijk}(x_j)\}, & \text{if } \tilde{r}_{ij} = 1, \\ \tilde{v}_{ij} + \max_{k=1,\cdots,\tilde{p}_{ij}} \{\tilde{u}_{ijk}(x_j)\}, & \text{if } \tilde{r}_{ij} = 0 \ \& \ \tilde{z}_{ij} = 1, \\ M, & \text{otherwise.} \end{cases} \quad (2.2)$$

A successful search process leads to transplantation if the patient is still alive when the search process is terminated. Suppose that $\tilde{l}_{ij}$ is the expected lifetime of patient $i$ to be estimated when admitted to centre $j$. Let us define a binary variable $\tilde{y}_{ij}(x_j)$ that takes 1 if the search for patient $i$ admitted to centre $j$ is unsuccessful and 0, otherwise. Thus, the relationship between the search outcome, the search completion time and the life expectancy of patient $i$ at centre $j$ can be expressed as follows:

$$\tilde{y}_{ij}(x_j) = \begin{cases} 0, & \text{if } \ \tilde{d}_{ij}(x_j) \leq \tilde{l}_{ij}, \\ 1, & \text{otherwise.} \end{cases}$$

Then, the number of unsuccessful searches conducted at centre $j$ can be easily computed as $\sum_{i=1}^{\tilde{I}_j} \tilde{y}_{ij}(x_j)$. The central authority needs to determine the capacity $x_j$ of each centre $j$, for $j = 1, \cdots, J$, such that total number of expected trans-

plantations (or unsuccessful searches) over all donation centres of the network is maximized (or minimized) in view of the budget restriction.

Given a unit-capacity cost $C_j$ of centre $j$ during the planning horizon $T$, we must ensure that the network capacity cost should not exceed the available budget $B$. This can be stated by a linear budget constraint as $\sum_{j=1}^{J} C_j x_j \leq B$. Then the stochastic capacity planning model (SCP) for a network of stem-cell donation centres can be formulated as follows:

$$\text{SCP:} \quad \min_{x_j \in \mathbb{Z}^+} \quad \sum_{j=1}^{J} \mathbb{E}\left[\sum_{i=1}^{\tilde{I}_j} \tilde{y}_{ij}(x_j)\right],$$
$$\text{subject to} \quad \sum_{j=1}^{J} C_j x_j \leq B.$$

This is a complex problem where the expectation in the objective function needs to be computed over all types of uncertainties given random number of arrivals. In order to do this, we adopt a scenario-based stochastic programming approach to determine the optimal capacities of the donation centres in a network.

## 2.5   Scenario-based Capacity Planning Model

In order to capture various events (including emergency situations) arising in the real life operations of a stem-cell donation network, we introduce a finite number of discrete scenarios (or may so-called cases) each of which represents a possible future realization of random patient arrivals. These scenarios are generated by using past data and statistics. Let $S$ denote total number of scenarios. Each scenario (represented by $s = 1, \cdots, S$) displays a sequence of patient arrivals with the corresponding probability $\omega_s$, where $\sum_{s=1}^{S} \omega_s = 1$. It also captures the information regarding total number of patients $I_{js}$ arriving to centre $j$ during the planning horizon. In this section we describe the scenario-based capacity planning model. The notation used for a specific scenario $s$ along with the operational

43

diagram is illustrated in Figure 2.2.

For each patient $i = 1, \cdots, I_{js}$ under scenario $s$, $p_{ijs}$ number of donors are assumed to be identified from the initial search and invited to the donation center $j$ for the advanced test. Indices $k \in \{1, \cdots, p_{ijs}\}$ label the blood samples that are received from $p_{ijs}$ donors for patient $i$ at centre $j$ under scenario $s$. Each donor (or blood sample $k = 1, \cdots, p_{ijs}$) of patient $i$ arrives to centre $j$ after $t_{ijks}$ periods from the time when the invitation is sent.

Figure 2.2: A description of the patient side operations along with the notation used for scenario $s$.



The service time for blood sample $k$ under scenario $s$ is denoted by $o_{ijks}$. Accordingly, we define the waiting time in advanced testing queue $W^s_{ijk}(x_j)$ of donor blood sample $k$ for patient $i$ at centre $j$ under scenario $s$. In addition, the search results obtained by the national and international sources are represented as $r_{ijs}$ and $z_{ijs}$, respectively, for patient $i$ in centre $j$ under scenario $s$. All these scenario-dependent parameters are deterministic under scenario $s$. On the other hand, the lifetime $\tilde{l}_{ijs}$ of patient $i$ at centre $j$ under scenario $s$ is assumed to be random, and follows a known distribution $f_{ijs}$. In general, one can only obtain probabilistic information of patient life-expectancy given his/her conditions during the search process.

The following rules express the cases when the search operations at centre

$j$ for patient $i$ under scenario $s$ lead to a transplantation or not:

$$\tilde{y}_{ijs}(x_j) = \begin{cases} 0, & \text{if } \max_{k=1,\cdots,p_{ijs}} \{t_{ijks} + o_{ijks} + W^s_{ijk}(x_j)\} - \tilde{l}_{ijs} < 0 \ \& \ r_{ijs} = 1, \quad (2.3) \\ 0, & \text{if } \max_{k=1,\cdots,p_{ijs}} \{t_{ijks} + o_{ijks} + W^s_{ijk}(x_j)\} + v_{ijs} - \tilde{l}_{ijs} < 0 \ \& \ z_{ijs} = 1, \\ 1, & \text{otherwise.} \end{cases}$$

Then, we can state the scenario-based capacity planning model (SCP$_{\text{scen}}$) as follows:

$$\text{SCP}_{\text{scen}} : \quad \min_{x_j \in \mathbb{Z}^+} \quad \sum_{s=1}^{S} \omega_s \sum_{j=1}^{J} \sum_{i=1}^{I_{js}} \mathbb{E}[\tilde{y}_{ijs}(x_j)],$$

$$\text{subject to} \quad \sum_{j=1}^{J} C_j x_j \leq B.$$

The size of SCP$_{\text{scen}}$ depends on the number of discrete scenarios, the number of centres, and total number of patients arriving to the centres during the planning horizon. In order to solve SCP$_{\text{scen}}$, we need to compute $\mathbb{E}[\tilde{y}_{ijs}(x_j)]$ in view of all scenarios. This involves the determination of the waiting time of each blood sample, $W^s_{ijk}(x_j)$.

As mentioned before, we model the advanced blood testing system as an incapacitated FCFS queue that involves multiple servers with general arrival and service time distributions. An assumption of general distribution prevents any possible inaccuracy or errors occurring due to an imprecise fitting of data of the underlying distributions. However, it is computationally challenging to derive a closed-form formulation for the waiting time for each blood sample, $W^s_{ijk}(x_j)$, even for given capacity decisions $x_j$ (Tijms et al., 1981). 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 (Tijms et al., 1981). Thus, we approximate the search success by considering the average waiting time of each blood sample in the system. However, the capacity decisions made in view of average waiting time of

the blood testing queue may cause severe delays during the high demand seasons. Therefore, we consider the upper bound of waiting time of blood samples in the queue for determining the optimal service capacity of donation centres. This basically implies the worst-case approach for the waiting times in blood testing queues. In this way, the donation centre can accommodate the worst outcome of uncertain waiting time. Next, we derive an approximate formulation of the maximum waiting time for blood samples under each scenario.

## 2.5.1  An Approximation to Maximum Time Spent in Queue

It is worthwhile to mention that there exists different approximation methods for the maximum waiting time in a queue; for instance, see Gupta and Osogami (2011). However, as pointed out by Bandi and Bertsimas (2012), these approximations do not lead to realistic results when the arrival process follows a distribution apart from Poisson. In order to overcome this problem, Bandi and Bertsimas (2012) proposed an alternative approximation method to compute an upper bound on the waiting time when the arrival and service times are independent and identically distributed (i.i.d) random parameters following an unknown distribution in an FCFS queue with $x$ servers. We now provide a brief overview of this approach and then explain how to apply it for the donor blood samples.

Let $T_k$ and $Y_k$ represent random interarrival and service times for samples $k = 1, \cdots, K$, respectively. The first moment (mean values $1/\mu$ and $1/\lambda$) of the random service and interarrival times are estimated from the past data. Assume that $T_k$ and $Y_k$ belong to uncertainty sets $U^{arrv}$ and $U^{serv}$, respectively. Moreover, the sizes of these uncertainty sets are determined by parameters $\Gamma^{arv}$ and $\Gamma^{serv}$ that basically measure the variability in the interarrival and service times, respectively. The uncertainty set $U^{arrv}$ for interarrival times $T_k$ of samples $k = 1, \cdots, K$

is defined as follows:

$$U^{arrv} = \left\{ (T_1, T_2, ..., T_K) \ \middle| \ \frac{\left| \sum_{k=m+1}^{K} T_k - \frac{K-m}{\lambda} \right|}{\sqrt{(K-m)}} \leq \Gamma^{arrv}, \ \ \forall m \leq m_0 \right\}, \quad (2.4)$$

where $m_0$ can be set to $K - 30$. Similarly, the uncertainty set $U^{serv}$ for service times $Y_k$ of samples $k = 1, \cdots, K$ is

$$U^{serv} = \left\{ (Y_1, Y_2, \cdots, Y_K) \ \middle| \ \frac{\left| \sum_{k=m+1}^{e} Y_{kx+b} - \frac{e-m}{\mu} \right|}{\sqrt{(e-m)}} \leq \Gamma^{serv}, \ \forall m \leq e-1, \ 0 \leq b < x \right\},$$

where the accumulated service times are calculated over the partitions of service times into $x$ groups with sizes $e = \lfloor K/x \rfloor$ to reflect the multi-server nature of the problem. The following proposition states the upper-bound $\overline{W}(x)$ on the waiting time $W(x)$ in view of these uncertainty sets.

**Proposition 1** *(Bandi and Bertsimas, 2012) Assume that the interarrival and service times for an FCFS queue with $x$ servers belong to the uncertainty sets $U^{arrv}$ and $U^{serv}$, respectively. The approximate upper bound $\overline{W}(x)$ on the waiting time in the queue can be calculated as,*

$$\overline{W}(x) = \frac{\lambda(\Gamma^{arv} + \Gamma^{serv}/\sqrt{x})^2}{4\big[1 - \lambda/(\mu x)\big]}. \quad (2.5)$$

**Proof.** Readers are referred to Bandi and Bertsimas (2012) for the proof and the details of the parameter estimation. ∎

Notice that in a stable queuing system, the traffic density must be smaller than unity, that is $\frac{\lambda}{\mu x} < 1$. In other words, the number of servers $x$ must be larger than $\lambda/\mu$.

Let's assume that the interarrival and service times of the donor blood samples are i.i.d. The interarrival times of blood samples are subtracted from

the patient arrival times and donor travel times $t_{ijks}$ in the generated data. The parameters, $\lambda$ and $\mu$, are estimated from the past data while the parameters $\Gamma^{arv}$ and $\Gamma^{serv}$, denoting the variation in the interarrival and service times, can be set to a fixed number times of the standard deviation of the interarrival and service times, respectively. We can then apply Proposition 1 to compute the upper bound $\overline{W}_j(x_j)$ on the waiting time $W_{ijk}^s(x_j)$ for donor blood sample $k$ of patient $i$ in centre $j$ with capacity $x_j$ under scenario $s$. In this case, the inequality $\overline{W}_j(x_j) > W_{ijk}^s(x_j)$ is ensured for all blood samples $k = 1, \cdots, p_{ijs}$ of any patient $i = 1, \cdots, I_{js}$ arrived to center $j$ under different scenarios $s = 1, \cdots, S$. By replacing $W_{ijk}^s(x_j)$ by $\overline{W}_j(x_j)$ in (2.3), $\tilde{y}'_{ijs}(x_j)$ is obtained as

$$
\tilde{y}'_{ijs}(x_j) = \begin{cases} 0, & \text{if } \overline{W}_j(x_j) - \tilde{l}_{ijs} + \max_{k=1,\cdots,p_{ijs}} \{t_{ijks} + o_{ijks}\} < 0 \ \& \ r_{ijs} = 1, \quad (2.6) \\ 0, & \text{if } \overline{W}_j(x_j) - \tilde{l}_{ijs} + \max_{k=1,\cdots,p_{ijs}} \{t_{ijks} + o_{ijks}\} + v_{ijs} < 0 \ \& \ z_{ijs} = 1, \\ 1, & \text{otherwise.} \end{cases}
$$

Notice that we have $\tilde{y}'_{ijs}(x_j) \geq \tilde{y}_{ijs}(x_j)$ since $\overline{W}_j(x_j) > W_{ijk}^s(x_j)$ holds. Then the scenario-based capacity planning model $\text{SCP}_{\text{scen}}$ can be rewritten as the following approximated optimization model

$$
\text{SCP}_{\text{appx}}: \quad \min_{x_j \in \mathbb{Z}^+} \quad \sum_{s=1}^{S} \omega_s \sum_{j=1}^{J} \sum_{i=1}^{I_{js}} \mathbb{E}[\tilde{y}'_{ijs}(x_j)],
$$
$$
\text{subject to} \quad \sum_{j=1}^{J} C_j x_j \leq B.
$$

It is worthwhile to mention that the optimal capacities of stem-cell donation centres obtained from $\text{SCP}_{\text{appx}}$ are more conservative towards the uncertainty in waiting times. Therefore, it leads to a higher objective function value than the one obtained from $\text{SCP}_{\text{scen}}$.

To be able to solve $\text{SCP}_{\text{appx}}$, an analytical form of $\mathbb{E}[\tilde{y}'_{ijs}(x_j)]$ needs to be derived. We assume that the lifetime expectancy $\tilde{l}_{ijs}$ of each patient $i$ arriving to centre $j$ under scenario $s$ follows a general distribution. Let $f_{ijs}(.)$ represent the

probability distribution function of $\tilde{l}_{ijs}$. For the sake of convenience, we introduce random parameters

$$\tilde{l}'_{ijs} = \tilde{l}_{ijs} - \max_{k=1,\cdots,p_{ijs}}\{t_{ijks} + o_{ijks}\} \ \ \text{and} \ \ \tilde{l}''_{ijs} = \tilde{l}_{ijs} - \max_{k=1,\cdots,p_{ijs}}\{t_{ijks} + o_{ijks}\} - v_{ijs}.$$

The following proposition states the derivation of the expected number of unsuccessful searches within a network of stem-cell donation centres and reformulates the approximated scenario-based capacity planning problem $\text{SCP}_{\text{appx}}$.

**Proposition 2** *Using the upper-bound of the blood waiting time in (2.5), the scenario-based stochastic capacity planning problem ($SCP_{appx}$) for a network of stem-cell donation centres under patient lifetime expectancy following a general distribution becomes an integer optimization model ($SCP_{gdist}$) as follows:*

$$SCP_{gdist}:$$
$$\min_{x_j \in \mathbb{Z}^+} \sum_{s=1}^{S} \omega_s \sum_{j=1}^{J} \sum_{i=1}^{I_{js}} r_{ijs}\left(1 - \int_{l_{ijs}:\ l'_{ijs} > \beta_{ijs},\ x_j \geq \phi(l'_{ijs})} f_{ijs}(l_{ijs})\mathrm{d}l_{ijs}\right)$$
$$+ z_{ijs}\left(1 - \int_{l_{ijs}:\ l''_{ijs} > \beta_{ijs},\ x_j \geq \phi(l''_{ijs})} f_{ijs}(l_{ijs})\mathrm{d}l_{ijs}\right) + (1 - r_{ijs})(1 - z_{ijs}),$$
$$subject\ to \quad \sum_{j=1}^{J} C_j x_j \leq B,$$

$$where \ \ \phi(w) = \frac{\left[-\Gamma_j^{arv}\Gamma_j^{serv} - \sqrt{(1 - \lambda_j^2)(\Gamma_j^{arv}\Gamma_j^{serv})^2 + \frac{4\lambda_j w}{\mu}(4w - \lambda_j(\Gamma_j^{arv})^2) + 4\lambda_j w(\Gamma_j^{serv})^2}\right]^2}{\left[\lambda_j(\Gamma_j^{arv})^2 - 4w\right]^2}.$$

**Proof.** The upper bound $\overline{W}_j(x_j) = \dfrac{\lambda_j(\Gamma_j^{arv} + \Gamma_j^{serv}/\sqrt{x_j})^2}{4\left[1 - \lambda_j/(\mu_j x_j)\right]}$ on the blood waiting time in centre $j$ can be rewritten as $\overline{W}_j(x_j) = \dfrac{\beta_j x_j + \gamma_j \sqrt{x_j} + \eta_j}{x_j - \pi_j}$ by using parameters $\beta_j = \dfrac{(\Gamma_j^{arv})^2\lambda_j}{4}$, $\gamma_j = \dfrac{\Gamma_j^{arv}\Gamma_j^{serv}}{2}$, $\eta_j = \dfrac{(\Gamma_j^{serv})^2\lambda_j}{4}$, and $\pi_j = \dfrac{\lambda_j}{\mu}$. Then,

the expected number of unsuccessful searches becomes

$$
\mathbb{E}[\tilde{y}'_{ijs}(x_j)] = \begin{cases} \Pr\left(\overline{W}_j(x_j) - \tilde{l}_{ijs} + \max_{k=1,\cdots,p_{ijs}} \{t_{ijks} + o_{ijks}\} > 0\right), & \text{if } r_{ijs} = 1, \\[2mm] \Pr\left(\overline{W}_j(x_j) - \tilde{l}_{ijs} + \max_{k=1,\cdots,p_{ijs}} \{t_{ijks} + o_{ijks}\} + v_{ijs} > 0\right), & \text{if } z_{ijs} = 1, \\[2mm] 1, & \text{otherwise.} \end{cases}
$$

(2.7)

From the first two conditions $\overline{W}_j(x_j) - \tilde{l}'_{ijs} > 0$ and $\overline{W}_j(x_j) - \tilde{l}''_{ijs} > 0$ in (2.7), we obtain the following inequalities: $\dfrac{(\beta_j - \tilde{l}'_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l}'_{ijs}\pi_j}{x_j - \pi_j} > 0$ and $\dfrac{(\beta_j - \tilde{l}''_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l}''_{ijs}\pi_j}{x_j - \pi_j} > 0$, respectively. Since $(x_j - \pi_j)$ is always positive due to the traffic intensity condition, we can rewrite (2.7) as follows;

$$
\mathbb{E}[\tilde{y}'_{ijs}(x_j)] = \begin{cases} \Pr\left((\beta_j - \tilde{l}'_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l}'_{ijs}\pi_j > 0\right), & \text{if } r_{ijs} = 1, \\[2mm] \Pr\left((\beta_j - \tilde{l}''_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l}''_{ijs}\pi_j > 0\right), & \text{if } z_{ijs} = 1, \\[2mm] 1, & \text{otherwise.} \end{cases}
$$

(2.8)

Let us define $\xi_j = \sqrt{x_j}$ in order to analyse the first condition in (2.8). In this case, we have a quadratic function $h(\xi_j) = (\beta_j - \tilde{l}'_{ijs})\xi_j^2 + \gamma_j\xi_j + \eta_j + \tilde{l}'_{ijs}\pi_j$.

- If $\beta_j - \tilde{l}'_{ijs} \geq 0$, then the quadratic function is always positive, that is $h(\xi_j) > 0$, since $\gamma_j$, $\xi_j$, $\eta_j$, $\tilde{l}'_{ijs}$, $\pi_j$ are all positive.

- On the other hand, if $\beta_j - \tilde{l}'_{ijs} < 0$, then $h(\xi_j)$ possesses the positive and negative roots (denoted by $\xi^+$ and $\xi^-$, respectively) as $\xi_j^{+,-} = \dfrac{-\gamma_j \pm \sqrt{\gamma_j^2 - 4(\beta_j - \tilde{l}'_{ijs})(\eta_j + \tilde{l}'_{ijs}\pi_j)}}{2(\beta_j - \tilde{l}'_{ijs})}$. Then $h(\xi_j)$ can be written in a factorized form as $h(\xi_j) = (\beta_j - \tilde{l}'_{ijs})(\xi_j - \xi_j^+)(\xi_j - \xi_j^-)$. We can note that $h(\xi_j) > 0$ is satisfied if and only if $\xi_j < \xi_j^+$ which implies that $x_j < (\xi_j^+)^2$. For the square of the positive root, we introduce $\phi(\tilde{l}'_{ijs}) = (\xi_j^+)^2$ that can be

explicitly written as $\phi(\tilde{l'}_{ijs}) = \dfrac{\left(-\gamma_j - \sqrt{\gamma_j^2 - 4(\beta_j - \tilde{l'}_{ijs})(\eta_j + \tilde{l'}_{ijs}\pi_j)}\right)^2}{4(\beta_j - \tilde{l'}_{ijs})^2}$.

Then, $x_j < \phi(\tilde{l'}_{ijs})$.

As a result, $h(\xi_j) > 0$ in the first probability of (2.8) is valid only when $\beta_j - \tilde{l'}_{ijs} \geq 0$ or $\beta_j - \tilde{l'}_{ijs} < 0$ and $x_j < \phi(\tilde{l'}_{ijs})$. Then we can easily show that

$$\Pr\left((\beta_j - \tilde{l'}_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l'}_{ijs}\pi_j > 0\right) = 1 - \Pr\left(\beta_j - \tilde{l'}_{ijs} < 0,\ x_j \geq \phi(\tilde{l'}_{ijs})\right)$$

by using the following relationship between probability functions

$$\Pr\left(\beta_j - \tilde{l'}_{ijs} \geq 0\right) + \Pr\left(\beta_j - \tilde{l'}_{ijs} < 0\right) \cdot \Pr\left(x_j < \phi(\tilde{l'}_{ijs})\right) = 1 - \Pr\left(\beta_j - \tilde{l'}_{ijs} < 0\right) \cdot \Pr\left(x_j \geq \phi(\tilde{l'}_{ijs})\right).$$

By applying the same procedure, equivalent conditions for

$$(\beta_j - \tilde{l''}_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l''}_{ijs}\pi_j > 0,$$

the second probability of (2.8) are obtained. Moreover, we have

$$\Pr\left((\beta_j - \tilde{l''}_{ijs})x_j + \gamma_j\sqrt{x_j} + \eta_j + \tilde{l''}_{ijs}\pi_j > 0\right) = 1 - \Pr\left(\beta_j - \tilde{l''}_{ijs} < 0,\ x_j \geq \phi(\tilde{l''}_{ijs})\right).$$

Then we can compute the expected number of unsuccessful searches in view of $\overline{W}_j(x_j)$ for patient $i$ in centre $j$ under scenario $s$ as follows;

$$\mathbb{E}[\tilde{y'}_{ijs}(x_j)] = \begin{cases} 1 - \Pr\left(\beta_j - \tilde{l'}_{ijs} < 0, x_j \geq \phi(\tilde{l'}_{ijs})\right), & \text{if } r_{ijs} = 1, \\ 1 - \Pr\left(\beta_j - \tilde{l''}_{ijs} < 0, x_j \geq \phi(\tilde{l''}_{ijs})\right), & \text{if } z_{ijs} = 1, \\ 1, & \text{otherwise.} \end{cases} \quad (2.9)$$

that can also be equivalently rewritten as

$$
\mathbb{E}[\tilde{y}'_{ijs}(x_j)] = \left[1 - \Pr\left(\beta_j - \tilde{l}'_{ijs} < 0, x_j \geq \phi(\tilde{l}'_{ijs})\right)\right] r_{ijs} + \\
\left[1 - \Pr\left(\beta_j - \tilde{l}''_{ijs} < 0, x_j \geq \phi(\tilde{l}'''_{ijs})\right)\right] z_{ijs} + (1 - r_{ijs})(1 - z_{ijs}).
$$

Using the probability distribution function, $f_{ijs}(l_{ijs})$, $\mathbb{E}[\tilde{y}'_{ijs}(x_j)]$ can be restated as follows:

$$
\mathbb{E}[\tilde{y}'_{ijs}(x_j)] = \left(1 - \int_{l_{ijs}:\, l'_{ijs} > \beta_{ijs},\, x_j \geq \phi(l'_{ijs})} f_{ijs}(l_{ijs}) \mathrm{d}l_{ijs}\right) r_{ijs} + \\
\left(1 - \int_{l_{ijs}:\, l''_{ijs} > \beta_{ijs},\, x_j \geq \phi(l''_{ijs})} f_{ijs}(l_{ijs}) \mathrm{d}l_{ijs}\right) z_{ijs} + (1 - r_{ijs})(1 - z_{ijs}).
$$

By substituting this into the optimization model $\text{SCP}_{\text{appx}}$, we obtain $\text{SCP}_{\text{gdist}}$ as presented in Proposition 2.

∎

Recall that the model $\text{SCP}_{\text{gdist}}$ is developed under the assumption of general distribution for uncertain lifetime expectancy of patients. As suggested by the World Health Organization (for instance see, Salomon et al. 2001), the uniform (discrete) distribution can be a reasonable assumption for the patients' lifetime expectancy. The following proposition states the derivation of a scenario-based formulation of the stochastic capacity planning problem in view of uniformly distributed random parameters for lifetime expectancy of patients arriving to the stem-cell centres.

Assume that random lifetime expectancy parameters, $\tilde{l}'_{ijs}$ and $\tilde{l}''_{ijs}$, follow a uniform discrete distribution and vary within intervals $[\underline{l}'_{ijs},\ \bar{l}'_{ijs}]$ and $[\underline{l}''_{ijs},\ \bar{l}''_{ijs}]$, respectively. Moreover, let us consider sets $\underline{\Theta} = \{\underline{l}'_{ijs}, \cdots, \bar{l}'_{ijs}\}$ and $\overline{\Theta} = \{\underline{l}''_{ijs}, \cdots, \bar{l}''_{ijs}\}$

consisting of finite number of values taken from the corresponding intervals.

Let $\Pr(l'_{ijs} = w)$ represent the probability of lifetime of patient $i$ admitted to donation center $j$ under scenario $s$ taking value of $w$ within minimum and maximum possible lifetimes that the patient can have. Moreover, let's define an indicator function as $\psi_{wj} = \mathbb{1}(w > \beta_j)$ for $w \in \underline{\Theta} \cup \overline{\Theta} = \{\underline{l}''_{ijs}, \cdots, \overline{l}'_{ijs}\}$, $j = 1, \cdots, J$ and $s = 1, \cdots, S$. Note that a characteristic (indicator) function $\mathbb{1}(*)$ takes 1 if the condition "$*$" holds and 0, otherwise.

**Proposition 3** *The scenario-based capacity planning problem $SCP_{appx}$ for a network of stem-cell donation centres under patient lifetime expectancy following a uniform (discrete) distribution can be formulated as an integer linear optimization model $SCP_{udist}$ as follows:*

$SCP_{udist}$ :
$$
\min_{x_j \in \mathbb{Z}^+} \sum_{s=1}^{S} \omega_s \sum_{j=1}^{J} \sum_{i=1}^{I_{js}} \left( r_{ijs} \sum_{w \in \underline{\Theta}} \frac{1 - \psi_{wj}\tau_{wj}}{\overline{l}'_{ijs} - \underline{l}'_{ijs}} + z_{ijs} \sum_{w \in \overline{\Theta}} \frac{1 - \psi_{wj}\tau_{wj}}{\overline{l}''_{ijs} - \underline{l}''_{ijs}} + (1 - r_{ijs})(1 - z_{ijs}) \right),
$$
$$
subject\ to\ \ \sum_{j=1}^{J} C_j x_j \leq B,
$$
$$
\phi(w) - x_j \leq M(1 - \tau_{wj}),\ w \in \{\underline{l}''_{ijs}, \cdots, \overline{l}'_{ijs}\},\ \forall j, s,\ i = 1, \cdots, I_{js},
$$
$$
\tau_{wj} \in \{0, 1\},\ w \in \{\underline{l}''_{ijs}, \cdots, \overline{l}'_{ijs}\},\ \forall j, s,\ i = 1, \cdots, I_{js},
$$

*where $M$ represents a sufficiently big number.*

**Proof.** Under the uniform (discrete) distribution assumption, the probabilities in (2.9) are computed as

$$
\Pr\left(\tilde{l}'_{ijs} > \beta_j,\ x_j \geq \phi(\tilde{l}'_{ijs})\right) = \sum_{w \in \underline{\Theta}} \Pr(l'_{ijs} = w)\mathbb{1}(w > \beta_j,\ x_j \geq \phi(w)),
$$
$$
= \sum_{w \in \underline{\Theta}} \frac{\mathbb{1}(w > \beta_j,\ x_j \geq \phi(w))}{\overline{l}'_{ijs} - \underline{l}'_{ijs}} = \sum_{w \in \underline{\Theta}} \frac{\psi_{wj}\,\mathbb{1}(x_j \geq \phi(w))}{\overline{l}'_{ijs} - \underline{l}'_{ijs}},
$$

$$(2.10)$$

and

$$\Pr\left(\tilde{l}''_{ijs} > \beta_j, \ x_j \geq \phi(\tilde{l}''_{ijs})\right) = \sum_{w \in \overline{\Theta}} \Pr(l''_{ijs} = w)\mathbb{1}(w > \beta_j, \ x_j \geq \phi(w)),$$

$$= \sum_{w \in \overline{\Theta}} \frac{\mathbb{1}(w > \beta_j, \ x_j \geq \phi(w))}{\overline{l}''_{ijs} - \underline{l}''_{ijs}} = \sum_{w \in \overline{\Theta}} \frac{\psi_{wj} \, \mathbb{1}(x_j \geq \phi(w))}{\overline{l}''_{ijs} - \underline{l}''_{ijs}}.$$

(2.11)

It is worthwhile to emphasize that $\phi(w)$ takes a fixed value for a given $w$. Thus, for a given capacity $x_j$ of center $j = 1, \cdots, J$, the probabilities in (2.10) and (2.11) become deterministic. In order to express $\mathbb{1}(x_j \geq \phi(w))$, we introduce binary variable $\tau_{wj}$ for $w \in \{\underline{l}''_{ijs}, \cdots, \overline{l}'_{ijs}\}$, and patient $i = 1, \cdots, I_{js}$ under scenario $s = 1, \cdots, S$ subject to

$$\tau_{wj} = \begin{cases} 1, & \text{if } x_j \geq \phi(w), \\ 0, & \text{otherwise.} \end{cases}$$

This relationship can be formulated as a set of constraints using the big $M$ approach;

$$\phi(w) - x_j \leq M(1 - \tau_{wj}), \ w \in \{\underline{l}''_{ijs}, \cdots, \overline{l}'_{ijs}\}, \ j = 1, \cdots, J, \ i = 1, \cdots, I_{js}, \ s = 1, \cdots, S.$$

The expected number of unsuccessful searches $\mathbb{E}[\tilde{y}'_{ijs}(x_j)]$ in view of uniformly distributed life expectancy of patients can be computed as

$$\mathbb{E}[\tilde{y}'_{ijs}(x_j)] = \sum_{s=1}^{S} \omega_s \sum_{j=1}^{J} \sum_{i=1}^{I_{js}} \left( r_{ijs} \sum_{w \in \underline{\Theta}} \frac{1 - \psi_{wj}\tau_{wj}}{\overline{l}'_{ijs} - \underline{l}'_{ijs}} + z_{ijs} \sum_{w \in \overline{\Theta}} \frac{1 - \psi_{wj}\tau_{wj}}{\overline{l}''_{ijs} - \underline{l}''_{ijs}} + (1 - r_{ijs})(1 - z_{ijs}) \right).$$

Then the stochastic capacity planning problem SCP$_\text{appx}$ can be reformulated as an integer (linear) programming model SCP$_\text{udist}$ as stated in Proposition 3. $\blacksquare$

## 2.6 Computational Experiments

In this section, we first describe the design and the input data used for the numerical experiments and then present the computational results of the stochastic capacity planning model of a network of stem-cell donation centres.

### 2.6.1 Design of Experiments and Data

We design a series of computational experiments in order to illustrate the performance of the $SCP_{udist}$ model. In particular, we aim to answer the following questions:

- How would a network constructed by stem-cell donation centres with the optimal capacity perform under uncertain real-life environment?

- How does the capacity planning model behave under different size of uncertainty sets of the interarrival and service times?

- How do the model parameters such as budget, demand, and unit capacity cost affect the optimum capacity levels and overall performance of the network?

- What is the impact of the size of a network (i.e. the number of stem-cell donation centres) on the overall performance in terms of the number of successful searches?

The mixed integer (linear) optimization model $SCP_{udist}$ was implemented in IBM ILOG CPLEX and solved by the Cplex solver. All computational experiments were carried out on a laptop with Windows XP operating system, CPU 2.26 GHz Intel Corei5 and 8 Gb of RAM.

In order to illustrate how a network of stem-cell donation centres structured with the optimal capacity performs under real-life conditions as well as to validate the stochastic capacity planning optimization model, we developed a discrete-

event simulation model in MATLAB. The simulation model explicitly performs the queuing activities and the advanced blood testing operations. The success of search operations depends on the waiting time of each blood sample of patients admitted to the donation centre. While the simulation model computes the real waiting time of blood samples using a queuing model, the stochastic optimization approach uses the upper bound approximation of the waiting time for each blood sample.

The simulation model generates the input data using in-sample and out-of-sample simulation approaches for performance comparison purposes. The in-sample data for the scenario-based parameters are randomly generated using specific distributions within the predefined ranges while the values of the deterministic parameters remain the same as introduced in Table 4.2. The results of the optimization model obtained by using the in-sample data are abbreviated as "in-sample optimization". We can report that the CPU time taken to solve the underlying optimization model with the in-sample data is about 15 minutes.

The output of the optimization model is also validated via the simulation model. The optimal capacities of the centres within the network, obtained by solving the optimization model with the in-sample data, are inserted into the simulation model. Then, we run the simulation model again using the data sets generated with in-sample and out-of-sample approaches regarding with stochastic parameters in view of optimal capacities of centers. The performance metrics measuring the expected rate of unsuccessful searches computed by the in-sample and the out-of-sample data are labelled as "in-sample simulation" and "out-of-sample simulation", respectively.

**Input Data**

For the numerical experiments, we consider two different network structures consisting of two and five stem-cell donation centres. The two-center network that consists of stem-cell donation centres located in two cities, namely Istanbul and Ankara in Turkey, is a real-case. The initial data set for the two-center network was gathered from different sources such as published research papers in the literature as well as an expert knowledge. The five-center network is artificially constructed on the basis of the data collected for the two-center network. Table 4.2 shows a description of the input data used for the numerical experiments and the corresponding sources from where the data were obtained.

Table 2.1: Input data for parameters used in the numerical experiments

| Description of Parameters | Value/Range | Source of Data | Distribution |
|---|---|---|---|
| Patient arrival rates for two centres, respectively | 4 & 3.8 days/patient | Expert knowledge | Exponential |
| Probability of finding a perfect match via national and international sources, respectively | 0.12 and 0.4 | Querol et al. (2009a) | Binomial |
| International search duration | [5, 15] weeks | Querol et al. (2009a) | Uniform |
| Travel time of donors (samples) | [1, 3] weeks | Expert knowledge | Uniform |
| Patients' remaining lifetime | [1, 60] weeks | Howard et al. (2008) | Uniform |
| Average service (blood-testing) time | [5, 5.1] days | DYBMS (2015) | Uniform |
| Number of donors found by initial search | [0, 6] donors | Expert knowledge | Uniform |
| Variabilities of interarrival and service times | 9 and 0.015 | Expert knowledge | – |
| Unit capacity cost of a center | $10 | Expert knowledge | – |
| Total weekly budget | $800 per week | User-specific | – |
| Number of scenarios generated for in-sample and out-of-sample experiments, respectively | 200 and 2000 | User-specific | – |

A planning horizon is set for three years where each period corresponds to a time length of one week. We generate $S = 200$ scenarios with equal probabilities (i.e., $\omega_s = 1/S$ for $s = 1, \cdots, S$) as the input to the optimization model and the in-sample experiments. Similarly, we randomly generate 2000 scenarios to be used in the out-of-sample experiments.

The data associated with the scenario-dependent parameters such as *patient arrivals*, *donar travel time*, *international search duration* and *number of suitable donors found by initial search*, are randomly generated by using various

distributions with the estimated moments shown in Table 4.2. For *patient inter-arrival times*, we consider exponential distribution with the arrival rates, 4 and 3.8 days/patient, respectively, for two centres in Turkey (Istanbul Tip Fakultesi Kemik Iligi Bankasi, 2016). The parameters related to *donor travel times, international search duration, service time (blood-testing duration)* and *number of suitable donors* are assumed to follow a uniform distribution.

We assume that *results of national and international searches* follow binomial distributions with the average values equal to *the probabilities of finding a perfect match via national and international sources*, 0.12 and 0.4, respectively. It is worthwhile to emphasise that during the data generation, the international search is never initiated for a patient if his/her national search is successful.

As suggested by Bertsimas and Bandi (2012), the interarrival time variability ($\Gamma^{arrv}$) is set to three times of the standard deviation in the generated interarrival times (3). Similarly, the service time variability ($\Gamma^{serv}$) is set to the three times of the standard deviation in the generated service times (0.005).

Finally, as mentioned before, *remaining patient life-times (weeks)* follow uniform (discrete) distribution in each scenario. Note that the bounds of this distribution may be different for each patient and the real data for these bounds are not available. Thus, we randomly generated the lower and upper bounds for each patient at each scenario assuming that they also follow uniform distributions.

## 2.6.2 Numerical Results

In this section, we present the results of the computational experiments of the capacity planning model. Specifically, we aim to show the performance of the optimization model for different network structures consisting of one, two and five donation centres. We also investigate the effect of the model parameters (such as budget and arrival variability) in the optimal decisions. The other parameters

such as success rates, international search duration and lifetimes do not affect the optimal decisions significantly, thus, the experiments related to these parameters are not presented. Total rate of unsuccessful searches and maximum waiting time in the advanced blood testing queue are used as the performance metrics.

**Model performance:** In order to examine the performance of the optimization model developed in the previous section, we consider two independent stem-cell donation centres with patient arrival rates of 4 and 3.8 days/patient assuming that the two centres are not inter-connected to each other. The weekly budget shown in Table 4.2 are used for both centres. Figure 2.3 displays the relative frequency histograms of the longest (maximum) blood waiting time (left plot) and expected number of successful searches (right plot) using the in-sample and out-of-sample simulation experiments. The *relative frequency* of a performance metric is defined as the ratio of *the real frequency of the corresponding criteria* to *the total number of observations*. Since two centres have very close arrival rates, their performances are almost same. Thus, we only show the results for one centre (with 4 days/patient).

The results of both donation centres obtained by the out-of-sample and in-sample simulation approaches confirm that the upper bound of the waiting time computed for the optimization model (25) is a good approximation to the longest waiting time (22) obtained by the simulation model. In Figure 2.3, we present the frequency of having different longest waiting times and number of successful searches (aggregated in groups of 10) in out-of-sample and in-sample simulations. The frequencies show the rate of observing the corresponding outcome in the simulation runs. Note that the levels of successful searches are grouped and labelled according to the upper-bound of the corresponding category, for example all the values between 171 and 180 are shown under 180. As shown in Figure 2.3, the minimum number of successful searches (169) is very close to the expected

Figure 2.3: Relative frequency histograms for the longest waiting time (left) and number of successful searches (right) obtained at single stem-cell donation centre using the in-sample and out-of-sample simulation approaches.

number of successful searches $(\mathbb{E}[\tilde{y}'_{ijs}(x_j^*)] = 171)$ computed by the optimization model. The optimum capacity of the center (with the larger arrival rate) obtained by solving the optimization model is $x_j^* = 11$. We observe that the in-sample approach provides a slightly better performance than the out-of-sample as the optimization model uses the in-sample data.

Next we design a numerical experiment that was originally motivated by a real situation where the current network of stem-cell donation centres needs to be expanded. For instance, according to a national newspaper (Milliyet, 2016), the Turkish Government aims to increase the number of the stem-cell donation centres in Turkey to improve the current number of transplantations. For this purpose, suppose that the government plans to have at least one center in five different geographical areas of Turkey. Given the new stem-cell donation network, we investigate how the performance of the network would be affected if the number

of centres and total capacity budget are increased.

**Effect of Network Size:** In order to establish the impact of the network size on the optimum capacities as well as the performance metrics, we extend our experiments to consider a five-center network. As mentioned before, the five-center network is artificially constructed on the basis of the real data collected for the two-center network. The centres are assumed to be located in five different areas of Turkey. The interarrival rates to the (artificial) donation centres are determined based on the region populations as 3.2, 3.7, 4, 4, and 4 days/patient, respectively. All parameters (apart from the patient arrival rates) that are input to the optimization model remain the same as specified for the two-center network.



Figure 2.4: Performance comparison of the two-center and five-center networks at varying weekly budget per centre obtained by the optimization model

Figure 2.4 displays the results of the optimization model in terms of the relative rates of the unsuccessful searches (at y-axis) as the weekly budget per center (at x-axis) varying between \$300 and \$700 for the two-center and five-center

61

networks. These results show that the five-center network provides significantly less unsuccessful searches than the two-center network for almost all budget levels. We also observe that the five-center (two-center) network produces almost same level of unsuccessful search rate at weekly budget higher than \$475 (\$500). This is because the budget is already at a very high level and no more improvements can be achieved in the unsuccessful search rate. Note that the optimization model does not take into account the fixed cost of opening a new centre as well as travelling costs of patients that may affect the overall cost of expanding the network structure.

The out-of-sample experiments for the five-center network are designed as follows. For each donation center in the network, we first compute the average number of unsuccessful searches conducted at fixed capacity (that is determined by solving the optimization model using the in-sample data) over all scenarios. The total number of (patient arrivals) unsuccessful searches of the network is then computed as the sum of (patient arrivals) average unsuccessful searches over all patients admitted to all donation centres of the network. This basically implies the objective function value $\left( \sum_{j=1}^{J} \sum_{s=1}^{S} \omega_s \sum_{i=1}^{I_{js}} \mathbb{E}[\tilde{y}'_{ijs}(x_j^*)] \right)$ given the optimal capacity $x_j^*$. The relative rates of unsuccessful searches for the network can be defined as ratio of *the total number of unsuccessful searches of the network* to *the total number of patient arrivals to the network.*

**Sensitivity Analysis of Model Parameters:** We are also concerned with the impact of various model parameters (such as budget, arrival variability and demand) on the optimal capacity decisions and different performance metrics. To examine this impact, we designed a set of controlled experiments where we only change one parameter at a time within a certain range while keeping the other model parameters at their base levels as defined in Table 4.2. We also investigated the model sensitivity towards other parameters such as unit-capacity cost. We

can report that there is no significant change observed in the capacity decisions in different unit-capacity cost levels. Therefore, we only present the computational results related to the weekly budget, arrival variability and demand in this section.

*Impact of Budget:* Figure 2.5 presents the optimum capacities of each centre within the two-center (left panel) and five-center (right panel) networks at various weekly budget levels reflecting different economic conditions. From these results, we observe that the optimal capacities of centres increase as the budget of centres increases, but remain the same after certain budget levels. This confirms that increasing the weekly budget of centres more than a certain level does not improve the overall performance of the network of stem-cell donation centres. On the other hand, there does not exist a feasible capacity solution when the budget levels are less than \$200 and \$1000 for two-center and five-center networks. In addition, the optimal capacity of a centre with a high arrival rate is higher than the capacities of other centres with low arrival rates as expected.

We also display 5% confidence intervals for the total rate of unsuccessful searches obtained by the out-of-sample experiments (using the respective optimum capacities found by solving the optimization model $\text{SCP}_{\text{udist}}$) in Figure 2.5. Total rate of unsuccessful searches monotonically decreases as the value of weekly budget increases.

*Impact of Arrival Variability:* As mentioned before, the variability parameters (denoted by $\Gamma^{arv}$ and $\Gamma^{serv}$) define the conservativeness of the underlying uncertainty sets for the interarrival and service times. In other words, a larger variability corresponds to a more conservative uncertainty set since it covers a larger number of possible realisations. To investigate how the performance metrics and centre capacities change as the level of conservativeness (which mostly depends on the modellers' preference) varies, we solve the optimization model for different arrival variabilities (i.e., $\Gamma^{arv}/\sigma^{arv} = 0, 1, 2, 3, 4$). It is worthwhile

63

Figure 2.5: Impact of weekly budget on the optimal capacities of donation centers and unsuccessful searches for the two-center (left) and five-center (right) networks.

to mention that the service-time variability does not have a significant effect on the performance since it is almost negligible. Notice that the case for $\Gamma^{arv} = 0$ corresponds to optimize based on the expected waiting time of blood samples. In addition, we do not consider the cases where $\Gamma^{arv}/\sigma^{arv} \geq 5$ since the interarrival uncertainty set covers almost all possible realizations of random interarrrival time of patients when $\Gamma^{arv} = 4\sigma^{arrv}$.

For this experiment, the budget levels are set as \$1000/week and \$7200/week for the two-center and five-center networks, respectively, while keeping the other model parameters at their base levels. Figure 2.6 presents the average rate of unsuccessful searches as well as the optimal capacity levels of the donation centres obtained by solving the optimization model at different arrival variabilities for the two-center (left panel) and five-center (right panel) networks. Notice that the capacities of centres 3, 4 and 5 coincide at each value of arrival variability (since their arrival rates are the same).

64

Figure 2.6: Impact of arrival variability on the capacity of centres and the rate of unsuccessful searches in two-center (left) and five-center (right) networks

We also display the results of the out-of-sample experiments in Figure 2.6. As we can see from Figure 2.6, the difference between the rates of unsuccessful searches obtained by the optimization model and the out-of-sample experiments raises up to 30% as the level of conservativeness (variability parameter) increases. This approves that the level of conservativeness plays an important role on the performance metrics and the optimum capacities of centers. In particular, when the variability parameter is set as zero, all blood samples are assumed to face the same average waiting time. In this case, the out-of-sample results provide the highest rate of unsuccessful searches in both network structures. We can therefore conclude that the use of average waiting time in the model leads to a biased objective value.

*Impact of Demand Change:* In the numerical experiments so far, we have assumed that the demand pattern does not change during the planning horizon. In order to investigate the possible impact of demand changes on the performance

65

metric, we consider a case of single donation center with varying demand rates while all other parameters remain the same. The trends in the demand can be observed due to the population increase, the change in population dynamics or the emergence of better therapies replacing stem-cell donation. For this purpose, we conduct out-of-sample experiments while the demand changes with a linear fashion. We consider four cases as the demand increases by 25%, 50%, and 75% as well as decreases by 50% until the end of the next 3 years and compare the average rate of unsuccessful searches.



Figure 2.7: Average rate of unsuccessful searches at varying capacity levels with various demand patterns

Figure 2.7 plots the rate of average unsuccessful searches for all demand patterns at different capacity levels. The results in Figure 2.7 show that the increase in the demand results in a larger effect compared to a decrease with the same rate. In the optimum capacity level (11) of the centre (computed with the base demand rate), when the demand increases by 50% and 75%, the average

unsuccessful search rate increases by around 20% and 50% respectively. This confirms that the decision-maker should increase the capacity of the donation centre significantly in case of a rising demand. By setting up the capacity of the centre higher than 16 does not make much difference in terms of the overall performance.

## 2.7 Conclusions

Stem-cell donation centres serve patients with an urgent need of transplantation. The search process for a suitable stem-cell donor consists of several steps and require time-consuming and expensive advanced blood tests. The capacity for the blood-testing service affects the waiting time to complete the donor search and its success. Besides, several exogenous uncertainties such as patient lifetimes or donor travel times arise during the search. In this research, we develop a scenario-based stochastic model to find the optimum capacities in a network of stem-cell donation centres maximizing the expected number of successful patient searches under a budget restriction. The advanced blood testing in each centre is modelled as a first-come first-served, multi-server queue with unknown service and arrival distributions. The upper-bound of the waiting time in this queue is replaced with a safe approximation. The resulting non-linear integer programming model is reformulated into a linear one.

The computational experiments show that increasing the number of centres within a stem-cell donation network improves the cost-effectiveness, but in contrary the budget increase more than certain amount does not contribute to the network's performance. Moreover, the sensitivity analysis reveal that the variabilities in patient arrivals have a significant impact on the optimum capacities and the search success rates. The capacity decisions made in view of the

average waiting time of blood samples lead to a biased result, especially in high demand scenarios. On the other hand, the worst-case approach for the waiting time permits to take into account the extreme case of uncertain patient arrivals.

# Chapter 3

# Resource Allocation for Healthcare Network with Outsourcing

## 3.1   Introduction

Outsourcing has emerged as a business approach in the service sector over the last few decades. It can be defined as the procurement of goods or services from an external provider under a contract. Complete outsourcing aims to serve all customers through a provider while partial outsourcing (so-called co-sourcing) targets specific customers on the basis of their strategic importance or for geographical reasons. Outsourcing may be preferred as a way of either reducing costs (Johnson, 2008) or increasing value of services (Kakabadse and Kakabadse, 2000). Outsourcing has been practiced in various service and support sectors such as call-centre and housekeeping services. A significant outsourcing trend has recently been observed in healthcare services towards information technology and clinical services, such as anaesthesia, emergency department staffing, dialysis, diagnostic imaging and hospital staffing (Punke, 2013). According to the survey conducted by a US outsourcing company in 2014, around 81% and 90% of

the US community hospitals outsource emergency and anaesthetic care services, respectively (Saunders and Westerink, 2014).

Outsourcing of medical services, so-called healthcare outsourcing, has been significantly rising in other countries as well as in the UK. According to the Centre for Health and Public Interest, total amount spent on healthcare outsourcing from private service providers has increased 50% between 2009-14, and the value of outsourcing contracts in 2014 was £22.6 billion that was a quarter of the entire the UK National Health Service (NHS) budget at that time (CHPI, 2015). Types of contracted clinical services provided by the healthcare contracts vary widely across the country while general practices including surgeries constitute over a third of the total value.

A healthcare outsourcing contract may either be "activity-based" where the provider is paid for each patient served, or "block" where a lump sum of money is paid to the provider for the delivery of services over a fixed period of time (usually one year). There are also contracts combining both types of payment structures. For example, a block contract can be used for a baseline activity while beyond a specific threshold, the payment takes place according to an activity-based contract. The NHS reported that the majority of outsourcing in the UK healthcare sector is based on the block contracts due to the easy payment structure (NHS, 2014). In addition to the payment structure, the contracts also specify target performance levels, expected patient volumes, penalties and validity duration. Moreover, contract parties must agree on the type of patients to be served by the provider. For example, all patients in a specific area may be directly allocated or only some patients with certain medical conditions may be referred to the service provider (Earwicker and Whynes, 1998). The contract design is a time-consuming and costly process (Monitor, 2013). Due to complex negotiations between parties, it is generally issued for at least one year validity period.

In particular, 211 local clinical commissioning groups are responsible for issuing contracts in the interests of the population of their respective regions in the UK (NHS, 2014). These commissioning groups may even co-source with several providers at the same time and develop an *outsourcing network* in which several providers share the same patient population (UK Department of Health, 2014). Capacity planning is an important aspect in the design of healthcare outsourcing networks (Harrogate and Rural District CCG, 2014; NHS Scotland, 2015; Milton Keynes CCG, 2017).

A central health authority such as the NHS in the UK has a fixed budget for outsourcing healthcare services in a network composed of several regions (NHS, 2013). The central health authority must assign patients in different regions to be treated by the contracted providers, and at the same time, must ensure that the service level in the network satisfies performance targets at the global level; for instance, patient waiting and access times. However, outsourcing networks involve several uncertainties such as the number of patient requests and the service durations that the central authority needs to take into account when making the capacity planning decisions. Due to these uncertainties, even for fixed levels of outsourced capacities, calculating the expected waiting times in such a network is challenging. Thus, finding the optimum outsourcing capacities and assigning patients to providers in an outsourcing network is a complex problem that requires rigorous mathematical modelling and appropriate solution approaches.

In this chapter, we consider a healthcare outsourcing network managed by a central authority which has a fixed budget to outsource service capacities from healthcare providers in several regions. We develop a mathematical optimization model to determine the optimum allocation of patients within the network and capacity levels to outsource from the providers. Each service provider is modelled by a first-come-first-served (FCFS) queue, assuming that both arrivals and ser-

vice times follow general distributions. The maximum access time to service in each queue is approximated using a robust optimization framework. The resulting model is a multivariate, non-linear, integer programming problem which is difficult to solve with exact methods. Thus, we introduce an alternating optimization based heuristic to solve the underlying model. The computational experiments are designed to illustrate the performance of the heuristic using the sets of both generated and real data. The numerical results show that the heuristic approach provides almost optimum solutions within a reasonable CPU time. The experiments conducted with the real data suggest that the service performance in a UK healthcare outsourcing network can be improved. Finally, our results show that the structure of the network plays an important role on the overall service performance.

This chapter is organized as follows. The next section provides a literature review regarding capacity planning in service outsourcing, facility and resource allocation problems in healthcare networks and healthcare outsourcing. Section 3.3 describes the details of the underlying capacity planning problem and presents a stochastic programming model. The structural properties of the model are analyzed in Section 3.4 which also introduces an alternating optimization based heuristic to solve the model. Section 3.5 presents the design and the results of the computational experiments.

## 3.2 Literature Review

The capacity planning problem briefly outlined above relates to three streams of the Operations Research literature. The first stream focuses on the service outsourcing from a mathematical modelling perspective while the second stream studies resource and facility allocation problems in healthcare networks. The

third stream is concerned with the research on healthcare outsourcing.

The overall performance in service outsourcing significantly depends on the outsourced capacity levels. Therefore, the capacity planning problems within service outsourcing have been widely studied; for instance see Aksin et al. (2008), Gurvich and Pery (2012), Schrieck et al. (2014), Kocaga et al. (2015), and Liu et al. (2015). Zhou and Ren (2010) provide a comprehensive review of the literature on service outsourcing. The main features of these papers are summarized in Table 3.1.

Service outsourcing can be practiced based on either of two main strategies. In the first strategy, the excess demand is served by the outsourced organisation while the rest is served by the outsourcer. The second strategy is just the opposite of the first one; the excess demand is served by the outsourcer in this case. The capacity planning problems arising within both strategies have been analyzed in the literature. These problems are usually modelled by several approaches such as game theory (Aksin et al., 2008; Liu et al., 2014) and two-stage programming (Kocaga et al., 2015). The game-theoretic models are solved by using an equilibrium analysis (Aksin et al., 2008; Liu et al., 2014) while two-stage programming models are solved by heuristics (Kocaga et al., 2014). Additionally, a common feature in service outsourcing environments is the existence of queues. The queues inherent in the outsourcing problems are modelled by queuing theory (Gurvich and Pery, 2012; Schrieck et al., 2014; Liu et al., 2015). The queuing models are usually solved by approximation rules such as square root staffing rule and Hayward's approximation rule (Gurvich and Pery, 2012; Schrieck et al., 2014).

Among various services, call-centre outsourcing has been main focus of the capacity planning papers in service outsourcing. One of the most comprehensive analysis for capacity planning in call-centre outsourcing is conducted by Aksin et al. (2008). They assume that a contractor with some pricing power offers

capacity- and volume-based contracts to a service provider where the demand is uncertain. They consider different contract schemes, such as subcontracting base demand or demand fluctuations and to find the optimum contract price and service capacities of the contractor and the provider. They use a game-theoretic approach to derive the optimum capacities when the price is fixed, and optimum prices when the capacities are fixed. Rather than developing a queuing model, they provide generic insights, assuming that the model is extended with a queue.

Table 3.1: A review of the literature on capacity planning in service outsourcing

| Research Papers | Modelling Approach | | | Decisions | | Outsourcing | | Solution Approach | |
|---|---|---|---|---|---|---|---|---|---|
| | GT | TS | QT | Capacity | Price | Base | Peak | Exact | Heuristic |
| Aksin et al. (2008) | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | |
| Gurvich and Pery (2012) | | | ✓ | ✓ | | | ✓ | | ✓ |
| Schrieck et al. (2014) | | | ✓ | ✓ | | | ✓ | | ✓ |
| Kocaga et al. (2015) | | ✓ | | ✓ | | | ✓ | | ✓ |
| Liu et al. (2015) | ✓ | ✓ | | ✓ | | Threshold-based | | ✓ | |
| Our approach | Stochastic Programming | | | ✓ | | Demand-based | | | ✓ |

*GT: Game Theory; TS: Two-stage stochastic programming, QT: Queuing theory

Unlike Aksin et al. (2008), several authors (Schrieck et al., 2004; Koacaga et al., 2015; Gurvich and Perry, 2012) consider a queue within the call-centre outsourcing. Schrieck et al. (2014) develop a mathematical model to find the optimum number of staff in a call-centre assuming that the demand is outsourced at busy times. They assume that the probability that a customer waits more than a certain duration should not pass a specific value. They use an extension of the square root staffing rule and Hayward's approximation rule to identify the best staff levels in terms of the loss probability. The first rule uses a policy function based on the square root of the mean arrival rate to calculate the required staffing level. This rule requires that the interarrival and service times follow exponential distributions. The second rule provides approximate performance measures for a queuing system, assuming that arrivals and service times follow general and exponential distributions, respectively.

Similarly, Kocaga et al. (2015) consider the capacity planning problem in a call centre which outsources when the system is too crowded. They develop a two-stage stochastic programming model to find the optimum staffing and outsourcing policy to minimize the total cost. The first-stage decision is the number of staff to employ, while the second-stage decision is the real-time call routing. The service system is modelled as a queue with exponentially distributed interarrival and service times, multiple servers and abandonment. They use square root safety staffing policies combined with routing to solve the resulting model. Gurvich and Perry (2012) model a network of outsourced call-centres as a multi-queue system. They develop several approximation rules to find the greedy staffing and routing policy.

Unlike most of the studies on the service outsourcing, Liu et al. (2015) consider a healthcare outsourcing problem. They develop a mathematical model to find the best mutual referral policy between a community and a city hospital minimizing the cost of outsourcing. Two types of relationship between hospitals, subordination, two integrated hospitals, and subcontracting, separate but collaborating hospitals, are analyzed independently. They consider several contract schemes such as fee-for-service with and without cost sharing. They assume that patients are referred to the contracted hospital according to their medical states. The optimal referral policy is found by using game theory. In summary, our review in service outsourcing indicates that the relevant papers have only considered two-player, game-theoretic relationships, while capacity and patient allocations in healthcare outsourcing networks remain unexplored. Considering the rapid increase in healthcare outsourcing, our research can contribute significantly to both practice and theory in this area.

The other related stream of literature focuses on the facility location/allocation problem with immobile servers, stochastic demand and congestion (see the reviews

by Berman and Krass (2002) and Boffey et al. (2007)). A similar area of study is capacity planning for hospitals with uncertain demand or lengths of stay (see Mousazadeh et al. (2016) for a detailed review). Specifically, the studies in these streams seek to identify the location and capacity of healthcare facilities and allocate uncertain demand to these locations in order to maximize profits. Queuing theory is used by only a few of the papers in this area (Marianov and Serra, 2002; Chao et al., 2003). Marianov and Serra (2002) model the hospitals in a network as M/M/m queues where the objective is to minimize the number of hospitals in the network and the number of staff assigned to each hospital. They assume that the probability that there are more than certain number of waiting patients is limited. The resulting non-linear integer programming model is solved with a heuristic concentration method.

Similarly, Chao et al. (2003) consider a resource allocation problem in a network of hospitals in which each region has a certain rate of switchable and non-switchable patients. A central decision maker allocates the available resources optimally between the hospitals based on the expected patient waiting times. Each hospital is modelled as an M/M/1 queue with exponentially distributed service and interarrival times. Naboureh and Safari (2016) aim to find the optimum location and capacity of a specialized service within a hospital chain, where patients may be diverted to other hospitals at a certain cost. Zhang et al. (2010) study a bi-level problem for capacity and patient allocation to preventive healthcare centres in a network. The lower level problem focuses on the user choice nature of the allocation decisions, while the number, locations, and capacities of the facilities are found by solving the upper level problem. They assume that the patients are not assigned to specific centres, but may choose the facility to be served based on the total expected waiting and travelling time. The objective of the model is to maximize the participation of the population. The lower level problem is

solved with an exact method (gradient projection), while the upper level problem is solved with a heuristic (tabu search).

The underlying capacity planning problem considered in this chapter differs from facility and resource allocation problems in healthcare networks in several ways. Most importantly, we assume that arrival and service times follow general rather than exponential distributions. This assumption forces to consider approximate formulations since the queuing literature mostly deals with exponential service and interarrival times (Bandi and Bertsimas, 2012). Also, we consider maximum rather than average patient waiting times that are not focused much in the queuing theory literature. Finally, we consider the capacity and patient referral decisions simultaneously, unlike the related papers. Therefore, the models developed in the literature are not applicable to the strategic planning problem considered in this chapter.

Several authors concentrate on the design and planning of contracts between medical staff and hospital management by using game theory (Lu and Donaldson, 2000; Fuloria and Zenios, 2001; Lee and Zenios, 2012). Lu and Donaldson (2000) deal with performance-based contracting of medical staff and its effect on the overall clinical outcomes. Fuloria and Zenios (2001) study an outcome-adjusted contracting problem between two parties, where the purchaser seeks to reimburse for optimal treatment types by defining the contract terms according to the observed outcomes. Lee and Zenios (2012) focus on the structure of a principal–agent model in Medicare's dialysis payment system using an empirical method. However, these papers mostly concentrate on the design of contracts rather than the capacity planning problem and its effect on patient access times.

Our contributions to the literature can be summarized in terms of modelling and solution approaches as follows.

- We develop a non-linear integer programming model for the capacity plan-

ning problem of a healthcare outsourcing network. Each provider is modelled as an FCFS queue where the service and interarrival times follow general distributions. The maximum waiting time in each queue is approximated with a robust optimization based approach. The resulting model is a non-linear integer programming model and difficult to solve with exact methods.

- Due to the computational difficulties encountered with the exact methods, we introduce an alternating optimization based heuristic to solve the resulting model. We illustrate the performance of the heuristic through several computational experiments. The numerical results show that the proposed heuristic has a better computational performance than the considered commercial solvers for real-sized instances. Finally, we investigate the impact of the model parameters on the overall performance of a healthcare outsourcing network. The results are used to provide several policy insights for the planning of healthcare outsourcing networks.

## 3.3   Problem Formulation

We consider a patient population residing in several independent regions. Within each region, healthcare providers can supply medical services based on a contract with a central healthcare authority such as the NHS in the UK. A schematic description of a healthcare outsourcing network is presented in Figure 3.1. We assume that there exists no relationship between the service providers apart from offering contract-based services to the same authority. The health authority needs to determine the outsourcing capacity to supply from each service provider at each region in view of the total outsourcing budget. Additionally, the health authority has to allocate the expected patient volume among the contracted providers.

Patient requests directly arriving to the health authority are then referred to the service providers. We do not consider a specific referral process, but assume that a fixed ratio of expected patient volume is allocated to each provider. The providers are required to serve the patients referred to them.



Figure 3.1: A schematic representation of a healthcare outsourcing network

Let's consider a health authority responsible for the patients divided into $R$ regions, labelled as $r \in \{1, \cdots, R\}$. Each region $r$ consists of $n_r$ number of service providers that are labelled as $i = 1, \cdots, n_r$. Let $x_{ir} \in \mathbb{Z}_{\geq 0}$ represent the capacity (in terms of the number of servers or staff) that the central authority outsources from provider $i \in \{1, \cdots, n_r\}$ in region $r$. Let $y_{ir}(x_{ir})$ denote a binary variable representing whether there is a contract between provider $i$ in region $r$ and the central authority. If service capacity is contracted in provider $i$ and region $r$ (i.e. $x_{ir} > 0$), then the provider is outsourced and therefore $y_{ir}(x_{ir}) = 1$. However, if no outsourcing contract exists between the provider and the central authority, in other words no capacity is assigned to the provider (i.e. $x_{ir} = 0$), then $y_{ir}(x_{ir}) = 0$. The maximum capacity that can be outsourced from provider $i \in \{1, \cdots, n_r\}$ in region $r$ is denoted by $C_{ir}$. Thus, the capacity allocated to provider $i$ in region

$r$ cannot exceed the available capacity level. This can be stated as the following constraint:

$$C_{ir}y_{ir}(x_{ir}) \geq x_{ir}. \tag{3.1}$$

Let $p_{ir}$ denote the variable cost depending on the (unit) capacity outsourced from provider $i$ in region $r$. There is also a fixed contracting cost, denoted by $f$, that will be paid for the time spent on negotiations or paperwork when the contract has been signed between the partners. Given the budget level $B'$, we impose the following linear constraint that restricts the total cost of service outsourcing over all regions not to exceed the available budget:

$$\sum_{r=1}^{R}\sum_{i=1}^{n_r} fy_{ir}(x_{ir}) + p_{ir}x_{ir} \leq B'. \tag{3.2}$$

Generally speaking, demand for any service in the network and service durations are not known in advance. Suppose that demand in region $r$ follows a known distribution with average $\lambda_r$ (average number of patents arriving in region $r$ during a day) that can be estimated from historical data. Let $\alpha_{ir} \in [0, 1]$ denote a percentage (or a fraction) of overall demand realised (patient arrivals) in region $r$ to be served by provider $i$. Then, the average number of patients allocated to provider $i$ in region $r$ (denoted by $\lambda_{ir}$) can be computed as $\lambda_{ir} = \lambda_r\alpha_{ir}$. All expected patient demand in region $r$ should be allocated among the providers in this region:

$$\sum_{i=1}^{n_r} \alpha_{ir} = 1, \ r = 1, \cdots, R. \tag{3.3}$$

We model the service system of provider $i$ in region $r$ as an FCFS queue facing no blocking and congestion. The number of servers in each queue is defined by the outsourced capacity $x_{ir}$. The average service rate (average number of patients

served by unit-capacity during a day) of provider $i$ in region $r$ is denoted by $\mu_{ir}$. For a stable queue, the utilization rate (traffic intensity) in the queue should be smaller than 1, that is $\dfrac{\lambda_{ir}}{x_{ir}\mu_{ir}} < 1$ for $i = 1, \cdots, n_r$ and $r = 1, \cdots, R$. In other words, the total service rate ($x_{ir}\mu_{ir}$) should be larger than the total arrival rate ($\lambda_{ir}$) such that the queue does not grow exponentially. Using $\lambda_{ir} = \lambda_r \alpha_{ir}$, we obtain the following condition

$$x_{ir}\mu_{ir} > \lambda_r \alpha_{ir}, \ i = 1, \cdots, n_r, \ r = 1, \cdots, R. \qquad (3.4)$$

The central health authority is responsible for patients to be served within a certain time. However, due to the variations in arrival and service times, the access times of patients for the service can vary. Since each patient is equally important and the worst-case, the patient death, should be avoided as much as possible, we assume that the health authority aims to minimize the worst-case access time within the network (NHS, 2017). Let's represent the maximum patient access time in provider $i$ and region $r$ with $W_{ir}(\alpha_{ir}, x_{ir})$ which depends on the demand allocated to this provider and the outsourced capacity. This dependency is due to that the waiting times in a queue are affected by the number of servers and the arrival rates to this queue. The health authority would like to determine the capacity to outsource and number of patients referred to each provider in parallel to the outsourcing contract, if signed, so that the maximum (patient) access time over all regions is minimized. Then, the capacity planning problem of the central authority can be formulated as:

$$
\text{CAP}: \quad \min_{\alpha_{ir}, x_{ir}} \quad \max_{r=\{1,\cdots,R\}, i=\{1,\cdots,n_r\}} W_{ir}(\alpha_{ir}, x_{ir}),
$$

$$
\text{s. t.} \quad \sum_{r=1}^{R} \sum_{i=1}^{n_r} f y_{ir} + p_{ir} x_{ir} \leq B',
$$

$$
\sum_{i=1}^{n_r} \alpha_{ir} = 1, \ r = 1, \cdots, R,
$$

$$
x_{ir} \mu_{ir} > \alpha_{ir} \lambda_r, \ i = 1, \cdots, n_r, \ r = 1, \cdots, R,
$$

$$
C_{ir} y_{ir} \geq x_{ir}, \ i = 1, \cdots, n_r, \ r = 1, \cdots, R,
$$

$$
x_{ir} \geq y_{ir}, \ i = 1, \cdots, n_r, \ r = 1, \cdots, R,
$$

$$
x_{ir} \in \mathbb{Z}_{\geq 0}, y_{ir} \in \{0, 1\}, \alpha_{ir} \in [0, 1], \ i = 1, \cdots, n_r, \ r = 1, \cdots, R.
$$

We assume that the arrival and service processes follow general distributions in all regions. Our aim by this assumption is to prevent any possible inaccuracy or errors occurring due to an imprecise fitting of data of the underlying distributions. However, it is computationally difficult to derive an exact formulation for the maximum waiting time in a queue with a general arrival distribution (Bandi and Bertsimas, 2012). Therefore, we consider an approximate formulation of the maximum access time in each provider.

**Approximation to Maximum Waiting Time in Queue**

There are various approximate formulations for the maximum waiting time in a queue where the arrival and service times follow general distributions (for instance, see Gupta and Osogami (2011)). However, these approximations may not lead to realistic results due to the underlying assumptions (Bandi and Bertsimas, 2012). As an alternative approach, Bandi and Bertsimas (2012) propose an approximation method based on robust optimization for the maximum waiting time in an FCFS queue. This approach adjusts the conservativeness of the model against the uncertainties in the arrival and service times without assigning any specific distributions to them. Particularly, Bandi and Bertsimas (2012) consider an FCFS queue where the service and interarrival times of the customers are

independent and identically distributed (i.i.d.) random numbers. By using the central limit theorem, they develop uncertainty sets that the service and interarrival times belong to. Readers are referred to Bandi and Bertsimas (2012) for the details and Section 2.5.1 for an overview regarding this approximation. Next, we explain how to apply this approach to obtain the approximate maximum access time of each patient admitted to provider $i$ in region $r$.

Consider the FCFS queue in provider $i$ in region $r$ with $x_{ir}$ number of servers. Let's assume that the interarrival and service times of patients in provider $i$ in region $r$ are i.i.d. and represented as $T_{pir}$ and $Y_{pir}$ for patients $p = 1, \cdots, P_{ir}$, respectively. The means $\mu_{ir}$ and $\lambda_{ir}$ of the random service and interarrival times in provider $i$ in region $r$ are estimated from the generated data. Assume that $T_{pir}$ and $Y_{pir}$ belong to uncertainty sets $U_{ir}^{arrv}$ and $U_{ir}^{serv}$, respectively. Moreover, the sizes of these uncertainty sets are determined by parameters $\Gamma_{ir}^a$ and $\Gamma_{ir}^s$ that basically measure the variability in the interarrival and service times, respectively. They are set by the modeller based on the desired conservativeness of the model against the uncertainties (Bandi and Bertsimas, 2012). As $\Gamma_{ir}^a$ and $\Gamma_{ir}^s$ are higher, the model considers a wider range of realizations for the interarrival and service times i.e. it gets more conservative. The uncertainty set $U_{ir}^{arrv}$ for interarrival times $T_{pir}$ of patients $p = 1, \cdots, P_{ir}$, $i = 1, \cdots, n_r$, and $r = 1, \cdots, R$ is defined as follows:

$$U_{ir}^{arrv} = \left\{ (T_1, T_2, ..., T_{P_{ir}}) \;\middle|\; \frac{\left| \sum_{p=m+1}^{P_{ir}} T_p - \frac{P_{ir}-m}{\lambda_{ir}} \right|}{\sqrt{(P_{ir} - m)}} \leq \Gamma_{ir}^a, \quad \forall m \leq m_0 \right\},$$

where $m_0$ can be set to $P_{ir} - 30$. Similarly, the uncertainty set $U_{ir}^{serv}$ for service times $Y_p$ of samples $p = 1, \cdots, P_{ir}$ is defined as

$$U_{ir}^{serv} = \left\{ (Y_1, Y_2, \cdots, Y_{P_{ir}}) \;\middle|\; \frac{\left| \sum_{p=m+1}^{e_{ir}} Y_{px_{ir}+b} - \frac{e_{ir}-m}{\mu_{ir}} \right|}{\sqrt{(e_{ir} - m)}} \leq \Gamma_{ir}^s, \; \forall m \leq e_{ir} - 1, \; 0 \leq b < x_{ir} \right\},$$

where the service times are computed over the partitions of service times into $x_{ir}$ groups with sizes $e_{ir} = \lfloor P_{ir}/x_{ir} \rfloor$ due to the multiple servers. In view of these uncertainty sets, we can then apply Proposition (2.5.1) to compute the approximate upper bound of the access time in provider $i$ in region $r$, denoted by $\overline{W}_{ir}(\alpha_{ir}, x_{ir})$, for provider $i$ in region $r$ as follows:

$$\overline{W}_{ir}(x_{ir}, \alpha_{ir}) = \frac{\alpha_{ir}\lambda_r \left( \Gamma_{ir}^a + \Gamma_{ir}^s \sqrt{\frac{1}{x_{ir}}} \right)^2}{4 \left( 1 - \frac{\alpha_{ir}\lambda_r}{\mu_{ir}x_{ir}} \right)}. \tag{3.5}$$

Note that the patient interarrival times in provider $i$ and region $r$ depend on the demand allocation to this provider (i.e. the decision variable $\alpha_{ir}$). Indeed, the arrival of patients to provider $i$ in region $r$ is similar to an arrival thinning process with fraction $\alpha_{ir}$. Based on the analysis provided in Bandi et al. (2015), we can formulate the variation in the patient interarrival times in provider $i$ and region $r$ as,

$$\Gamma_{ir}^a = \Gamma_r^a \sqrt{\frac{1}{\alpha_{ir}}}, \quad r = 1, \cdots, R, \ i = 1, \cdots, n_r,$$

where $\Gamma_r^a$ is the variation in the interarrival times in region $r$ and can be driven from the historical interarrival times. For example, it can be set to the double of the standard deviation of the interarrival times to cover a wide range (around 95%) of possible interarrival times, as suggested by Bandi and Bertsimas (2012). Then, the approximated capacity planning problem can be formulated as follows,

$$\text{ACAP}: \quad \min_{\alpha_{ir}, x_{ir}} \quad \max_{r=\{1,\cdots,R\}, i=\{1,\cdots,n_r\}} \frac{\alpha_{ir}\lambda_r \left( \Gamma_r^a/\sqrt{\alpha_{ir}} + \Gamma_{ir}^s/\sqrt{x_{ir}} \right)^2}{4 \left( 1 - \frac{\alpha_{ir}\lambda_r}{\mu_{ir}x_{ir}} \right)},$$

$$\text{s. t.} \quad \text{Constraints } (3.1), \cdots, (3.4),$$

$$x_{ir} \in \mathbb{Z}_{\geq 0}, y_{ir}(x_{ir}) \in \{0, 1\}, \alpha_{ir} \in [0, 1], \ r = 1, \cdots, R, \ i = 1, \cdots, n_r.$$

Note that the size of ACAP depends on the number of regions as well as service providers that basically increase the number of constraints and, most importantly,

the number of integer and binary decision variables. On the other hand, the non-linear approximate upper-bound of access time also adds complexity to the problem.

## 3.4 Structural Properties and Solution Method for the Approximated Model

In this section, we first present an analysis of the model structure and then introduce a solution method based on the alternating optimization. For the sake of computational convenience, we consider a special case of the approximated capacity planning problem under an assumption that the fixed contract costs are removed. In other words, we fix the decision variables representing the choice of service providers $i = 1, \cdots, n_r$ in region $r = 1, \cdots, R$ as $y_{ir}(x_{ir}) = 1$. Then, the fixed contract cost $f \times y_{ir}(x_{ir})$ in the budget constraint becomes a constant and we obtain the modified capacity constraint as

$$C_{ir} \geq x_{ir}, \quad r = 1, \cdots, R, \, i = 1, \cdots, n_r, \tag{3.6}$$

and the budget constraint as

$$\sum_{r=1}^{R} \sum_{i=1}^{n_r} p_{ir} x_{ir} \leq B, \tag{3.7}$$

where $B = B' - \sum_{r=1}^{R} \sum_{i=1}^{n_r} f y_{ir}(x_{ir})$. The approximated problem ACAP becomes:

$$\text{SACAP}: \quad \min_{x_{ir}, \alpha_{ir}} \quad \max_{r=\{1,\cdots,R\}, i=\{1,\cdots,n_r\}} \frac{\alpha_{ir} \lambda_r \left(\Gamma_r^a / \sqrt{\alpha_{ir}} + \Gamma_{ir}^s / \sqrt{x_{ir}}\right)^2}{4\left(1 - \frac{\alpha_{ir}\lambda_r}{\mu_{ir}x_{ir}}\right)},$$

$$\text{subject to} \quad \text{Constraints } (3.3), (3.4), (3.6), (3.7),$$

$$x_{ir} \in \mathbb{Z}_{\geq 0}, \alpha_{ir} \in [0, 1], \, r = 1, \cdots, R, \, i = 1, \cdots, n_r.$$

Although the model is simplified significantly with this modification, the objective function, describing the approximate upper-bound of the waiting time, is still non-linear even with continuous capacity variables. The following proposition states the convexity of this function.

**Proposition 4** *For $x_{ir} \in \mathbb{R}_{\geq 0}$ and $\alpha_{ir} \in [0,1]$, the approximate upper-bound of the waiting time $\overline{W}(x_{ir}, \alpha_{ir})$ for the service provider $i$ and region $r$*

$$\overline{W}(x_{ir}, \alpha_{ir}) = \frac{\alpha_{ir}\lambda_r \left(\Gamma_r^a/\sqrt{\alpha_{ir}} + \Gamma_{ir}^s/\sqrt{x_{ir}}\right)^2}{4\left(1 - \frac{\alpha_{ir}\lambda_r}{\mu_{ir}x_{ir}}\right)}, \tag{3.8}$$

*is neither a convex nor a pseudo-convex function.*

**Proof.** For the sake of convenience, we drop indices $i$ and $r$ in the formulation of $\overline{W}$ in (3.8) which is then defined as

$$f(x, \alpha) = \frac{\alpha\lambda\left(\Gamma^a/\sqrt{\alpha} + \Gamma^s/\sqrt{x}\right)^2}{4\left(1 - \frac{\alpha\lambda}{x\mu}\right)} = \frac{\mu\lambda\left(\Gamma^a\sqrt{x} + \Gamma^s\sqrt{\alpha}\right)^2}{4(\mu x - \alpha\lambda)}.$$

For $y = \sqrt{x}$, $\beta = \sqrt{\alpha}$, and, $m = \lambda/\mu$, we obtain

$$g(y, \beta) = \frac{\lambda\mu\left(\Gamma^a y + \Gamma^s \beta\right)^2}{4\left(y^2\mu - \beta^2\lambda\right)} = \frac{\lambda\left(\Gamma^a y + \Gamma^s \beta\right)^2}{4\left(y^2 - m\beta^2\right)}.$$

Since $\lambda$ is a constant, we drop it from the derivative calculations. The first order partial derivatives of $f(x, \alpha)$ with respect to $x$ and $\alpha$, respectively, can be written as,

$$\frac{\partial f(x, \alpha)}{\partial x} = \frac{1}{2y}\frac{\partial(g(y, \beta))}{\partial y} = -\frac{y\beta^2(m(\Gamma^a)^2 + (\Gamma^s)^2) + m\beta^3\Gamma^s\Gamma^a + y^2\Gamma^a\Gamma^s\beta}{4y(y^2 - m\beta^2)^2},$$

86

and

$$\frac{\partial f(x,\alpha)}{\partial \alpha} = \frac{1}{2\beta}\frac{\partial(g(y,\beta))}{\partial \beta} = \frac{y^2\beta((\Gamma^a)^2 + (\Gamma^s)^2) + \Gamma^s\Gamma^a y^3 + m\Gamma^a\Gamma^s\beta^2 y}{\beta(y^2 - m\beta^2)^2}.$$

Let $H = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$ denote the Hessian matrix of function $f(x,\alpha)$, where the second order derivatives of $f(x,\alpha)$ are denoted by $a$, $b$ and $c$. The second order partial derivatives of $f(x,\alpha)$ can be computed follows:

$$a = \frac{\partial^2 f(x,\alpha)}{\partial x^2} = \frac{\beta[-m^2\Gamma^a\Gamma^s\beta^4 + 6m\Gamma^a\Gamma^s\beta^2 y^2 + 3y^4\Gamma^a\Gamma^s + 4y^3\beta(m(\Gamma^a)^2 + (\Gamma^s)^2)]}{2y^3(y^2 - m\beta^2)^3},$$

$$c = \frac{\partial^2 f(x,\alpha)}{\partial \alpha^2} = \frac{y[4m\beta^3 y(m(\Gamma^a)^2 + (\Gamma^s)^2) + 6my^2\beta^2\Gamma^a\Gamma^s + 3m^2\beta^4\Gamma^a\Gamma^s - \Gamma^a\Gamma^s y^4]}{2\beta^3(y^2 - m\beta^2)^3}, \text{ and,}$$

$$b = \frac{\partial^2 f(x,\alpha)}{\partial x \partial \alpha},$$
$$= \frac{-[6my^2\beta^2\Gamma^a\Gamma^s + 2\beta y^3(m(\Gamma^a)^2 + (\Gamma^s)^2) + \Gamma^a\Gamma^s y^4 + m^2\beta^4\Gamma^a\Gamma^s + 2m\beta^3 y(m(\Gamma^a)^2 + (\Gamma^s)^2)]}{2\beta y(y^2 - m\beta^2)^3}.$$

Following Lau (1978), we need to show that $H$ is a positive semi-definite matrix for $f(x,\alpha)$ to be a convex function. In other words, all principal minors, $a$, $c$, $(ac - b^2)$, should be non-negative. The third principal minor can be rewritten as:
$$ac - b^2 = \frac{a'c' - (b')^2}{2\beta^2 x^2(y^2 - m\beta^2)^6}, \text{ where}$$

$$a' = -m^2\Gamma^a\Gamma^s\beta^4 + 6m\Gamma^a\Gamma^s\beta^2 y^2 + 3y^4\Gamma^a\Gamma^s + 4y^3\beta(m(\Gamma^a)^2 + (\Gamma^s)^2),$$

$$b' = 6my^2\beta^2\Gamma^a\Gamma^s + 2\beta y^3(m(\Gamma^a)^2 + (\Gamma^s)^2) + \Gamma^a\Gamma^s y^4 + m^2\beta^4\Gamma^a\Gamma^s + 2m\beta^3 y(m(\Gamma^a)^2 + (\Gamma^s)^2),$$

$$c' = 4m\beta^3 y(m(\Gamma^a)^2 + (\Gamma^s)^2) + 6my^2\beta^2\Gamma^a\Gamma^s + 3m^2\beta^4\Gamma^a\Gamma^s - \Gamma^a\Gamma^s y^4.$$

Note that since the denumerator is always positive, the sign of the principal minor $(ac - b^2)$ is equal to that of $a'c' - (b')^2$ which can be rewritten as:
$$a'c' - (b')^2 = (\Delta a' + \Delta c')b' + \Delta a'\Delta c', \text{ where,}$$

$$\Delta a' = a' - b',$$

$$= 2\Gamma^a\Gamma^s(y^4 - m^2\beta^4) + 2\beta y(m(\Gamma^a)^2 + (\Gamma^s)^2)(y^2 - m\beta^2),$$

and

$$\Delta c' = c' - b',$$

$$= 2\Gamma^a\Gamma^s(m^2\beta^4 - y^4) + 2\beta y(m(\Gamma^a)^2 + (\Gamma^s)^2)(m\beta^2 - y^2).$$

Note that $\Delta a' + \Delta c' = 0$ and $\Delta a' > 0$, $\Delta c' < 0$ due to the traffic intensity condition (i.e. $y^2 - m\beta^2 > 0$). This leads to the third principal minor to become $a'c' - (b')^2 \leq 0$. This shows that the approximate upper-bound of waiting time is not a convex function even with relaxed capacity variables.

For pseudo-convexity, we need to show that all the leading principal minors of the bordered Hessian matrix (denoted as $H^p$) of the approximate upper-bound of waiting time must be negative (Crouzeix and Ferland, 1982; Avriel and Schaible, 1978). For $d = \frac{\partial f(x,\alpha)}{\partial x}$ and $e = \frac{\partial f(x,\alpha)}{\partial \alpha}$, the bordered Hessian matrix is denoted by

$$H^p = \begin{bmatrix} 0 & d & e \\ d & a & b \\ e & b & c \end{bmatrix}.$$

The first leading principal minor of $H^p$, $-(d)^2$, is always negative. The second leading principal minor can be written as,

$$-d(dc - eb) + e(db - ea). \tag{3.9}$$

For the simplicity purposes, let's denote the denumerator of $d$ and $e$ by $d' = y\beta^2(m(\Gamma^a)^2 + (\Gamma^s)^2) + y^2\Gamma^a\Gamma^s\beta + m\beta^3\Gamma^a\Gamma^s$ and $e' = y^2\beta(m(\Gamma^a)^2 + (\Gamma^s)^2) + \Gamma^a\Gamma^s y^3 + m\Gamma^a\Gamma^s\beta^2 y$, respectively. By using the relationship $a' + c' = 2b'$, we can see that the sign of (3.9) is equivalent to that of

$(c'yd' - a'\beta e')(yd' - \beta e')$. Note that $(yd' - \beta e')$ is always zero. Thus, the principal minor is non-negative. This shows that the approximate upper-bound of the access time for provider $i$ in region $r$ (3.8) is not a pseudo-convex function. ∎

The convexity analysis shows that even the relaxed version of the problem is non-convex. However, it is crucial to find the best solution close to the global one. For large size non-convex problems, finding the global optima is computationally expensive. In this chapter, we introduce an alternating optimization method for solving the capacity planning of a healthcare outsourcing network.

In the alternating optimization method, the underlying problem is decomposed into several subproblems with respect to the decision variables. Then these subproblems are iteratively solved to find the optimum solution of the original problem. Initially, the algorithm requires to identify a feasible solution for one set of decision variables. This solution is used as the input to the other subproblem that is solved to optimality. It is proved that the alternating optimization algorithm is convergent when the variables are partitioned into two sets (Bezdek and Hathaway, 2002). Next we consider the optimization problem, SACAP, and apply an alternating optimization method to find the optimal capacity planning strategy for a healthcare outsourcing network. We first transform SACAP by introducing a positive variable $z$ representing the maximum waiting time. Thus, the inner maximization problem can be written as a constraint as follows:

$$
\begin{aligned}
\text{SCPZ}: \quad & \min_{z, \alpha_{ir}, x_{ir}} && z, \\
& \text{subject to} && z \geq \overline{W}(\alpha_{ir}, x_{ir}), \ i = 1, \cdots, n_r, r = 1, \cdots, R, \\
& && \text{Constraints } (3.3), (3.4), (3.6), (3.7), \\
& && z \in \mathbb{R}_{\geq 0}, x_{ir} \in \mathbb{Z}_{\geq 0}, \alpha_{ir} \in [0, 1], \ r = 1, \cdots, R, \ i = 1, \cdots, n_r.
\end{aligned}
$$

This model can be decomposed into two subproblems in each of which the capacity and demand allocation decisions ($x_{ir}$ and $\alpha_{ir}$) are fixed, respectively. The following propositions display how to obtain these subproblems which are to be

solved iteratively at each iteration of the alternating optimization algorithm.

**Proposition 5** *For fixed capacity level $\bar{x}_{ir}$ for provider $i$ in region $r$, the optimization model, SCPZ, becomes:*

$SCPZ(\bar{x}_{ir})$:

$$\min_{z, \alpha_{ir} \in [0,1]} z,$$

$$\text{s. t. } 4z \geq \lambda_r (\Gamma_r^a)^2, \ r = 1, \cdots, R, \tag{3.10}$$

$$\left[ \frac{-\sqrt{\bar{x}_{ir}} \mu_{ir} \Gamma_r^a \Gamma_{ir}^s \lambda_r + 2\sqrt{\lambda_r^2 (\Gamma_r^a)^2 z \mu_{ir} \bar{x}_{ir} - \bar{x}_{ir} \mu_{ir}^2 (\Gamma_{ir}^s)^2 z \lambda_r - 4z^2 \bar{x}_{ir} \mu_{ir} \lambda_r}}{\lambda_r (\mu_{ir} (\Gamma_{ir}^s)^2 + 4z)} \right]^2 \geq \alpha_{ir},$$

$$r = 1, \cdots, R, i = 1, \cdots, n_r. \tag{3.11}$$

$$\sum_{i=1}^{n_r} \alpha_{ir} = 1, \ r = 1, \cdots, R,$$

$$\bar{x}_{ir} \mu_{ir} > \alpha_{ir} \lambda_r, \ r = 1, \cdots, R, \ i = 1, \cdots, n_r.$$

**Proof.** By substituting $x_{ir} = \bar{x}_{ir}$ in constraint, $z \geq \overline{W}(\alpha_{ir}, \bar{x}_{ir})$, we obtain,

$$\overline{W}(\alpha_{ir}, \bar{x}_{ir}) - z = \frac{\mu_{ir} \bar{x}_{ir} \alpha_{ir} \lambda_r \left( \Gamma_r^a \sqrt{\frac{1}{\alpha_{ir}}} + \Gamma_{ir}^s \sqrt{\frac{1}{\bar{x}_{ir}}} \right)^2 - 4z(\mu_{ir} \bar{x}_{ir} - \alpha_{ir} \lambda_r)}{4\left( \mu_{ir} \bar{x}_{ir} - \alpha_{ir} \lambda_r \right)} \leq 0. \tag{3.12}$$

Note that the denominator of (3.12) is always positive due to the traffic intensity condition. For $\omega_{ir} = \sqrt{\alpha_{ir}}$, $u_{ir} = \mu_{ir} \lambda_r (\Gamma_{ir}^s)^2 + 4\lambda_r z$, $v_{ir} = 2\sqrt{\bar{x}_{ir}} \mu_{ir} \Gamma_r^a \Gamma_{ir}^s$, and $h_{ir} = \mu_{ir} \lambda_r (\Gamma_r^a)^2 \bar{x}_{ir} - 4\bar{x}_{ir} \mu_{ir} z$, the numerator of (3.12) can be written in a quadratic form: $A(\omega_{ir}) = u_{ir} \omega_{ir}^2 + v_{ir} \omega_{ir} + h_{ir}$. Notice that $u_{ir}$ and $v_{ir}$ are always positive and depending on the sign of $h_{ir}$, $A(\omega_{ir})$ has either two negative roots ($h_{ir} > 0$) or one positive and one negative root ($h_{ir} \leq 0$). We also know that $\omega_{ir} = \sqrt{\alpha_{ir}}$ is always positive. Thus, we can only consider $h_{ir} \leq 0 \Leftrightarrow \lambda_r (\Gamma_r^a)^2 \leq 4z$ that is equal to constraint (3.10). $A(\omega_{ir})$ is not positive if $\omega_{ir}$ is smaller than or

equal to the positive root. This leads to

$$\left[\frac{-\sqrt{x_{ir}}\mu_{ir}\Gamma_r^a\Gamma_{ir}^s\lambda_r + 2\sqrt{\lambda_r^2(\Gamma_r^a)^2 z\mu_{ir}x_{ir} - x_{ir}(\mu_{ir})^2(\Gamma_{ir}^s)^2 z\lambda_r - 4z^2 x_{ir}\mu_{ir}\lambda_r}}{\lambda_r\big(\mu_{ir}(\Gamma_{ir}^s)^2 + 4z\big)}\right]^2 \geq \alpha_{ir},$$

which is equivalent to constraint (3.11). ∎

**Proposition 6** *For fixed values of patient allocation decisions, $\bar{\alpha}_{ir}$, for service provider i in region r, the optimization model SCPZ becomes:*

$SCPZ(\bar{\alpha}_{ir})$:

$$\min_{z, x_{ir} \in \mathbb{Z}_{\geq 0}} z,$$

s. t.,

$$4z \geq \lambda_r(\Gamma_r^a)^2, \ r = 1, \cdots, R,$$

$$x_{ir} \geq \left[\frac{-\sqrt{\bar{\alpha}_{ir}}(\mu_{ir}\Gamma_r^a\Gamma_{ir}^s\lambda_r + 2\sqrt{-\lambda_r^2\mu_{ir}(\Gamma_r^a)^2 z + \mu_{ir}^2\ \lambda_r(\Gamma_{ir}^s)^2 + 4\mu_{ir}\lambda_r z^2})}{\mu_{ir}(\lambda_r(\Gamma_r^a)^2 - 4z)}\right]^2,$$

$$r = 1, \cdots, R, \ i = 1, \cdots, n_r, \hspace{3cm} (3.13)$$

$$x_{ir}\mu_{ir} > \bar{\alpha}_{ir}\lambda_r, \ r = 1, \cdots, R, \ i = 1, \cdots, n_r,$$

$$x_{ir} \leq C_{ir}, \ r = 1, \cdots, R, \ i = 1, \cdots, n_r,$$

$$\sum_{r=1}^{R}\sum_{i=1}^{n_r} p_{ir}x_{ir} \leq B.$$

**Proof.** By substituting $\alpha_{ir} = \bar{\alpha}_{ir}$ in constraint, $z \geq \overline{W}(\bar{\alpha}_{ir}, x_{ir})$, we obtain,

$$\overline{W}(\bar{\alpha}_{ir}, x_{ir}) - z = \frac{\mu_{ir}x_{ir}\bar{\alpha}_{ir}\lambda_r\left(\Gamma_r^a\sqrt{\frac{1}{\bar{\alpha}_{ir}}} + \Gamma_r^s\sqrt{\frac{1}{x_{ir}}}\right)^2 - 4z(\mu_{ir}x_{ir} - \bar{\alpha}_{ir}\lambda_r)}{4\left(\mu_{ir}x_{ir} - \bar{\alpha}_{ir}\lambda_r\right)} \leq 0. \hspace{0.5cm} (3.14)$$

We observe that the denominator of (3.14) is always positive due to the traffic intensity condition. For $k_{ir} = \sqrt{x_{ir}}$, $\eta_{ir} = \mu_{ir}\lambda_r(\Gamma_r^a)^2 - 4\mu_{ir}z$, $\tau_{ir} = 2\mu_{ir}\Gamma_r^a\Gamma_{ir}^s\lambda_r\sqrt{\bar{\alpha}_{ir}}$, and $\kappa_{ir} = \bar{\alpha}_{ir}\mu_{ir}\lambda_r(\Gamma_{ir}^s)^2 + 4\bar{\alpha}_{ir}\lambda_r z$, its numerator can be written in a quadratic

form as $P(k_{ir}) = \eta_{ir}k_{ir}^2 + \tau_{ir}k_{ir} + \kappa_{ir}$. Notice that both $\tau_{ir}$ and $\kappa_{ir}$ are always positive. The increasing or decreasing pattern of function $P(k_{ir})$ depends on the sign of $\eta_{ir}$. It has either two negative roots (when $\eta_{ir} > 0$) or one positive and one negative root (when $\eta_{ir} \leq 0$). Also, note that $k_{ir} = \sqrt{x}_{ir}$ is always non-negative. Thus, $\eta_{ir} > 0$ is not feasible, and $\lambda_r(\Gamma_r^a)^2 \leq 4z$ should always be satisfied. In this case, since $\eta_{ir} \leq 0$, $P(k_{ir})$ is not positive when $k_{ir}$ is larger than or equal to the positive root that leads to constraint (3.13):

$$x_{ir} \geq \left[ \frac{-\sqrt{\bar{\alpha}_{ir}}\Big( \mu_{ir}\Gamma_r^a\Gamma_{ir}^s\lambda_r + 2\sqrt{-\lambda_r^2\mu_{ir}(\Gamma_r^a)^2 z + \mu_{ir}^2\ \lambda_r(\Gamma_{ir}^s)^2 + 4\mu_{ir}\lambda_r z^2}\Big)\right]^2}{\left[\mu_{ir}(\lambda_r(\Gamma_r^a)^2 - 4z)\right]^2}.$$

∎

Note that the objective functions of both subproblems SCPZ($\bar{\alpha}_{ir}$) and SCPZ($\bar{x}_{ir}$) are monotonic. Also, for a fixed value of $z$, they become linear programming problems. Thus, we consider a section search method to solve these subproblems.

A section search method does not require differentiation and converges to the optimum solution when the objective function is monotonic (Burden and Faires, 1993). In particular, a section search method narrows the feasible region of the variables by systematically comparing the objective function values. The most widely used section search method is bisection search (Waeber et al., 2013). The bisection search method divides the feasible region into two halves at each iteration. By using a bisection search, we can obtain the unique optimum of subproblems SCPZ($\bar{\alpha}_{ir}$) and SCPZ($\bar{x}_{ir}$). Bertsekas (1999) showed that when each subproblem in an alternating optimization algorithm attains a unique minimum, the convergence point of the algorithm is a stationary point.

In this chapter, we combine the alternating optimization based heuristic with the bisection search to solve the optimization model SACAP. The pseudo

code of the alternating optimization algorithm is presented in Algorithm 1. At each iteration, the algorithm solves two subproblems SCPZ($\bar{\alpha}_{ir}$) and SCPZ($\bar{x}_{ir}$) with a bisection search using linear programming (e.g. simplex) and integer linear programming (branch-and-bound) methods, respectively. Then, it compares the objective function values computed at the optimal solution of the subproblems. The heuristic stops when the difference between these objective values is smaller than a certain tolerance level (e.g. $\delta = 10^{-8}$).

The main steps of the algorithm are described as follows. In the initial step, the heuristic requires to find a feasible solution at which the objective function value is represented by $W'$. Note that this value can be set to a very large number. We also determine the initial feasible solution $\bar{x}_{ir}$ by setting the minimum capacity levels that satisfy the traffic intensity constraints, i.e. $\bar{x}_{ir} = \lambda_{ir}/\mu_{ir}$, assuming that all the service providers in region $r$ share total number of patients $\lambda_r$ arriving to region $r$ equally, that is $\lambda_{ir} = \lambda_r/n_r$ for $i = 1, \cdots, n_r$ and $r = 1, \cdots, R$. The fixed capacity levels, $\bar{x}_{ir}$, and the objective function value $z = W'$ are given as inputs to subproblem SCPZ($\bar{x}_{ir}$). Then, SCPZ($\bar{x}_{ir}$) is solved by narrowing the distance between a feasible and infeasible objective values, denoted by $z_f^d$ and $z_{inf}^d$, respectively, at iteration $d$. If this subproblem has a feasible solution, then the distance between $z_f^d$ and $z_{inf}^d$ is halved and the problem is solved again with the new objective value, while the iteration counter $d$ is increased. This process goes until the difference between $z_f^d$ and $z_{inf}^d$ is less than a tolerance level $\delta$. Then, the solution obtained in the last iteration and the objective value $z^{fc}$ are given as the inputs to SCPZ($\bar{\alpha}_{ir}$) where the objective value $z = z^{fc}$. The same process repeats for this subproblem while the interval for the objective value is halved at each iteration $m$ until there is no further change. If the difference between the last objective values obtained in two subproblems, $z_{fa}^*$ and $z_{fc}^*$ is smaller than $\delta$, then the algorithm stops. Otherwise, the optimum capacity levels obtained by

solving SCPZ($\bar{\alpha}_{ir}$) and the maximum of two objective values $\max\{z_{fc}^*, z_{fa}^*\}$ are given as the inputs to SCPZ($\bar{x}_{ir}$) and the whole process is repeated again.

Note that each subproblem is solved with bisection search that has computational complexity of $O(\log h)$ where $h$ is the size of feasible space, i.e. $[0, W']$. There is no worst-case iteration complexity analysis for alternating optimization when the objective function is not convex. The computation time of the heuristic depends on the starting point (the initial capacity set) as well as the convergence rates i.e. how fast the subproblems converge to a solution in terms of the number of iterations. We analyze the computational performance of the heuristic in more detail in the next section.

## 3.5   Computational Experiments

This section is concerned with the design and data structure used for the numerical experiments and also presents the computational results obtained by solving the capacity planning problem using the alternating optimization algorithm. With these experiments, we also aim to display the sensitivity of the solution towards the model parameters. Specifically, we would like to answer the following questions during the computational experiments:

- How does the alternating optimization algorithm perform with respect to other approaches?

- How does the factors such as budget and network size affect the performance of the heuristic?

- How does the service performance of the outsourcing network change with respect to the model parameters?

94

---

**Algorithm 1** Alternating Optimization combined with Bisection Search

---

Initialize $\delta$ (tolerance), $z_{fa}^* = W'$ and $z_{fc}^* = 0$.

Determine feasible $x_{ir}$ values $(\overline{x}_{ir})$ for $r = 1, \cdots, R$ and $i = 1, \cdots, n_r$.

**while** $|z_{fa}^* - z_{fc}^*| > \delta$, **do**

    Set $z_f^1 = \max\{z_{fa}^*, z_{fc}^*\}$, $z_{inf}^1 = 0$, and $d, m = 1$.

    **while** $|z_f^d - z_{inf}^d| > \delta$ **do**

        **if** the following optimization model has a feasible solution,

$$\min_{\alpha_{ir} \in [0,1]} z^d = (z_{inf}^n + z_f^d)/2,$$

$$\text{s. t. } z^d \geq \overline{W}_{ir}(\alpha_{ir}, \overline{x}_{ir}), \qquad r = 1, \cdots, R, i = 1, \cdots, n_r,$$

$$\sum_{i=1}^{n_r} \alpha_{ir} = 1, \qquad\qquad\qquad r = 1, \cdots, R,$$

$$\overline{x}_{ir}\mu_{ir} > \alpha_{ir}\lambda_r, \qquad\qquad r = 1, \cdots, R, i = 1, \cdots, n_r,$$

        **then**

          $z_f^{d+1} = (z_{inf}^d + z_f^d)/2,$

        **else**

          $z_{inf}^{d+1} = (z_{inf}^d + z_f^d)/2.$

        **end if**

        $d := d + 1.$

    **end while**

    Set $z_{fa}^* = z_f^{d-1}$, $\overline{\alpha}_{ir} = \alpha*_{ir}$, the optimum decision allocation levels at $d - 1$, $z_{inf}^1 = 0$ and $m = 1$.

    **while** $|z_f^m - z_{inf}^m| > \delta$. **do**

        **if** the following optimization model has a feasible solution,

$$\min_{x_{ir} \in \mathbb{Z}_{\geq 0}} z^m = (z_{inf}^m + z_f^m)/2,$$

$$\text{s. t. } z^m \geq \overline{W}_{ir}(\overline{\alpha}_{ir}, x_{ir}), \qquad r = 1, \cdots, R, i = 1, \cdots, n_r,$$

$$x_{ir}\mu_{ir} > \overline{\alpha}_{ir}\lambda_r, \qquad\qquad r = 1, \cdots, R, i = 1, \cdots, n_r,$$

$$C_{ir} \geq x_{ir}, \qquad\qquad\qquad r = 1, \cdots, R, i = 1, \cdots, n_r,$$

$$B \geq \sum_{r=1}^{R}\sum_{i=1}^{n_r} p_{ir}x_{ir},$$

        **then**

          $z_f^{m+1} = (z_{inf}^m + z_f^m)/2,$

        **else**

          $z_{inf}^{m+1} = (z_{inf}^m + z_f^m)/2.$

        **end if**

        $m := m + 1.$

    **end while**

    **return** Set $z_{fc}^* = z_f^{m-1}$ and $\overline{x}_{ir} = x*_{ir}$, the optimum capacity levels obtained at $m - 1$.

**end while**

---

- How does the capacity planning model perform in terms of the patient access time when real data are considered?

The alternating optimization algorithm was implemented in MATLAB while the integer linear programming problem (the second subproblem) is solved by Cplex (branch-and-bound). All computational experiments were carried out on a laptop with Windows XP operating system, CPU 2.26 GHz Intel Corei5 and 8 Gb of RAM.

## 3.5.1 Design of Experiments and Data

We design three sets of computational experiments in order to illustrate the performance of model SACAP developed in the previous section. Initially, we investigate computational performance of the proposed heuristic and compare it with the commercially available exact and local solvers. We randomly generate a data set for a network consisting of three regions with 10 service providers. The specifications of this data set, so called small network, is presented in Table 3.2. In order to illustrate effect of different network structures on the performance of the algorithm, we also generate other artificial data sets with more regions and providers (labelled as R and P, respectively). We abbreviate these networks as N(R, P) in the rest of the chapter. In particular, we consider different network structures with 3, 6, and 12 regions and 10, 20 and 40 providers (that are abbreviated as N(3, 10), N(6, 20) and N(12, 40), respectively). These networks are constructed as duplication of identical small networks. For example, the network N(12, 40) is combination of four small networks of N(3, 10), where the data related to the providers and regions remain the same as presented in Table 3.2.

In the second part of the experiments, we conduct a sensitivity analysis by using the artificial data sets to investigate effect of the model parameters on the results. We also investigate impact of real data setting on the performance of the

capacity planning problem of a healthcare outsourcing network. The real data are obtained from the online NHS sources (Monitor, 2016). Since the real data do not differentiate between the unit-prices of providers, we generate the artificial data sets, N(R, P).

We assume that the interarrival and service times follow exponential distributions with the rates presented in Table 3.2. The interarrival and service times are generated by simulation. Then, the variation parameters are set to the double of the standard deviations of the corresponding simulated data, as suggested by Bandi and Bertsimas (2012).

Table 3.2: Description of the data set specified for a small network N(3, 10)

| Parameter | Level |
|---|---|
| Budget | 500$ |
| Number of providers in the regions | 3, 4, 3 |
| Average patient arrival rates for each region (patient/day) | 24, 24, 18 |
| Mean service times for each region (patient/day) | 1.5, 1.5, 1.5 |
| Available capacities ($C_{ir}$) of service providers at each region | [100, 45, 15], [100, 45, 35, 35], [135, 35, 25] |
| Unit-capacity costs for each provider at each region ($) | [1.2, 1, 1], [1.2, 1, 1, 2], [1.1, 1, 1] |

## 3.5.2   Computational Results

**Computational Performance of the Heuristic**

In order to illustrate the performance of the heuristic, we choose a local solver (Bonmin in GAMS) and a global solver (Couenne in Julia) as benchmarks. There are few commercial solvers available for non-linear integer programming problems. Bonmin and Couenne are chosen based on their computational performance. The solution quality of these solvers and the heuristic is measured in terms of the optimality gap and the CPU time taken to find a solution. To investigate the computational performance of the heuristic, we consider different network structures, i.e. number of providers and regions, that affect the problem size. Note that the global solver cannot find the optimum for some instances and had to be

stopped at 20,000th second. Thus, the gap (i.e. the normalized difference) between local and global solutions is not always the optimality gap and, therefore, named as 'lower bound gap'.

To test effect of the integrality constraints, we also solve the relaxed version of model SACAP (by fixing the capacity decisions, $x_{ir}$, to be continuous) for all solvers. In this case, the lower bound gap obtained by the heuristic is almost negligible, e.g. $6.3e^{-4}$ for the network N(12, 40). Thus, in the rest of the experiments, the initial capacity set given as input to the heuristic, $\bar{x}_{ir}$ for $i = 1, \cdots, n_r$ and $r = 1, \cdots, R$, is found by solving the relaxed version of the corresponding problem instance with the heuristic.

Table 3.3: Impact of network structure on performance of solution methods

| Network structure: N(region, provider) | N(3, 10) | N(6, 20) | N(12, 40) |
|---|---|---|---|
| *Problem Size* | | | |
| Number of constraints | 24 | 47 | 93 |
| Number of continuous variables | 10 | 20 | 40 |
| Number of integer variables | 10 | 20 | 40 |
| *Lower bound gap (%)* | | | |
| Local Solver | 0.09 | 0.09 | 0.01 |
| Heuristic | 0.12 | 0.2 | 0.4 |
| *CPU time (seconds)* | | | |
| Local Solver | 0.35 | 55 | 281 |
| Heuristic | 3.6 | 6.3 | 9.22 |
| Global Solver | 81.975 | > 20,000 | >20,000 |

Table 3.3 shows the problem sizes, computation times and lower bound gap for different network structures and solvers. The computation time of the global solver increases exponentially with respect to the problem size. For example, in the networks N(6, 20) and N(12, 40), the global solver is stopped after 20,000 seconds when the optimality gap reported by the solver, between the upper and lower bounds, was negligible ($2.63 \times 10^{-5}$). Moreover, the heuristic and the local solver display similar performance pattern for the network N(3, 10). On the other hand, for the networks N(6, 20) and N(12, 40), the heuristic has a much shorter

computation time but a larger gap than the local solver. However, occasionally, the local solver reported the function evaluation or infeasibility errors while the heuristic always manages to find a solution. Also, we should note that the lower bound gap of the heuristic is almost negligible in all network structures.

From the computational experiments, we observe that the choice of model and input parameters plays an important role on the performance of the heuristic. Therefore, we test the heuristic for different levels of the parameters such as the tolerance of the heuristic, budget, and available capacities. In this section, we only report the numerical results of experiments obtained with the parameters such as tolerance of the heuristic (denoted by $\delta$ in Algorithm 1) and budget, that show significant impact on the heuristic performance. We use the network N(12, 40) with one instance in these experiments.

As Table 3.4 shows, the lower bound gap (the CPU time) obtained by 0.01 tolerance is smaller (higher) than the gap (the CPU time) obtained at the levels 0.1 and 0.05. We also observe that as the tolerance is decreased further from 0.01, the heuristic's performance does not change significantly. Thus, the tolerance is set to 0.01 in the rest of the experiments.

Table 3.4: Impact of the tolerance level on the heuristic's performance obtained using the network N(12, 40)

| Tolerance | 0.1 | 0.05 | 0.01 |
|---|---|---|---|
| Lower bound gap (%) | 0.8 | 0.6 | 0.2 |
| CPU time (seconds) | 3.98 | 6.1 | 9.22 |

To investigate the effect of budget on the heuristic's performance, we consider tight (900$) and loose (2100$) budget levels for the network N(12, 40). Note that the level of the budget ($B$) influences the computation time of the heuristic since it determines the search space for the feasible solutions. As displayed in Table 3.5, with the tight budget, the computation time of the local solver increases

99

considerably (cannot reach to a solution and stopped at 20.000th seconds). On the other hand, the tight budget does not affect the computation time of the heuristic significantly but increases its the lower bound gap.

Table 3.5: Performance comparison of different solution methods using the outsourcing network N(12, 40) and different budget levels

| Budget | 900$ | 2100$ |
|---|---|---|
| *Lower bound gap (%)* | | |
| Heuristic | 2 | 0.4 |
| Local Solver | 2* | 0.01 |
| *CPU time (seconds)* | | |
| Heuristic | 12.1 | 9.22 |
| Local Solver | >20.000 | 281 |
| Global Solver | >20.000 | >20.000 |

*stopped

In summary, we can conclude that the heuristic performs better than the local and global solvers in terms of the computation time and lower bound gap, especially when the budget constraint is tight. Thus, we use the heuristic to obtain the results in the rest of the computational experiments.

**Sensitivity Analysis**

In order to investigate how the model parameters affect the capacity planning strategies as well as the performance metric (patient access time) in a network, we design controlled experiments where one parameter is varied within a certain interval while other model parameters remain the same level as initially defined. The numerical experiments show that several model parameters such as network structure, available capacity, budget and arrival rate play an important role (in comparison to other parameters) on the maximum access time. Thus, in this section we present the results of the sensitivity analysis for these parameters.

**Effect of Network Structure on Maximum Access Time:** As mentioned before, the network structure, i.e. the number of regions and providers, can

affect the maximum access time in the network in several ways. For example, the number of regions (regional boundaries) in an outsourcing network is usually identified based on the geography as well as some ad-hoc decisions. Large regions may be preferred to decrease the bureaucratic burden but may be difficult to manage. In order to investigate the effect of the number of regions in a network on the maximum access time, we generate an artificial network by modifying the network $N(12, 40)$ as a network $N(5, 40)$, while total number of providers and total arrival rate in the network are kept the same. Thus, in the new instance, each region has more providers and a larger arrival rate compared to the original network $N(12, 40)$. In other words, we just enlarge the boundries of the existing regions. With this change, the maximum access time in the network decreases by 53% than the one obtained in the network $N(12, 40)$. This result suggests that larger regions are more advantageous than the smaller ones. However, note that the management of operations and contracting may get more complex with larger regions.

We also measure the change in the maximum waiting time when the budget is identified for a different network size that consists of less number of providers as well as regions. Unlike the previous test, here, we change the number of providers in addition to the number of regions in the network.

We also identify three budget levels for each network structure to investigate the effect of the budget at the same time. For the network $N(12, 40)$, we consider different budget levels, 900$, 1200$, and 1800$, that are labelled as 'low budget', 'medium budget' and 'high budget', respectively. For other instances, the budget is halved proportional to their sizes, e.g. the high budget corresponds to 900$ for the network $N(6, 20)$. Figure 3.2 displays the maximum access times in different network structures and budget levels.

101

Figure 3.2: Maximum access times obtained for different network structures at various budget levels

As Figure 3.2 shows, the maximum access time decreases as the size of a network increases. This observation is more prevalent with a tighter budget. It suggests that defining a common budget for a larger network is more advantageous. This can be explained by increased level of risk sharing in a larger network. Also, we see that the decline in the maximum access time gets smaller as the network size is increased.

**Effect of Available Capacities:** The available capacity of each provider depends on the annual strategic plans which can be revised in the coming years. To test the effect of different capacity levels, we solve the instance for the network N(12, 40). We consider the double and half of the original available capacities and three different budget levels: high (1800$), medium (1200$), and low (900$).

The results shown in Figure 3.3 indicate that the maximum access time does not have a linear relationship with the available capacities; when they

Figure 3.3: Impact of different capacity and budget levels on the maximum access time of the network N(12, 40)

are halved, the maximum access time does not double. Also, the decrease in the available capacities has a larger effect on the maximum access time than an increase of the capacities. This difference is more prevalent with a tighter budget. However, we see that when the capacities are halved, even doubling the original budget does not affect the maximum waiting time.

**Effect of Budget and Arrival Variation:** The budget level may change according to the economic conditions. Additionally, the risk-aversion of the decision makers, which defines the conservativeness of the model, can vary among different decision-makers. In other words, the ratio of $\Gamma_r^a$ to the standard deviation of the interarrival times, denoted by $\sigma_r^a$, may be different for all regions $r = 1, \cdots, R$. Thus, we investigate the effect of the budget and the interarrival time variation $(\Gamma_r^a/\sigma_r^a)$ on the maximum access time. Note that $\Gamma_r^a/\sigma_r^a = 2$ for $r = 1, \cdots, R$ in the previous experiments. Since the variation in the service times is quite low, its effect is not presented here.

103

Figure 3.4: Regional outsourced, remaining capacities and maximum access times in the network N(12, 40) when the budget is 300$

For this purpose, we solve the model with the network N(12, 40) for a wide range of budget levels. We found that the problem is infeasible when the budget is 200$. Figure 3.4 presents the maximum access times and the outsourced/remaining capacity levels out of the available capacities in each region when the budget is 300$. We see that the maximum waiting times in the regions are almost equal to each other. This is possibly because the budget is very tight which makes the capacity constraints redundant. In this case, the heuristic optimizes the budget distribution in such a way that the maximum waiting times in the regions are equal to each other.

Figure 3.5 shows the maximum access times obtained at different levels of budget and interarrival time variation for the network N(12, 40). The maximum access time does not have a linear relationship with the budget and interarrival variation. Thus, increasing the budget more than a certain level is not advantageous. Also, the maximum access time is very sensitive to the changes in $\Gamma_r^a/\sigma_r^a$. As expected, the outsourced capacity levels increase when

Figure 3.5: Maximum access time in the network N(12, 40) with varying budget and $\Gamma_r^a/\sigma_r^a$ levels

the budget and arrival variation are increased and decreased, respectively.

**Analysis of an Audiology Outsourcing Network**

To investigate how the heuristic performs with real data, we solve a real-life instance obtained from the NHS audiology services network in the UK (Monitor, 2016). In this case, there are four regions, namely Essex, North durham, Hartlepool and Newcastle, as summarized in Table 3.6. We should note that the budget and the capacities offered by the providers are not available in the online data sources. Thus, we solve the model for a wide range of these parameters. However, the numerical results indicate that the optimum capacity levels change slightly for different available capacity levels. Thus, we present the maximum access time for only the available capacities given in Table 3.6. In order to show impact of future demand realisations (as possible changes of population dynamics) on the capacity planning strategies of the audiology services, we consider four demand

patterns where the original patient arrival rates (in Table 3.6) are increased by 25% and 50%, and decreased by 25%.

Table 3.6: Real data obtained from the NHS audiology network

| Regions | Number of Providers | Price (£) | Patient arrival per day ($\lambda$) | Maximum Capacity ($C$) | Ratio $C/\lambda$ |
|---------|---------------------|-----------|-------------------------------------|------------------------|-------------------|
| Essex | 2 | 294 | 1.35 | 10 | 7.4 |
| N. Durham | 13 | 283 | 3.38 | 36 | 10.65 |
| Hartlepool | 6 | 283 | 2.9 | 21 | 7.24 |
| Newcastle | 9 | 283 | 2.74 | 18 | 6.56 |

The experiments show that when the budget constraint is not too tight, in one of the regions, all available capacity is contracted and some available capacity is left over in the others. This *bottleneck* region identifies the maximum access time within the whole network. Thus, in such a case, the maximum access time within a network can only be decreased by improving the performance in the bottleneck region. For this, the bottleneck region should be first identified from the available data. One intuitive way of identifying the bottleneck region is dividing total available capacity ($C$) with the total arrival rate ($\lambda$) in all regions. The last column of Table 3.6 shows this ratio for the regions in the NHS Audiology network. The results show that although Newcastle has the smallest $C/\lambda$, the bottleneck region is Essex. This counter-intuitive result is possibly due to the number of providers in regions that affect the dispersion of the arrival variation among the providers.

Figure 3.6 shows the maximum access times in different budget and arrival levels. When the budget is lower than £12150 and £11150 and the arrivals are increased by 25% and 50% than the original, respectively, the problem becomes infeasible (shown as blank sections in the graph). According to the NHS statistics, 95% of the patients wait at least 8 weeks for the audiology service in these regions (NHS UK, 2016). Figure 3.6 shows that the current maximum waiting time may be reduced significantly if the budget is increased to £12150. However, when the budget is larger than £12150, the maximum access time within the network does

106

Figure 3.6: Maximum access time for the Audiology outsourcing network with varying budget and patient arrival rates

not change significantly. As the budget is increased to £16150, the maximum access time is not affected significantly by an increase in the patient arrivals. Also we see that the maximum access time does not vary much when the arrivals are decreased by 25% from the original and the budget is larger than £12150.

## 3.6   Conclusions

Outsourcing has been increasing within the healthcare sector mainly due to the increasing demand for more cost-effective services. This increase results in healthcare outsourcing networks in which several providers share a patient population based on contractual relationships. This chapter focuses on the strategic capacity planning problem in a healthcare outsourcing network. We consider a health authority which buys healthcare services from available providers based on a fixed price. The authority has a limited budget to be used for the contracts.

We develop a mathematical model to find the optimum capacities to out-

107

source from each provider. A robust queuing approach is used to approximate the maximum access time in the providers. The resulting model is intractable due to the non-linear and integer formulations. To solve the problem, we propose an alternating optimization based heuristic combined with bisection search. The computational experiments show that the heuristic performs better than the available commercial solvers, especially when the budget is tight. The sensitivity analysis reveals that it is more advantageous to identify larger regions with more number of providers rather than the smaller ones. Also, defining the budget for a smaller network with less providers and regions results in a higher maximum access time. The computational experiments with the real data show that the current waiting times for the NHS Audiology services can be reduced.

The developed model and solution method can be applied to several healthcare services such as surgery or imaging services. Besides, although the model is developed for healthcare outsourcing, it can be applied to any kind of service outsourcing network, with slight modifications if needed. As future work, the shortcomings of some assumptions used in modelling of the problem can be studied. For instance, in the current model, the contract prices between the providers and the central authority are fixed. On the other hand, the negotiations between the providers and the central authority can be studied by using a game-theoretic approach. Another possible extension is to consider the effect of patient choice during the referral on the overall performance. This would require to add choice models onto the original capacity planning study. In the future work, contract types in terms of payment structures such as activity-based or hybrid payment systems can be included into the decision making process.

# Chapter 4

# Real-time Surgery Planning under Uncertainty in Surgery Suite

## 4.1  Introduction

A surgery suite can be seen as the engine of a hospital. Surgeries generate around 40 per cent of the revenues in the UK hospitals (HFMA, 2005). They are among the most profitable healthcare services, with prices of up to a hundred thousand pounds (Carey et al., 2011). Surgeries also account for the majority of hospitals' operational capacity (Macario et al., 1995). They consume a significant amount of physical resources, such as beds and equipment, as well as human resources with different levels of expertise. Given its impact on the revenue and resource usage, surgery management is one of the most crucial tasks for healthcare professionals. However, a recent study indicates that the current management practices are not able to reach the performance targets such as average time to get service (Department of Health, 2016).

Surgery management involves four major decision-making stages: strategic case-mix planning, development of a master surgery schedule, scheduling of indi-

vidual elective cases and reactive surgery scheduling. The strategic case-mix plan-

ning identifies the surgical blocks i.e. how the available time in operating rooms'

(ORs) is distributed among different specialities (for instance, see Yahia et al.

(2015)). The second stage, also called master surgery scheduling, consists of the

assignment of surgical specialities to surgical blocks over the scheduling horizon

(typically one week) in order to maximize the resource utilization (for instance,

see Beliën et al. (2007)). The third stage, which has been most widely studied,

involves the assignment of specific surgeries directly to the surgery blocks identi-

fied in the previous step. Reactive surgery scheduling, the last stage of surgery

management, is defined as the real-time management and revision of surgery suite

schedules as disruptions such as non-elective admissions occur (for instance, see

Stuart and Kozan (2012)). Readers are referred to Erdogan et al. (2011) for more

details of the decision-making stages in surgery management.

The last two stages of surgery management involve offline and online (real-

time) tasks, depending on whether the patient is elective or non-elective. A non-

elective and elective patients are simply referred as 'non-elective' and 'elective',

respectively, throughout the chapter. Elective surgeries, which constitute most

of the demand, are scheduled days or weeks ahead. There are different rules

to sequence elective surgeries such as *the longest surgery first* and *the shortest

surgery first* in which the surgeries are ordered according to the *descending* or

*ascending expected duration*, respectively. More sophisticated sequencing rules

can be listed as *longest waiting time, earliest start time* and *latest start time*.

This chapter is not concerned with the elective surgery scheduling; readers are

referred to Guerriero and Guido (2011) for a detailed review of elective surgery

scheduling.

In reactive surgery scheduling, various types of disruptions such as patient

no-shows or the staff unavailability need to be considered. The most crucial one of

these disruptions can be observed on the non-elective arrivals and the variations in surgery durations (Van Riet & Demeulemeester, 2015). Besides, the realization time of the disruptions and possible effects are not known in advance. Thus, the reactive surgery scheduling requires a real-time decision-making process in a very short time by considering many criteria such as costs, and patient and staff satisfaction. Most importantly, the decisions made at any time may affect the future schedules significantly.

In order to handle these kinds of disruptions, hospitals in general develop two main strategies (Van Riet & Demeulemeester, 2015). In the first approach, non-electives are treated separately from the scheduled cases by reserving dedicated room(s) based on the predicted demand. However, this results in an inefficient schedule due to the uncertainty in non-elective arrivals (Van Riet & Demeulemeester, 2015). For example, when there is no non-elective arrival, this strategy may result in revenue losses since the dedicated rooms and medical staff have to stay idle. The other strategy accommodates the non-electives by allowing staff overtime and cancellation of electives in a reactive fashion (e.g. see Ozkarahan, 2000; Blake et al., 2002). This (online) approach brings patient discontent and extra cost due to cancellations and overtime (Hosseini, 2012). However, Wullink et al. (2007) showed that the online approach usually results in a better performance in terms of waiting time, staff overtime, and OR utilization. This chapter focuses on decision-making problems encountered in the online approach. Recently, two approaches are combined in some hospitals (Van Riet & Demeulemeester, 2015). In this so-called hybrid approach, some buffers are left within the elective schedules to accommodate the non-electives (Van Riet & Demeulemeester, 2015).

The online approach requires to take critical interdependent decisions through-

111

out a day. At the beginning of a typical work day, the schedule of elective surgeries in a surgery suite for that day is usually ready. Each elective/non-elective surgery is associated with an estimated duration (number of required time slots). As the surgeries are carried out, non-electives in different health conditions arrive randomly to the hospital. In particular, when a non-elective arrives, the responsible staff of the hospital must decide whether or not to accept the patient and assign him/her to one of the operating rooms. Since non-electives usually require urgent treatment, accepting a non-elective may lead to the postponement of pre-scheduled surgeries. On the other hand, the delays in the scheduled start times may result in the deterioration of their health conditions. In addition to this, hospitals are generally concerned with the extra staff overtime and consequently the additional cost. Besides, the last surgery in a day need to be completed by the end of shift time at that day (i.e. 24 hours). If the completion time of the last surgery exceeds the planning horizon, the decision maker may also cancel some electives. Cancellations do not only create anger and discontent for elective patients (Schofield et al., 2005), but also impact on the future schedules. Due to its effect on the cancellations and overtime, the acceptance of a non-elective is a dynamic decision and should consider all future possibilities. According to the NHS statistics, 20,464 elective surgeries were cancelled at the last minute by English hospitals at the last quarter of 2014/15 (NHS, 2014). Moreover, 3,567 emergency patients were rejected from the same hospitals in 2014. The main reason for these cancellations and rejections was capacity shortages (Campbell and Arnett, 2015).

Due to the uncertain and dynamic nature of the surgeries, as well as various criteria such as overtime, cancellations and rejections, the real-time management of a surgery suite is a challenging and important decision-making process for hospitals. Moreover, hospital managers are under the pressure of reducing patient

dissatisfaction and operational costs. This complex process requires an elaborate mathematical approach. However, the literature related to the real-time surgery management is scarce (Van Riet & Demeulemeester, 2015). The related studies lack the comprehensive analysis of all relevant uncertainties and criteria which are crucial to obtain rigorous solutions. Also, the evolving and dynamic nature of the real-time surgery management is overlooked in most of the papers. Thus, this research aims to fill this gap and find an optimum policy to manage non-elective arrivals, such that patient satisfaction is maximised while overall operational costs are minimised. For this purpose, we develop a stochastic dynamic programming model of the real-time surgery management problem. Due to the problem size, it is computationally intractable to find exact solutions for practical instances. For such cases, approximate dynamic programming (ADP) is used. The policies obtained through ADP are then compared with a myopic approach for different cost schemes. Finally, we investigate the effect of different elective scheduling strategies on the overall cost by using generated data inspired from real data.

The remainder of this chapter is organised as follows. Section 4.2 presents a literature review of the real-time surgery management problem and ADP applications in healthcare. Section 4.3 first describes the underlying problem in more detail and then introduces a stochastic dynamic programming model for the daily management of a surgery suite. We extend the model by considering the uncertainty in surgery durations and multiple surgery rooms. Section 4.4 explains the solution approach, ADP, in more detail. Section 4.5 presents the design as well as the findings of the computational experiments.

## 4.2 Literature Review

The Operations Research community has conducted extensive studies for various strategic and tactical decisions as well as offline and online approaches arising within the surgery management (Cardoen et al., 2010). However, the studies on the real-time surgery management is limited. In this section, first, we review the literature on the real-time surgery management. This is then followed by an examination of the ADP applications in healthcare.

### 4.2.1 Modelling of Real-time Surgery Planning Problem

The closest task to surgery planning is that of appointments. However, appointment management is not as complex as the real-time surgery management since there are no significant disruptions to appointments such as random non-elective arrivals or task durations (Gupta and Denton, 2008). Thus, the research related to the appointment management is not included in this review.

A real-time production scheduling problem in manufacturing has been widely studied by many researchers (Billaut and Roubellat, 1996; Wu et al., 1999; Aloulou and Portmann, 2003). However, the nature of the healthcare service i.e. involving the lifes of patients as well as significant operational uncertainties such as surgery durations, distinguishes the healthcare applications of real-time management from those in manufacturing area.

Real-time surgery management has not received enough attention in the literature, as supported by Erdogan et al., (2011), Guerriero and Guido (2011) and Van Riet and Demeulemeester (2015). Table 4.1 summarises previous research on the real-time surgery management. The related papers have considered various performance measures such as overtime and patient-related costs including those of cancellation and rejection; however, none has studied all measures together.

Given their effects on decision making, these measures should be considered together for a more rigorous analysis (Van Riet and Demeulemeester, 2015). Similarly, although many authors have considered the uncertainty either in surgery durations or non-elective arrivals, only Hosseini (2012) and Borgman (2017) has considered both of them together. Two-stage stochastic programming has been used to model these uncertainties (Batun, 2011; Zhang et al., 2013; Heydari and Saoudi, 2016). In this method, the first-stage decisions are taken before the realisation of uncertainty. After the uncertainty has unfolded, the second-stage decisions, or so-called recourse decisions, are taken to adjust the earlier decisions. However, the uncertainties affecting the real-time surgery management, namely the non-elective arrivals and surgery durations, are revealed throughout the day, not only once a day. Therefore, stochastic dynamic programming, which models multi-stage stochastic decision-making problems, is a more suitable technique for modelling of the real-time surgery management problem.

In the real-time surgery management, analytical approaches can be applied each time when there is a disruption to the schedule (reactively) or to derive an optimum policy consisting of the optimum action for each possible case (proactively). Reactive decision-making models for real-time surgery management are developed by Stuart and Kozan (2012), Van Essen et al. (2012), Duma and Aringhieri (2015) and Erdem et al., (2012). Stuart and Kozan (2012) focus on the reactive surgery scheduling problem as random non-elective arrivals occur. They model the problem as a single machine scheduling problem with sequence dependent processing times and due dates by including the priorities of elective and non-elective cases. Two conflicting objectives considered in the paper are to maximise the number of non-electives inserted into the schedule and to minimise the number of cancelled electives. They assume that surgery durations are sequence-dependent and follow independent log-normal distributions. The opti-

mization model is solved by an exact method, the branch-and-bound algorithm, which produces a list of surgeries that are expected to be late.

Due to the computational difficulties in reactive surgery scheduling problem (Van Essen et al. 2012), several authors (Van Essen et al. 2012; Duma and Aringhieri, 2015; Erdem et al., 2012) have developed approximate solution methods. Van Essen et al. (2012) model a reactive surgery scheduling problem as an integer linear program to minimise the deviances from stakeholders' preferences. Their model can be used as a decision support system to determine a new elective surgery schedule in case of a disruption. They show that the reactive surgery scheduling problem for multiple rooms is NP-hard. Unlike us, they assume that non-elective patients are always accepted for the operation. Similarly, Erdem et al. (2012) construct a deterministic mixed-integer linear programming model that reschedules elective patients to the operation and post-anaesthesia clinical care units when a non-elective patient arrives. They assume that non-elective patients may be rejected or accepted. The objective function consists of the costs of postponing electives and rejecting non-electives, whereas it disregards the cost of overtime. They use a genetic algorithm to obtain approximate solutions for real-sized instances. Simulation is another approximate technique used for reactive surgery scheduling (Duma and Aringhieri, 2015).

Table 4.1: A classification of the research papers on real-time surgery management

| Research papers | Modelling | | | Decisions | | | Uncertainty | | Costs | | Solution Method | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IP | TS | MDP | Reject | Cancel | Assign | Arrival | Duration | OR & Staff | Patient-related | Exact | Heuristic |
| Batun (2011) | | ✓ | | | | ✓ | ✓ | | ✓ | ✓ | ✓ | |
| Hosseini (2012) | | | ✓ | | | ✓ | ✓ | ✓ | ✓ | | | ✓ |
| Stuart & Kozan (2012) | ✓ | | | ✓ | ✓ | | | ✓ | | ✓ | ✓ | |
| Erdem et al.(2012) | ✓ | | | ✓ | | | | | ✓ | ✓ | | ✓ |
| VanEssen et al. (2012) | ✓ | | | | | ✓ | | | ✓ | ✓ | | ✓ |
| Zhang et al.(2013) | | ✓ | | | ✓ | | | ✓ | ✓ | | ✓ | ✓ |
| Duma & Aringhieri (2015) | Simulation | | | | ✓ | | | ✓ | ✓ | | ✓ | ✓ |
| Addis et al. (2016) | ✓ | | | | | ✓ | | ✓ | | ✓ | ✓ | |
| Heydari & Saoudi (2016) | | ✓ | | | | ✓ | ✓ | | ✓ | | ✓ | |
| Our approach | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

*IP: Integer programming, *TS: Two-stage stochastic programming, *MDP: Markov Decision Process

The reactive decision-making approach ignores the dynamic nature of the real-time surgery management problem. Therefore, several authors (Zhang et al., 2014; Heydari and Saoudi, 2016) have considered the proactive modelling approach in which all uncertainties and multiple stages of decision-making are included to obtain a proactive policy. For example, Zhang et al. (2014) are concerned with the dynamic assignment of surgeries to multi-rooms. They develop a two-stage stochastic programming model which identifies the next surgery to assign to any room that is currently available. However, they do not consider the non-elective arrivals and the disruptions in the schedule related to these arrivals. The random non-elective arrivals are considered in Heydari and Saoudi (2016) who develop a two-stage stochastic model that reschedules the surgeries upon the arrival of an emergency. They do not model the operational decisions such as cancellation of electives or rejection of non-electives, i.e. they assume that non-elective surgeries are always accepted. Borgman (2017) also considers a simulation approach to evaluate different non-elective surgery management strategies such as dedicating emergency rooms or combining electives and non-electives. On the other hand, the acceptance or rejection of non-elective patients is not considered in this analysis.

Considering the dynamic nature of the real-time surgery management problem, the current literature lags behind in terms of the methodologies considered. Although stochastic dynamic programming has been widely applied to model different stages of surgery management such as case-mix planning and elective surgery scheduling (Min and Yih, 2014; Gerchak et al., 1996; Lamiri et al., 2008), only Hosseini (2012) has used it for the real-time surgery management. The author models two FCFS queues for the elective and non-elective surgeries separately. The surgeries are assumed to have a random service duration. A Markov decision process (MDP) model is developed to minimise the system-wide, long-

run average costs relating to patient allocations. The state space consists of the length of queues and operating room idleness while the action space consists of assigning a room to a patient from either queue.

Several tactical and strategic decisions regarding the surgery management i.e. number of surgery rooms, number of hours in a shift, etc. can be incorporated with the modelling of operational decisions in real-time surgery management, such as cancellations, rejections and scheduling. For example, Batun (2011) considers the surgery rescheduling problem in a surgery suite. The author develops a two-stage stochastic formulation in which the first-stage decision is to determine the number of operating rooms to open in a day while the second-stage decision is to reschedule surgeries in the middle of the day. L-shaped decomposition and progressive hedging algorithms are used to solve the two-stage stochastic mixed-integer programming model. The algorithms solve each scenario sub-problem independently and enforce the non-anticipativity constraints progressively. However, the author assumes that rescheduling happens only once a day, which is not applicable to most cases.

Instead of modelling real-time surgery management decisions for a short planning period like a day, a longer planning horizon may also be considered as in Addis et al. (2016). Specifically, they study a rolling-horizon elective surgery scheduling problem. In the first step, the optimal schedule for several weeks is found by solving an integer linear programming model. After the first week is implemented, they re-optimise the schedule by adding new arrivals and cancellations from the previous week. They also employ a robust optimization approach and develop uncertainty sets around the surgery arrivals.

We develop a stochastic dynamic programming model for the real-time surgery management problem for multiple surgery rooms during a day. Unlike the related models in the literature, our model takes into account dynamic actions,

uncertainties on non-elective arrivals and surgery durations as well as different criteria such as non-elective rejections, overtime, elective cancellations and waiting time. However, the resulting model has a large state space and thus, the real-sized instances are not solvable in reasonable time with the exact solvers. For this reason, we develop an approximate solution approach based on ADP, a solution framework suitable for large size stochastic dynamic programming formulations. We design a series of computational experiments to illustrate the performance of the ADP algorithm. The numerical results show that the proposed algorithm provides a good approximation to optimum policy. Moreover, we compare the performances of the ADP algorithm and a myopic heuristic for different levels of cost components. Finally, we investigate the impact of various surgery scheduling rules on the overall operational cost.

## 4.2.2 Approximate Dynamic Programming for Healthcare Applications

ADP is a very popular approximation technique to solve large dynamic programming problems encountered in many real-life problems. For example, it has been applied to resource planning (Erdelyi and Topaloglu, 2010; Schütz and Kolisch, 2011), inventory control (Simao and Powell, 2009; Roy et al., 1997), inventory routing (Adelman, 2004), option pricing (Tsitsiklis and van Roy, 2001), game playing (Yan et al., 2004), revenue management (Adelman, 2007) and transportation (Topaloglu and Powell, 2006). Healthcare related strategic and operational problems, such as resource allocation and patient scheduling in particular, have also been active application areas of ADP.

Linear programming based ADP is used by several authors to solve the patient scheduling problem (Patrick et al., 2008; Gocgun and Puterman, 2014; Barz and Rajarm, 2015). In these papers, further reduction techniques such as

column generation are employed to solve the resulting linear programming formulations (Patrick et al. 2008; Saure et al., 2012; Gocgun and Puterman 2014). For example, Barz and Rajarm (2015) solve the elective patient admission/scheduling and resource usage problem with a linear programming based ADP. They assume that emergency patients use some of the capacity reserved for elective patients, and that patients' health conditions follow a random trajectory. They consider three types of cost: overtime, loss of patient goodwill, and inferior care. However, the resulting approximate linear problem is still computationally expensive and requires relaxation to be solved.

Due to the computational difficulties encountered in linear programming based ADP, simulation based ADP methods have been preferred by several authors to solve the patient scheduling problem (e.g. Hulshof et al. (2013), Lin et al., (2011)). Among the simulation based methods, basis function approximation has been used by Hulshof et al. (2013), while state aggregation is preferred by Lin et al. (2011). Specifically, Hulshof et al. (2013) model elective patient admissions and intermediate-term resource allocation in hospitals under uncertain treatment paths and patient arrivals. The state space consists of the number of patients waiting in the service queues and their waiting times. At each epoch, the decision maker must decide on the number of patients to treat from each queue. The total cost depends on the patients' waiting time. The computational experiments show that a basis function based ADP algorithm with value iteration performs better than two other heuristics. Lin et al. (2011) develop an MDP model for a sequential scheduling problem to optimise the performance of a clinic in the presence of overbooking and no showing up of patients. They divide the call-in period into discrete time intervals, such that the scheduler can handle no more than one call within an interval. The action set is composed of scheduling a patient to a future day or transferring the request to another day. The model is solved by a

simulation based ADP in large instances.

ADP is especially useful for large size healthcare problems such as ambulance redeployment and scheduling problem (Maxwell et al., 2010; Schmid, 2012) which may even have infinite state spaces. Schmid (2012) solves the ambulance relocation and dispatching problem by introducing a simulation based ADP method. The decision-maker needs to identify which ambulance to send to an incoming call and where to locate this ambulance after it has left the patient. The time unit between two epochs is assumed to be constant, although the time of the triggering event (arrival of a request) is random. In other words, requests arriving during a period are kept in a waiting list and served in the next decision epoch. They aggregate the elements of state to overcome the computational difficulties. Maxwell et al. (2010) claim that the state aggregation may result in multiple optimum actions, instead, they develop a basis function based ADP to solve the ambulance redeployment problem. Specifically, they develop an event-driven model to establish where to redeploy idle ambulances to maximise the number of calls served within a certain threshold.

Other than patient and ambulance scheduling, ADP has been applied to solve medical decision making problems (Mason, 2012). Mason (2012) develops two MDP models to prevent adverse events and to determine the optimum timing of adherence-improving interventions after a treatment is started. The author develops lookup table and basis function based ADP algorithms and compares their computational performances.

Our review shows that, considering the computational difficulties encountered in the real-sized healthcare problems, the number of ADP applications in this area is very limited. Also, it has not been used to solve a real-time surgery management problem. Thus, one of our contributions in this chapter is to implement a tailored ADP algorithm to solve the real-time surgery management

122

problem. The next section describes the underlying problem in more detail and provides a stochastic dynamic programming formulation.

## 4.3 Stochastic Dynamic Programming Model for Real-time Operating Room Planning

In this section, we present a stochastic dynamic programming model for the real-time management of the surgery schedule of single operating room in the presence of random non-elective arrivals. We consider one day planning horizon. The surgery room is well equipped to serve different types of surgeries such as orthopaedic or heart operations. For each surgery type, the expected duration has been estimated in terms of the number of required time slots where one time slot is equal to half an hour. In the start of the day, the initial schedule of the surgery room is available to the decision-maker e.g. a surgery suite manager. We assume that, during the day, non-elective patients arrive randomly.

In particular, when a non-elective patient arrives, the surgery suite manager needs to make a rapid decision on either rejecting or accepting the patient for the surgery at that day. We are especially concerned with the patients who need an urgent treatment. Moreover, the rejection of a non-elective causes loss of potential profit that could have been obtained from this surgery if s/he was accepted. If the patient is accepted, the scheduled surgeries may face delay, since non-elective surgeries are usually more urgent than electives. Delays may cause in the deterioration of patients' health as well as patient discontent and a decrease in resource utilization. Staff overtime, that counts for the work-hours after a fixed time threshold (overtime threshold), may be used to accommodate the extended shifts, which also entails additional costs. However, in any case, the last surgery scheduled to the room must be completed before a certain time limit. When the

expected completion time exceeds the time limit, some of the scheduled surgeries need to be cancelled with a certain cost. In summary, the delays due to the admission of non-electives are reflected in the model as a source of cost. On the other hand, rejection of non-electives are represented by an opportunity cost. The model aims to minimize total costs of cancellation, rejection, overtime and delays in view of the resource limitations.

Since elective patients are advised to arrive to the hospital long before the scheduled surgery time, the delays in the starting time of elective surgeries are usually negligible. Similarly, even though an elective patient does not show up on the surgery day, the next surgery starts on time as initially scheduled. This is because the surgery crew for the next surgery may be busy until the scheduled starting time of this surgery.

## 4.3.1 Real-time Operating Room Planning Model under Deterministic Surgery Duration

In this section, we introduce the formulation of a real-time operating room planning problem. The following assumptions are made for the development of the mathematical model.

- We assume that non-elective patients have priority to be scheduled over elective ones. In other words, if a non-elective patient is accepted, then s/he will be operated as soon as the current surgery in the operating room is finished. Consequently, the starting times of the scheduled semi-urgent surgeries are delayed (Van Riet and Demeulemeester, 2015).

- The duration of a surgery covers the time taken from the initial preparation of the operating room for the current patient until the preparation of the surgery room for the next patient. We initially assume that the durations

of all types of surgeries are certain and known in advance. However, in the next section, this assumption is relaxed to take uncertain surgery durations into account for modelling the real-time surgery management problem.

- The overtime cost (refers to the cost of staff working overtime) linearly depends on the overtime incurred (Talluri et al., 2006). Note that the other types of costs such as cancellation and rejection are based on per surgery rather than the duration of surgeries.

- Non-elective patients arriving after the overtime threshold are not considered since the acceptance/rejection of these non-electives usually depends on whether the room is empty or not; if the room is empty, then the non-elective is accepted.

The surgery types are classified into $M$ groups, represented by $m = 1, \cdots, M$. Each surgery type is associated with an expected surgery duration denoted by $d_m$. The planning horizon is considered as one day and discretized into $T$ time points (so called as decision epochs) represented with $t = 1, \cdots, T$. The time between decision epochs $t-1$ and $t$ is called as time period $t$ for $t = 1, \cdots, T$. Note that a *period* is the time interval when a non-elective patient may arrive. On the other hand, a *decision epoch* represents the time point where a decision is made. We assume that an equal duration length $\Delta$ is used for all time periods and selected by the modeller in a way that there will be at most one non-elective arrival during $\Delta$, e.g. 15 minutes. The final epoch $T$ represents the overtime threshold where the overtime cost starts to incur. No action is taken on and after epoch $T$. In other words, we assume that the surgery schedule does not change after the overtime threshold. Although a non-elective patient arrives any time during period $t$, a decision on his/her acceptance or rejection is made at decision epoch $t$. The overall time limit that the last surgery should be completed is denoted by $T^{max}$.

The cost of cancelling a surgery is assumed to be the same for all surgeries and shown with $\theta^c$. On the other hand, $\theta^m$ represents the cost of rejecting a non-elective patient with type $m$ for $m = 1, \cdots, M$. This is considered as an opportunity loss that depends on surgery type $m$.

Suppose that $e$ number of surgeries are initially scheduled to the operating room for a specific day. Since there can be at most one non-elective arrival during a time period, the operating room can serve at most $(e + T - 1)$ surgeries during the day. Let's represent $(e + T - 1)$ with $I$ and a surgery with $i = 1, \cdots, I$. Note that indices $i = 1, \cdots, e$ correspond to the pre-scheduled electives. Thus, reasonably, a non-elective accepted for an operation at time $t$ can be assigned to index $i = e + t$. To formulate this assignment, we introduce a binary parameter $z_i^t$ taking 1 at time $t$ if $i = e + t$ and 0, otherwise,

$$
z_i^t =
\begin{cases}
1, & \text{if } i = e + t, \\
0, & \text{otherwise.}
\end{cases}
$$

The use of parameter $z_i^t$ will become more clear shortly. Next, we will describe a mathematical programming formulation of the real-time surgery management problem in terms of states and actions.

**States:** At the beginning of epoch $t$ (before any decision is taken), the following information regarding the state of the system is available to the decision-maker:

- Let vector $\mathbf{C}^t = (C_1^t, \cdots, C_I^t) \in \{0, \cdots, T^{max}\}$ define the completion times of surgeries $i = 1, \cdots, I$, where $C_i^0 = 0$ for $i = e + 1, \cdots, I$. If surgery $i$ is cancelled at epoch $t'$, then $C_i^t = 0$ for $t = t', \cdots, T$. We assume that the surgery with completion time $t$ finishes during epoch $t$ before any decision is made in this epoch, but after the state information is accrued. For example,

126

if $C_1^2 = 2$, then surgery $i = 1$ will be finished before a decision is made for the non-elective arriving during period 2 (if there is any). Note that the completion time list is not necessarily in an increasing order.

- The type of the non-elective arrived during time period $t$ is denoted by $m^t \in \{1, \cdots, M\}$ for $t = 1, \cdots, T$. The expected duration of this non-elective is represented with $d_{m^t} \in \{0, \cdots, d^{max}\}$, where $d^{max}$ is the maximum possible duration that a non-elective can have. It is also possible that there is not any non-elective arrival in this period which corresponds to $d_{m^t} = 0$.

In view of all information, we denote a state space $(\mathbf{C}^t, m^t)$ consisting of completion times of surgeries and the type of non-elective arrival.

**Actions:** The binary decision variable $x^t$ at epoch $t$ shows whether we accept a non-elective (if there is any) arrival during period $t$ ($x^t = 1$) or not ($x^t = 0$).

The order of the activities during decision epoch $t$ can be summarized as follows. First, the decision maker learns the type of (any) non-elective arrival. Second, the surgery that has completion time $t$ (if any) finishes or not. Third, the decision maker accepts (any) non-elective arrival in period $t$ or not on the basis of the current state information. Finally, if there is any surgery finishing in the current epoch, then the next surgery (if any) starts.

**Update in state space:** At the end of each decision epoch $t$, the completion times of all surgeries need to be updated according to action $x^t$ and the type of the non-elective arrival $m^t$. In other words, the completion times in $t+1$ depend on information about $x^t, \mathbf{C}^t$, and $m^t$ and, thus, denoted by $\mathbf{C}^{t+1}(\mathbf{C}^t, x^t, m^t)$. If the non-elective arrived during the last period $t$ is rejected ($x^t = 0$), then the completion times of the surgeries accepted so far do not change. When a non-elective arrival is accepted ($x^t = 1$), then the following updates take place:

(a) The surgery of the accepted non-elective has given a priority in the waiting list due to emergency of the patient's medical condition. Thus, the completion times of the waiting surgeries are increased by the expected duration of the non-elective, $d_{m^t}$. For this, we first need to identify the uncompleted surgeries at period $t$ by checking whether its completion time is larger than the current period or not. Let's define binary variable $a_i^t$ taking 1, if the completion time of surgery $i$ is larger than or equal to $t$ and 0, otherwise,

$$a_i^t = \begin{cases} 1, & \text{if } C_i^t \geq t, \\ 0, & \text{otherwise.} \end{cases} \tag{4.1}$$

(b) A completion time should be assigned to the accepted non-elective patients. To formulate this, we follow a two-step procedure. In the first step, we compute the operation starting time for this surgery. Let's represent the starting time of the operation for the non-elective with $K^t \in \{0, \cdots, T^{max}\}$. If there is any surgery under operation at period $t$, then the starting time of the non-elective, $K^t$, is equal to the completion time of the ongoing surgery. If no surgery is in progress, it is simply equal to the current epoch $t$.

In the second step, we determine the completion time of the surgery under operation by identifying its indice. For this purpose, we define $b_i^t$ that is equal to 1 if surgery $i$ is currently served in the room at time $t$; 0, otherwise. Note that the surgery under operation at epoch $t$ always has the smallest completion time among the uncompleted surgeries $(i|a_i^t = 1)$:

$$b_i^t = \begin{cases} 1, & \text{if } C_i^t = \min_{k \in \{1, \cdots, I\}} \{C_k^t | a_k^t = 1\}, \\ 0, & \text{otherwise.} \end{cases} \tag{4.2}$$

Then, we can formulate the starting time of the accepted non-elective surgery,

128

$K^t$, as follows:

$$K^t = \sum_{i=1}^{I} C_i^t b_i^t + t\left(1 - \sum_i b_i^t\right), \ t = 1, \cdots, T-1. \quad (4.3)$$

The first term in (4.3) computes the completion time of the surgery under operation at epoch $t$. The second term ensures that if the surgery room is idle at time $t$, then the operation of the non-elective patient admitted at period $t$ starts immediately.

(c) When the overall completion time exceeds the daily time limit, $T^{max}$, the surgeries from the end of the list should be cancelled. In other words, their completion times should be assigned to 0. In order to achieve this, we first need to formulate (a) and (b) as described above and compute the (temporary) completion times represented with $C_i^{t+1,y} \in \{0, \cdots, T^{max} + d^{max}\}$ for $i = 1, \cdots, I$ and $t = 1, \cdots, T-1$. We can then check whether the latest completion time is larger than the time limit or not. In other words, $C_i^{t+1,y}$ displays the updated completion times before the cancellations and can be formulated as follows:

$$C_i^{t+1,y} = \left(1 - (a_i^t - b_i^t)\right)C_i^t + (a_i^t - b_i^t)\left[C_i^t + x^t d_{m^t}\right] + x^t z_i^t \left[d_{m^t} + K^t\right],$$
$$i = 1, \cdots, I, \ t = 1, \cdots, T-1. \quad (4.4)$$

Note that the first term ensures that the completion times of the **completed surgeries** do not change while the second one delays the completion times of the **waiting surgeries** and the last term assigns a completion time to the **non-elective** (if accepted).

Finally, we identify the surgeries with a completion time larger than the time limit. Therefore, we define a new binary variable $j_i^t$ taking 1 if the completion

time of $i$ is larger than $T^{max}$; 0, otherwise:

$$
j_i^t = \begin{cases} 1, & \text{if } C_i^{t+1,y} > T^{max}, \ i = 1, \cdots, I, \ t = 1, \cdots, T-1, \\ 0, & \text{otherwise.} \end{cases} \tag{4.5}
$$

Then, we obtain

$$
C_i^{t+1} = C_i^{t+1,y}(1 - j_i^t), \ i = 1, \cdots, I, \ t = 1, \cdots, T-1. \tag{4.6}
$$

To facilitate the explanation of the model, we present a timeline of the decision-making process along with the corresponding model notation in Figure 4.1.



Figure 4.1: A description of the decision-making process along with notation

**Costs:** The immediate cost of action $x^t$ consists of the costs of rejecting, cancelling and waiting of surgeries at time $t$ and is computed as follows:

$$
\eta(\mathbf{C}^t, x^t, m^t) = \theta^{m^t} \left[ 1 - x^t \right] + \theta^c \sum_{i=1}^{I} j_i^t + \theta^w x^t \sum_{i=1}^{I} (a_i^t - b_i^t - j_i^t) d_{m^t}, \ t = 1, \cdots, T-1, \tag{4.7}
$$

where $\theta^m$, $\theta^c$, and, $\theta^w$ denote, respectively, the unit cost of rejecting a surgery

130

with type $m$, and cancelling and waiting a surgery.

Given a list of completion times in the beginning of the day, $\mathbf{C^0}$, the objective is to minimize the value function at epoch $t = 0$ which can be written as

$$v_0(\mathbf{C}^0) = \sum_{m=1}^{M} \Pr(\tilde{m}^1 = m)v_1(\mathbf{C}^1, m). \tag{4.8}$$

Note that at epoch $t = 0$, there is no non-elective arrival. Thus, no action is taken and consequently the completion times of the surgeries in the list remain the same: that is $\mathbf{C}^0 = \mathbf{C}^1$. In other words, the computation of the value of the initial state is not affected by the decisions, but still presented for the sake of clarity. Given state $(\mathbf{C}^t, m^t)$ at time $t$, the value function, $v_t(\mathbf{C}^t, m^t)$, can be formulated in a recursive form as follows:

$$v_t(\mathbf{C}^t, m^t) = \min_{x^t \in \{0,1\}} \left\{ \eta(x^t, \mathbf{C}^t, m^t) + \sum_{m=1}^{M} \Pr(\tilde{m}^{t+1} = m)v_{t+1}(\mathbf{C}^{t+1}(x^t, \mathbf{C}^t, m^t), m) \right\},$$

$$t = 1, \cdots, T - 1. \tag{4.9}$$

Note that the value function for period $t = 0$ is different than the value function for $t = 1, \cdots, T - 1$. For any state $(\mathbf{C}^t, m^t)$, the value function computes the optimum action that minimizes the overall expected cost. The expected cost of each action (acceptance or rejection of a non-elective surgery) is computed as the sum of the immediate cost of the action and the expected future cost after one period assuming that the corresponding action is taken. Note that, at epoch $t$, the type of the non-elective arrival, $\tilde{m}^{t+1}$, at the next period $t + 1$ is not known. The expected cost after one period for each non-elective type is calculated as the multiplication of the probability of type $m$ non-elective arrival, $\Pr(m)$, with the value function of the state $v_{t+1}(\mathbf{C}^{t+1}(x^t, \mathbf{C}^t, m^t), m)$ given that the corresponding type of non-elective ($m$) arrived in the next period, $t + 1$.

The value function in the final epoch $T$ consists of only the overtime cost

131

and is formulated as

$$v_T(\mathbf{C}^T) = \theta^o \max \left\{ \max_{i=1,\cdots,I} \{C_i^T\} - T, \ 0 \right\}, \tag{4.10}$$

where $\theta^o$ denotes the overtime cost of one time period. If the maximum completion time, $\max_{i=1,\cdots,I} \{C_i^T\}$, is larger than the overtime threshold $T$, then the overtime is defined as the difference between the maximum completion time and $T$. Otherwise, there is no overtime surgery.

The optimum policy, consisting of the optimum action for each possible state, can be found by using the traditional backward value iteration technique (Boyan and Littman, 2000). This technique starts by calculating the values of all possible final states and moves to previous time periods iteratively. At a certain state $(\mathbf{C}^t, m^t)$, for each $x^t \in \{0, 1\}$, we conduct a set of operations to find the optimum state value and action. First, the set of future states that would be attained by the action $x^t$ is found by first computing $\mathbf{C}^{t+1}(\mathbf{C}^t, m^t, x^t)$ based on Equations (4.1), $\cdots$, (4.6), and then combining it with all possible $m^{t+1} \in \{1, \cdots, M\}$. Since $\mathbf{C}^t$ does not change for rejecting the non-elective (i.e. $x^t = 0$), it is enough to find the (next period) completion times for accepting the non-elective (i.e. $x^t = 1$). Then, the value of each future state is multiplied with the probability of $m^{t+1}$ and summed up to find the expected (future) cost. This expected cost is aggregated with the immediate cost $\eta(x^t, \mathbf{C}^t, m^t)$ to compute the overall cost of $x^t$. Finally, the value of the state $(\mathbf{C}^t, m^t)$ becomes as equal to the minimum of the overall costs of two actions.

## 4.3.2   Real-time Operating Room Planning Model under Surgery Duration Uncertainty

In this section, we will extend the real-time surgery management model (presented in the previous section) by taking into account the uncertainty in surgery durations. We assume that the duration of surgery type $m = 1, \cdots, M$ follows a probability distribution, $\text{Pr}_{[m]}$ with an expected value $d_m$. The expected completion time $\mathbf{C^0}$ at $t = 0$ is initialized as the expected duration of the surgeries of elective patients. We represent the overall duration and the type of surgery under operation during period $t$ as $\tilde{D}^t$ and $\overline{m}^t$, respectively.

Let's consider a surgery that is expected to be completed in the next epoch, $t + 1$: that is $C_i^t = t + 1$. The decision-maker collects information from the surgery crew whether the current surgery will be delayed or completed on time as expected. Let $\delta^t$ denote whether there is any delay ($\delta^t = 1$) or not ($\delta^t = 0$) in the surgery under operation at $t$. Note that the completion time is just an estimation and not possible to know the real surgery duration until it is realized. Thus, it is sensible to consider only one time period as the duration of delay. If a delay occurs in the current period, then the decision-maker is informed again in the next period whether the surgery will be completed or not. This process continues until the surgery is completed. We make the following assumptions regarding with the surgery duration uncertainty:

- Even if a surgery is finished earlier than the initial expected duration, the next surgery cannot start immediately due to pre-surgery operations. Also, the surgery crew for the next surgery may be unavailable until the planned start time. For this reason, we assume that the schedule remains the same when a surgery is finished earlier than expected. In other words, we do not consider the cases where a surgery finishes before its expected completion

time.

- The uncertainty regarding with the delay of a surgery is revealed only when it is expected to be completed in the next epoch, because the surgery duration is realized only when it is finished. The probability of delay in a surgery depends on the duration passed since the surgery has started.

Since the probability of delay in a surgery depends on the type of the surgery, we need to track the types of scheduled surgeries by inserting them into the state definition. Let $Q_i^t \in \{1, \cdots, M\}$ represent the type of surgery $i$ in epoch $t$ and $\mathbf{Q}^t = (Q_1^t, \cdots, Q_I^t)$ denote a vector of all surgeries at time $t$. In addition to the list of surgery types, we need to track the information regarding the delay since it affects the expected completion time list, and thus, the optimum action. We make the following changes for the description of the MDP formulation with uncertain surgery durations:

- **System state** at epoch $t$ consists of $(\mathbf{C}^t, \mathbf{Q}^t, m^t, \delta^t)$, where $\delta^t$ denotes whether there is a delay or not for the surgery in progress at period $t$.

- The type of surgery under operation at period $t$, $\overline{m}^t$, is identified by using information $\mathbf{Q}^t$, as $\overline{m}^t = \sum_{i=1}^{I} b_i^t Q_i^t$, that is used to identify the probability of delay in the current surgery. Similar to the completion times, we need to update the list of surgery types by inserting the type of the accepted non-elective:

$$Q_i^{t+1} = m^t x^t z_i^t + Q_i^t, \; i = 1, \cdots, I, \; t = 1, \cdots, T-1. \qquad (4.11)$$

- In case of a delay, the completion times of the waiting surgeries should be

134

increased by one period. Thus, $C_i^{t+1,y}$ is reformulated as follows:

$$C_i^{t+1,y} = \left[1 - (a_i^t - b_i^t)\right]C_i^t + b_i^t\delta^t + (a_i^t - b_i^t)\left[C_i^t + \delta^t + x^t d_{m^t}\right] + x^t z_i^t[d_{m^t} + K^t + \delta^t],$$
$$i = 1, \cdots, I, \ t = 1, \cdots, T - 1. \tag{4.12}$$

- As mentioned before, the probability of delay on the surgery under operation in period $t$ also depends on the duration passed since this surgery has started. Let $\overline{D}^t \in \{1, \cdots, d^{max}\}$ represent the time passed since the surgery under operation at epoch $t$ has started. $\overline{D}^{t+1}$ is 1 if a surgery is finished in epoch $t$. However, when there is a delay, the duration of the current surgery is increased by one period as $\overline{D}^t + 1$:

$$\overline{D}^{t+1} = \begin{cases} \overline{D}^t + 1, & \text{if } b_i^{t+1} = b_i^t, \ i = 1, \cdots, I, \\ 1, & \text{otherwise.} \end{cases} \tag{4.13}$$

This is due to the fact that when a surgery finishes, the next one, i.e. the surgery under operation for the next epoch, starts immediately. Then, the probability of delay at period $t$ can be defined as

$$\Pr(\tilde{\delta}^t = 1) = \begin{cases} \Pr\left(\tilde{D}^{t+1} > \overline{D}^t + 1 | \tilde{D}^t \geq \overline{D}^t\right), & \text{if } \overline{D}^t \geq d_{\overline{m}^t} - 1, \\ 0, & \text{otherwise,} \end{cases} \tag{4.14}$$

which is a discrete function since the surgery durations are defined as discrete units.

Given an initial state $(\mathbf{C^0}, \mathbf{Q^0})$, the objective is to minimize the value

135

function at epoch $t = 0$, which can be written as

$$v_0(\mathbf{C^0}, \mathbf{Q^0}) = \sum_{m=1}^{M} \Pr(\tilde{m}^1 = m)\Big[\big(1 - \Pr(\tilde{\delta}^1 = 1)\big)v_1(\mathbf{C^1}, \mathbf{Q^1}, m, 0) + \Pr(\tilde{\delta}^1 = 1)v_1(\mathbf{C^1}, \mathbf{Q^1}, m, 1)\Big].$$
(4.15)

Finally, for $t = 1, \cdots, T - 1$, the optimality equation (4.9) is reformulated as follows:

$$\begin{aligned}
v_t(\mathbf{C}^t, \mathbf{Q}^t, m^t, \delta^t) &= \min_{x^t \in \{0,1\}} \left\{ \sum_{m=1}^{M} \Pr(\tilde{m}^{t+1} = m)\mathbb{E}\Big[v_{t+1}\big(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \tilde{\delta}^{t+1}\big) + \eta(x^t, \mathbf{C}^t, m^t)\Big] \right\}, \\
&= \min_{x^t \in \{0,1\}} \left\{ \sum_{m=1}^{M} \Pr(\tilde{m}^{t+1} = m)\Big[ \Pr\big(\tilde{\delta}^{t+1} = 1\big)v_{t+1}\big(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, 1\big) \right. \\
&\qquad \left. + \big(1 - \Pr\big(\tilde{\delta}^{t+1} = 1\big)\big)v_{t+1}(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, 0) + \eta(x^t, \mathbf{C}^t, m^t)\Big] \right\}. \quad (4.16)
\end{aligned}$$

Note that we need to consider the possible delay in the current surgery when computing the set of future states. Since the random variables representing the type of new non-elective arrival and possible delay in the ongoing surgery are independent, we can separate the probabilities of these two cases in order to calculate the overall probability of having a possible future state. Because there are only two possibilities regarding the delay in the current surgery, the future states in the next period can be written explicitly as $\big(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, 1\big)$ and $\big(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, 0\big)$ for all possible surgery types $m \in \{1, \cdots, M\}$. The value function in the final period remains the same as in (4.10):

$$v_T(\mathbf{C}^T) = \theta^o \max \left\{ \Big( \max_{i=1,\cdots,I} \{C_i^T\} - T \Big), 0 \right\}. \tag{4.17}$$

The dynamic programming formulation of the surgery room planning problem under uncertain surgery durations can be solved by using the backward value iteration, similar to the MDP formulation presented in the previous section. Recall

136

that the only difference is observed when identifying the set of all possible future states. In other words, we have to update $\mathbf{Q}^t$ (in addition to $\mathbf{C}^t$) and also consider all possible cases regarding the delay in the current surgery (in addition to the non-elective types).

### 4.3.3 Real-time Multiple Operating Room Planning under Surgery Duration Uncertainty

The MDP formulation of the single operating room planning problem presented in Section (4.3.2) can be extended to the dynamic multiple operating room planning problem. We consider $R$ surgery rooms that are represented with $r = 1, \cdots, R$. In this case, we need to ensure that a non-elective patient arrived during a period is assigned to at most one room. For this purpose, the feasible set of actions is defined as $\mathcal{X} = \left\{ x_r : \sum_{r=1}^{R} x_r \leq 1, x_r \in \{0,1\}, r = 1, \cdots, R \right\}$. In addition, all decision variables have an additional index representing the room index that the corresponding variable belongs to. Let the uncertain parameter $\tilde{\delta}_r^{t+1}$ represent whether there exists any delay or not in room $r$ in period $t+1$. We denote the delay information for all operating rooms at time period $t+1$ by $\tilde{\boldsymbol{\delta}}^{t+1} = \left( \tilde{\delta}_1^{t+1}, \cdots, \tilde{\delta}_R^{t+1} \right)$. The value function presented in (4.16) can be reformulated as follows:

$$v_t(\mathbf{C}^t, \mathbf{Q}^t, m^t, \boldsymbol{\delta}^t) = \min_{\boldsymbol{x}^t \in \mathcal{X}} \left\{ \sum_{m=1}^{M} \Pr(\tilde{m}^{t+1} = m) \mathbb{E}\left[ v_{t+1}\left( \mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \tilde{\boldsymbol{\delta}}^{t+1} \right) \right] + \eta(\boldsymbol{x}^t, \mathbf{C}^t, m^t) \right\}.$$

Let $\boldsymbol{\delta} = \left( \delta_1, \cdots, \delta_R \right)$ denote a possible scenario regarding the delays in the current surgeries. Then we can compute the expected value function at $t+1$ as $\mathbb{E}\left[ v_{t+1}\left( \mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \tilde{\boldsymbol{\delta}}^{t+1} \right) \right] = \Pr(\tilde{\boldsymbol{\delta}}^{t+1} = \boldsymbol{\delta}) v_{t+1}\left( \mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \boldsymbol{\delta} \right)$. Note that the probability of delay in room $r$, $\Pr(\tilde{\delta}_r^{t+1} = \delta_r)$, is independent from possible delays that may occur in other rooms. Therefore, we can calculate the overall probability of having $\boldsymbol{\delta} = \{\delta_1, \cdots, \delta_R\}$:

$$\Pr(\tilde{\boldsymbol{\delta}}^{t+1} = \boldsymbol{\delta}) = \prod_{r=1}^{R} \Pr(\tilde{\delta}_r^{t+1} = \delta_r).$$

Finally, we obtain the value function as

$$v_t(\mathbf{C}^t, \mathbf{Q}^t, m^t, \boldsymbol{\delta}^t) = \min_{\boldsymbol{x}^t \in \mathcal{X}} \left\{ \sum_{m=1}^{M} \Pr(\tilde{m}^{t+1} = m) \sum_{\boldsymbol{\delta} \in \{0,1\}^R} \Pr(\tilde{\boldsymbol{\delta}}^{t+1} = \boldsymbol{\delta}) v_{t+1}\big(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \boldsymbol{\delta}\big) \right.$$
$$\left. + \eta(\boldsymbol{x}^t, \mathbf{C}^t, m^t) \right\}, \ t = 1, \cdots, T-1.$$

The formulations of the final period value function (4.17) and one-step action cost (4.7) do not change in the multi-rooms case. The optimum policy can be found with the backward value iteration, similar to the single room MDP formulations presented previously.

The traditional dynamic programming (DP) algorithm uses the backward value iteration where the optimal decisions and value functions are calculated iteratively starting from the final period and stepping backwards in time. Although this produces the exact solution, it is affected by the curse of dimensionality since the value function is computed at each state and all possible actions are evaluated. On the other hand, ADP is designed to reduce action and state spaces by adopting an approximation technique for the value function.

## 4.4 Simulation-based Approximate Dynamic Programming Approach

The real-time surgery management model under uncertainty with single room is computationally expensive to solve due to the large state space. For example, in a very small instance with 1 surgery room, 4 time periods, 2 initial electives, and 2 possible surgery types, the state space can be as large as 18,432. Besides, the state space increases by around 6 times as the number of rooms is increased to

two. Note that the action space consists of only two cases at each time period, and, thus, it is relatively small.

ADP overcomes the curse of dimensionality of traditional backward DP by using approximations. The ADP approach uses simulation to find approximate state values and policies. These estimates are generated based on Monte-Carlo simulation of state trajectories. The value functions are evaluated for all visited states and updated using an aggregation structure for states (such as a single entity in a lookup table) or regression models. Simulation-based ADP algorithm is very suitable to solve MDP problems (Powell, 2009). Thus, to solve the real-time operating room planning model under uncertainty, we consider a simulation based ADP method with double-pass and a lookup table. In this section, we provide the details of the proposed ADP algorithm. A linear programming based ADP is not applied since the value function (4.9) is complex (Powell, 2009). Also, we implemented a value iteration instead of a policy iteration based algorithm since the problem has a large state space and a comparatively small action set (Sun et al., 2013). Note that the overtime cost is realized at the end of the planning horizon. Therefore, the overtime cost should be added into the values of the previous states in the simulated trajectory. Thus, we need a double pass to update the state values at each iteration rather than a single-pass.

An initial list of completion times $\mathbf{C}^0$ and surgery types $\mathbf{Q}^0$, as well as the probability distributions for the non-elective arrival types and surgery durations are given as the inputs to the algorithm. Note that the initial state is the same in all iterations. As we run the algorithm, each visited state and its approximate value are inserted into the lookup table. The simulation based ADP algorithm has five main modules (Powell, 2007): initialization, generator, simulator, decision generator and value function approximator. The details of these modules are explained below.

- Initialization: Let $n$ denote the iteration counter and $N$ be the maximum number of iterations. We set $N$ and initialize $n = 1$. The value of the initial state $(\mathbf{C}_1^0, \mathbf{Q}_1^0)$ is defined as zero.

- Generator: Based on the probability distributions of non-elective types, at each iteration $n$, a sample path for the type of non-elective arrivals, $\mathbf{m}_n$, is randomly generated. In other words, a random non-elective type (including the possibility of no arrival) for each time period $t = 1, \cdots, T - 1$ in the sample path is created.

- Simulator: At each iteration $n$ and time period $t = 1, \cdots, T - 1$, the algorithm simulates an occurrence of a delay, $\delta_n^t$, based on $\mathbf{C}_n^{t-1}, \mathbf{Q}_n^{t-1}$ and the probability distributions of surgery durations.

- Decision generator: For the generated state $(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t)$, it finds the greedy action $x_n^t \in \{0, 1\}$ as well as the state value, represented with $\hat{v}_t^n$, by using the optimality equation (4.16) and the approximate values stored in the lookup table. If the value of a future state does not exist in the lookup table, then it is simply assumed to be 0. Then $\mathbf{C}_n^{t+1}, \mathbf{Q}_n^{t+1}$ can be computed based on the greedy action $x_n^t, \mathbf{C}_n^t$, and $\mathbf{Q}_n^t$.

- The value function approximator: This module does the second pass and stores the visited states and the computed values of these states into the lookup table.

Algorithm 2 shows the pseudo-code of the ADP algorithm with a value iteration, lookup table and double pass for the MDP formulation with single room and uncertain surgery durations. In each iteration $n = 1, \cdots, N$, after all states in the planning horizon are visited, the algorithm goes backward in time and recursively adds the values of the future states (in the sample path) into $\hat{v}_t^n$ for

$t = T - 1, \cdots, 1$. If a state, $(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t)$, is visited for the first time by the algorithm, then its computed value $\hat{v}_t^n$ is directly added to the lookup table. This approximate value stored in the lookup table is shown with $\overline{V}_t^n(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t)$. Otherwise, $\overline{V}_t^n(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t)$ is computed by summing $\hat{v}_t^n$ and the value of that state most recently stored in the lookup table (shown with $\overline{V}_t^k(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t)$ in the pseudo-code), after weighting them by a smoothing parameter represented with $\alpha_n$. Since the state values are expected to approach to their exact values through iterations, $\alpha_n$ is formulated as a positive linear function of $n$. The linear form is selected because it is simple and also converges eventually (Powell, 2007).

Note that the nature of the problem allows multiple optima since different actions may result in the same state value. In order to increase the number of explored states (visited by the algorithm), we employ different exploration strategies. For instance, we randomly select a greedy action in case of multiple optima. As shown in computational experiments section (4.5), we also consider a random strategy for the action selection in order to improve the exploration process. Specifically, we check whether a randomly generated number is lower than a fixed constant, represented with $\Gamma$. If so, then the greedy action is randomly selected among the feasible actions. Otherwise, it is selected randomly among the optimum actions. However, this strategy may result in suboptimal policies and decreases the exploitation (that is defined as the degree of approach of the approximate values to the exact ones). Therefore, we only apply it for the first half of the iterations, i.e. for $n = 1, \cdots, N/2$.

In addition to the lookup table based ADP algorithm, we consider other approaches such as a basis function based ADP algorithm. The basis function based ADP algorithm can be used to estimate the values of the states that are not visited by the lookup table. There are two main differences between the lookup table and basis function based algorithms. Instead of using the value

---

**Algorithm 2** Pseudo-code of the double-pass ADP algorithm with a lookup table

---

**Step 0.** Set $n = 1$ and maximum number of iterations $N$ and initialize the value of $(\mathbf{C}_n^0, \mathbf{Q}_n^0)$, i.e. $\overline{V}_0^1$, .

**Step 1.** Generate a sample path of $\mathbf{m}_n$.

**for** $t = 1, 2, \cdots, T - 1$, **do**
   Generate $\delta_n^t$ based on $\mathbf{C}_n^{t-1}, \mathbf{Q}_n^{t-1}$.
   Generate a random number $\gamma$ and find the greedy action $x_n^t$ and $\hat{v}_t^n$ by,
   **if** $n \leq N/2$, and $\gamma \leq \Gamma$ **then**
     Randomly select $x_n^t$ among the feasible action set $\{0, 1\}$.
   **else**
     Solve (4.16) based on;
     **if** the value function of state $(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \delta)$ exists in the lookup table, represented with $v_j(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \delta)$ **then**
       $v_{t+1}(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \delta) = v_j(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \delta)$.
     **else**
       $v_{t+1}(\mathbf{C}^{t+1}, \mathbf{Q}^{t+1}, m, \delta) = 0$.
     **end if**
   **end if**
   Update the state variables,

$$\mathbf{C}_n^{t+1} = \mathbf{C}^{t+1}(\mathbf{C}_n^t, x_n^t, m_n^t, \delta_n^t) \text{ and } \mathbf{Q}_n^{t+1} = \mathbf{Q}^{t+1}(\mathbf{Q}_n^t, x^t, m_n^t)$$

**end for**
**for** $t = T - 1, \cdots, 1$ **do**
   Compute;

$$\hat{v}_t^n = \eta(\mathbf{C}_n^t, x_n^t, m_n^t) + \hat{v}_{t+1}^n(\mathbf{C}_n^{t+1}, \mathbf{Q}_n^{t+1}, m_n^{t+1}, \delta_n^{t+1})$$

   **if** the state, $(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t)$, is visited at iteration $k < n$, **then**
     Update the stored value of this state, $\overline{V}_t^k$, as

$$\overline{V}_t^n(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t) = (1 - \alpha_{n-1})\overline{V}_t^k(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t) + \alpha_{n-1}\hat{v}_t^n,$$

   **else**
     $\overline{V}_t^n(\mathbf{C}_n^t, \mathbf{Q}_n^t, m_n^t, \delta_n^t) = \hat{v}_t^n$.
   **end if**
**end for**

**Step 2.** $n := n + 1$. If $n \leq N$, go to Step 1. Otherwise, go to Step 3.

**Step 3.** Return the value function approximations $(\overline{V}_t^n)$ for $t = 1, \cdots, T$ and $n = 1, \cdots, N$.

---

approximations stored in a lookup table, an on-policy basis function method uses the basis function approximations to find the greedy action. Also, instead of updating the lookup table at the end of each iteration, it updates the weights of the basis functions based on the state values computed at that iteration. To find the best set of basis functions, we run the algorithm with various possible basis functions. The experiments conducted for the basis function based ADP algorithm are presented in the next section.

## 4.5    Computational Experiments

In this section, we first describe the design and data structure used for the numerical experiments and then present the computational results of the ADP algorithm applied to the real-time operating room planning model under surgery duration uncertainty.

We design a series of computational experiments in order to illustrate the performance of the ADP algorithm. In particular, we are concerned with finding the approximate policies for the operating room allocation with different parameter settings by using generated data inspired from real data. A myopic heuristic is selected as a benchmark to compare with the performance of the ADP algorithm. The computational experiments also aim to investigate the effect of different elective scheduling strategies on the overall cost using generated data. All the experiments are conducted by using the ADP algorithm with lookup table for the MDP model with single room and uncertain surgery durations, unless stated otherwise. All computational experiments are carried out on a laptop with Windows XP operating system, CPU 2.26 GHz Intel Corei5 and 8 Gb of RAM.

### 4.5.1 Parameter Selection and Modelling Features for the ADP Algorithm

We first aim to investigate the structural features of the ADP algorithm and illustrate the impact of parameter selection and action selection strategies on its performance. The performance of the ADP algorithm is measured by two criteria: the convergence speed and optimality gap. The *convergence speed (rate)* is defined as the number of iterations required to achieve the stability in state values. The stability refers to a state value staying within 10% of the exact value (Hulshof et al, 2013). The *optimality gap* is defined as the average relative difference between the optimum and approximate values of each state visited by the algorithm. Note that, for the final states (at the end of the decision horizon), the approximate and optimum values are the same since they do not include the evaluation of the future expected value functions. Thus, for a more rigorous comparison, we exclude the final states from the optimality gap calculation.

A **small problem instance** to be used for this experiment consists of a single operating room with

- 4 time periods while the daily time limit is 7 time periods,

- 2 initial electives with completion times [1, 2],

- a non-elective surgery having $\{0,1,2\}$ durations with equal probabilities, and,

- two surgery types.

In this experiment, all cost parameters are fixed to unity. The dynamic programming formulation for this small problem instance involves 4484 possible states and the optimal policy is obtained by using the backward value iteration method

within 3 hours.

***Smoothing parameter*** *($\alpha_n$)* is defined as a linear function of iteration counter $n$ as $\alpha_n = a + b \times n$. Note that $a$ and $b$ are continuous and can take a wide range of values. In order to establish the effect of the smoothing parameter on the convergence speed, we test the algorithm with different values of $a$ and $b$. We initially conducted some preliminary trial-and-error experiments and selected only four instances to present for an illustrative purpose. In these cases, $\alpha_n$ varies within intervals $[0.6, 0.9], [0.3, 0.9], [0.1, 0.9]$ and also takes a fixed value of $0.5$. We only show the value of the initial state since the smoothing parameter affects the convergence of other states in a similar fashion. Figure 4.2 presents the value of the initial state with respect to the number of iterations for different smoothing parameters. The black solid line displays the optimum (exact) state value. We observe that when $\alpha_n$ varies in $[0.6, 0.9]$ and is fixed at $0.5$, the state value is not within 10% of the optimum value. The convergence is obtained when $\alpha_n$ varies within $[0.3, 0.9]$ and $[0.1, 0.9]$.

***Random greedy action selection from the feasible set*** *($\Gamma$)*: As explained before, in order to increase the exploration of the algorithm, we added randomness into the action selection based on parameter $\Gamma$. In order to investigate the effect of $\Gamma$ on the performance of the ADP algorithm, we run the algorithm with fixed $\Gamma$ levels as $0.3$ and $0.5$ as well as zero. Note that the case of $\Gamma = 0$ corresponds to no randomness in the action selection procedure.

Figure 4.3 presents the performance of the algorithm in terms of the number of states explored as well as the optimality gap, that is obtained at fixed levels of $\Gamma$ and number of iterations. During these runs, $\Gamma$ is set to 0 for the second half of the iterations to prevent the suboptimality. As these results indicate,

(a) $\alpha_n$ varies in [0.1, 0.9]    (b) $\alpha_n$ varies in [0.3, 0.9]

(c) $\alpha_n$ varies in [0.6, 0.9]    (d) $\alpha_n$ is fixed at 0.5
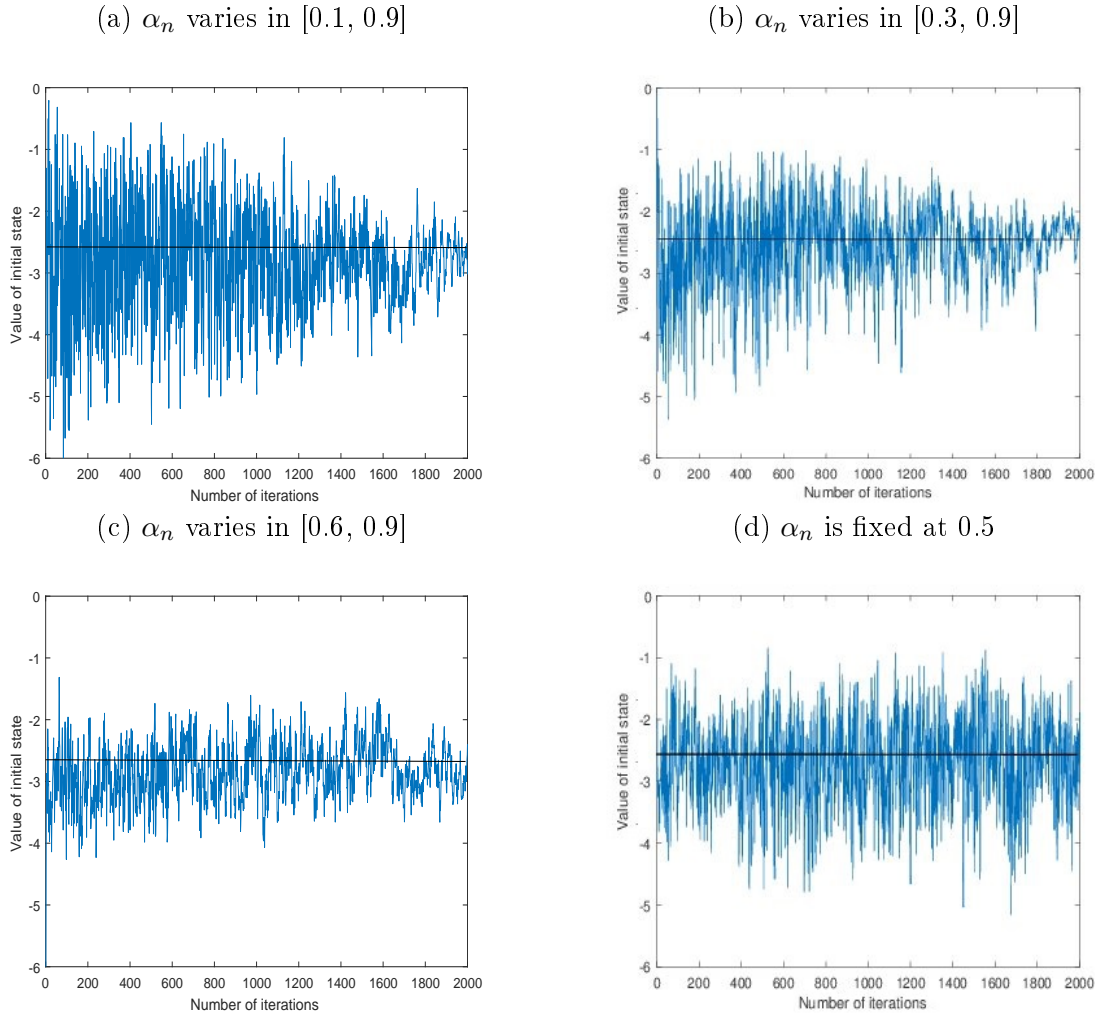
Figure 4.2: Value of the initial state at different number of iterations obtained by various smoothing parameters

the number of explored states increases considerably when the action selection is affected by randomness (i.e. for $\Gamma > 0$). The smallest optimality gap is obtained when $\Gamma = 0.3$. However, the optimality gap increases for $\Gamma > 0.3$. In other words, as more states are explored, the degree of approach to the true state values (exploitation) decreases. This is a well-known phenomenon within the ADP literature, called as exploration vs. exploitation trade-off (Powell, 2007).

***Basis function approximation approaches:*** The results displayed so far are obtained by using a lookup table based ADP algorithm. We also intend to estab-
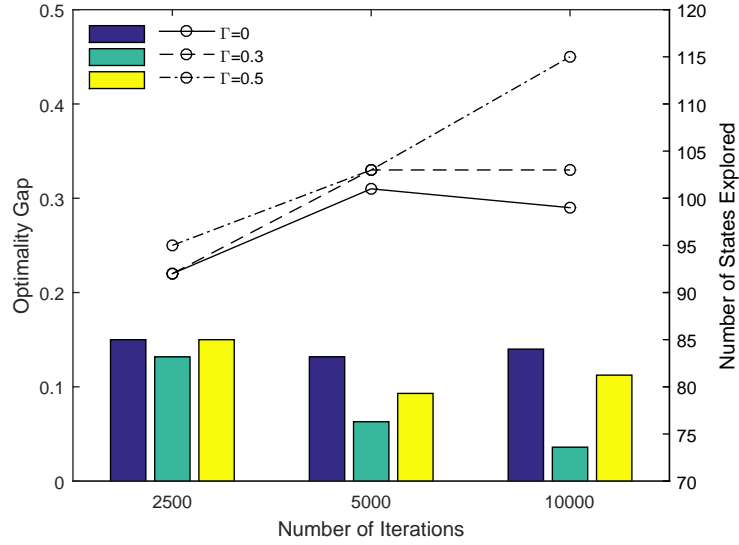
Figure 4.3: Optimality gap (columns) and number of states explored (lines) by the ADP algorithm by using different Γ levels and number of iterations

lish the impact of different approximation approaches such as basis function with respect to the lookup table. We first implemented an on-policy basis function algorithm with different basis function structures such as the number of surgeries waiting, time period, and type of non-elective arrival. However, our results show that the state values obtained by on-policy basis function approximations do not converge even with a wide range of different basis functions, as also encountered in the literature (for instance, see Powell, 2007).

We also considered an off-policy basis function aggregation. A brief description of this approach is follows. First, we obtain the approximate state values with the lookup table based ADP algorithm, then, apply regression on these approximate values to find the best weights of the basis functions. We tested several basis functions and found the best fit for *the number of waiting surgeries, the current time period, the duration of the recent non-elective arrival* and *occurrence of possible delay*. In other words, the approximate value function at state $(\mathbf{C}^t, \mathbf{Q}^t, m^t, \delta^t)$ is formulated as

147

$$\overline{V}_t(\mathbf{C}^t, \mathbf{Q}^t, m^t, \delta^t) = w_1 \sum_{i=1}^{I}(a_i^t - b_i^t) + w_2 t + w_3 d_{m^t} + w_4 \delta^t, \qquad (4.18)$$

where $w_j$ for $j = 1, \cdots, 4$ are the weights of the basis functions, respectively. The best weight levels are found as $\mathbf{w} = [0.7, 0.24, 0.88, 0.85]$ with the regression fitting tool in MATLAB. To validate the accuracy of the regression fitting, we computed the approximate values of 1000 sampled states by using the basis function structure (4.18) as well as the lookup table approximation. The t-test showed that the basis function approximations and the state values in the lookup table are not statistically different (with a p-value of 0.39). In other words, the off-policy basis function method based on (4.18) approximates the value functions as good as the lookup table method. Therefore, we conclude that the off-policy basis function approach based on (4.18) can be used to approximate the values of the states that are not visited with the lookup table method.

**A Brief Summary of Findings:** The experiments presented in this section indicate that the performance of the ADP algorithm highly depends on the choice of the algorithm's parameters that need to be tailored according to the underlying problem. Our main findings to be used for the remaining of the experiments are summarized as follows.

- The smoothing parameter should be selected within an interval [0.1-0.9] or [0.3-0.9] in order to achieve the best approximation to the state values.

- The exploration and exploitation are well balanced for $\Gamma = 0.3$.

- An off-policy basis function structure based on *the number of waiting surgeries, the current time period, the duration of the recent non-elective arrival* and *occurrence of delay* provides the best value function approximation for the real-time surgery management problem.

**Performance Comparison of Approximate and Backward DP Approaches:**

As mentioned before, we cannot obtain an optimal policy for a real-sized problem instance by using the traditional backward value iteration approach (as an exact algorithm) since it suffers from the curse of dimensionality. For instance, recall that the problem instance (as allocation of 2 types of patients to a single operating room over 4 time periods) consists of 4484 possible states, and can be solved by the exact method (backward value iteration) in 3 hours. On the other hand, the ADP algorithm with the best selection of parameters identified can solve the same problem in less than a minute. Moreover, we reach 96% of match between the actions determined in the optimal and greedy policies (by using the exact and ADP approaches, respectively) over all states. This can also be interpreted as follows: the ADP approach misses the optimal action only for 4% of all states. Also, the optimality gap (the average relative difference between the state values found by two approaches) is 5%.

## 4.5.2 Performance Comparison of the ADP Algorithm and A Myopic Approach

In this section, we illustrate the performance of the approximate policy obtained by the ADP algorithm using generated data inspired from real data. For performance comparison purposes, we implement a myopic heuristic that is used as a benchmark strategy in MATLAB. The myopic strategy selects a greedy action based on the lowest one-step-cost, $\eta(\mathbf{C}^t, x^t, m^t)$, for each non-elective arrival. However, it does not take into account the expected future cost.

We design computational experiments for a single room using generated data inspired from the publicly available sources (University of Twente, 2017). A description of the data set is presented in Table 4.2. This data set is referred as the base data in the rest of the experiments. One time period (slot) is assumed to be half an hour which is found to be reasonable for at most one non-elective

arrival. The cost parameters are not available in the real data set, thus, they are estimated based on the literature sources (Zonderland et al., 2010). The statistical analysis on the real data shows that log-normal distribution fits well to the surgery durations as also stated in the literature; for instance, see Strum et al., (1998), Strum et al., (2000), and May et al. (2000). The first two moments of log-normal distributions for each type of patient, that are also estimated from the real data, are used to compute the delay probabilities for surgery types. Both the myopic heuristic and the ADP algorithm are run with 1000 sample paths generated by using the real data. Next, we analyze the performances of the myopic and approximate policies in terms of the solution time and the overall cost.

**Efficiency of Solution Approaches:** We can report that the CPU time taken to obtain a policy by the ADP algorithm is about 30 minutes. On the other hand, the myopic heuristic provides a greedy action within a couple of seconds. As the nature of methods, the ADP policy provides the greedy actions for every decision epoch of the whole planning horizon whereas the myopic heuristic, as a reactive approach, is run whenever a non-elective patient arrives to the hospital in the generated scenarios. Still, the average CPU time taken to solve a problem instance by the myopic heuristic is much smaller than the one by the ADP algorithm. On the other hand, the overall cost (the value of the initial state) computed via the myopic approach is much higher than the one obtained by the ADP algorithm as presented in the next set of experiments.

Table 4.2: Input data

| Description of parameters | Fixed level |
|---|---|
| Type of surgeries | 3 |
| Initial elective schedule (completion times in terms of slot) | [2, 4, 6, 10, 14, 16] |
| Mean (and std. deviation) of surgery duration distribution for each surgery type | [2 (1), 4 (1), 6 (1)] |
| Overtime threshold (slots) | 16 |
| Daily time threshold (slots) | 22 |
| Cost of cancellation, rejection, overtime, waiting time (£) | [1, 1, 1, 0.2] |
| Non-elective arrival rates for each surgery type (surgery/slot) | 0.125, 0.1875, 0.0625 |

**Overall Cost:** In order to show how each cost component, namely overtime, rejection, cancellation and waiting time, affects the overall cost obtained by the policies, we consider 3 cases for each cost component where only one cost component is changed at a time while the others are kept at their base levels as given in Table 4.2. For those cases labelled as 'Decreased and Increased', the corresponding cost component is decreased and increased, respectively, by 50% from its base level. The other case (labelled as 'Base-case') uses the same cost levels as given in Table 4.2. Note that the cost of waiting is smaller than the others. For that reason, when it is increased and decreased by 50%, we could not observe a significant change in the overall cost. Therefore, the waiting time cost component is varied in a larger interval so that the decreased and increased levels correspond to 0.01 and 1, respectively.

Figure 4.4 shows the box plots for the overall costs obtained with the ADP policy and the myopic heuristic for different levels of the cost components. The statistical analysis shows that the ADP policy performs significantly better than the myopic heuristic in most of the cost levels, except when the cancellation cost is low. As the level of a cost component increases, the difference between the overall costs obtained by the myopic and ADP policies gets larger. Another observation is that the rejection cost coefficient has the biggest effect on the overall cost, while

the waiting time cost coefficient has the smallest effect.

**Structural Features of the Approximate and Myopic Policies:** We analyze the structural characteristics of both (ADP and myopic) policies in terms of the acceptance rates (of non-electives) during the planning horizon and the effect of the non-elective duration on the acceptance decision. With this experiment, we aim to derive insights of different policies, such as when and what type of non-elective patients to accept, that will support the decision-making process of the operating room manager. The base data are used to obtain both policies.

*Non-elective admissions:* To investigate the effect of the (non-elective) arrival rates on the acceptance decision, we considered two cases where the arrival rates are doubled and reduced to the half of the base rate: those cases are abbreviated as 'overloaded' and 'underloaded' cases, respectively. Figure 4.5 shows the relative frequency of the accepted non-elective patients at each time period that is measured as the number of iterations where the non-elective arrival is accepted divided by the overall number of iterations with a non-elective arrival, both at the corresponding time period. The graphs show that the ADP policy tends to accept more non-electives at early and later time periods during the day and more likely to reject towards the mid-day. When the arrivals are overloaded, the acceptance frequencies do not change much, unlike the underloaded case where the acceptance rates around the mid-day are larger than the base case. We can interpret this situation as follows. For the base arrival case, there are more available slots for non-electives to schedule during the early time periods. As more non-electives are accepted early in the day, the schedule is filled up quickly. Consequently, the risk of delay increases towards the mid-day and rejections become more likely to be realised. On the other hand, when the number of waiting surgeries and the risk of overtime decrease over time, the ADP policy starts accepting more non-electives arriving towards the end of the day.
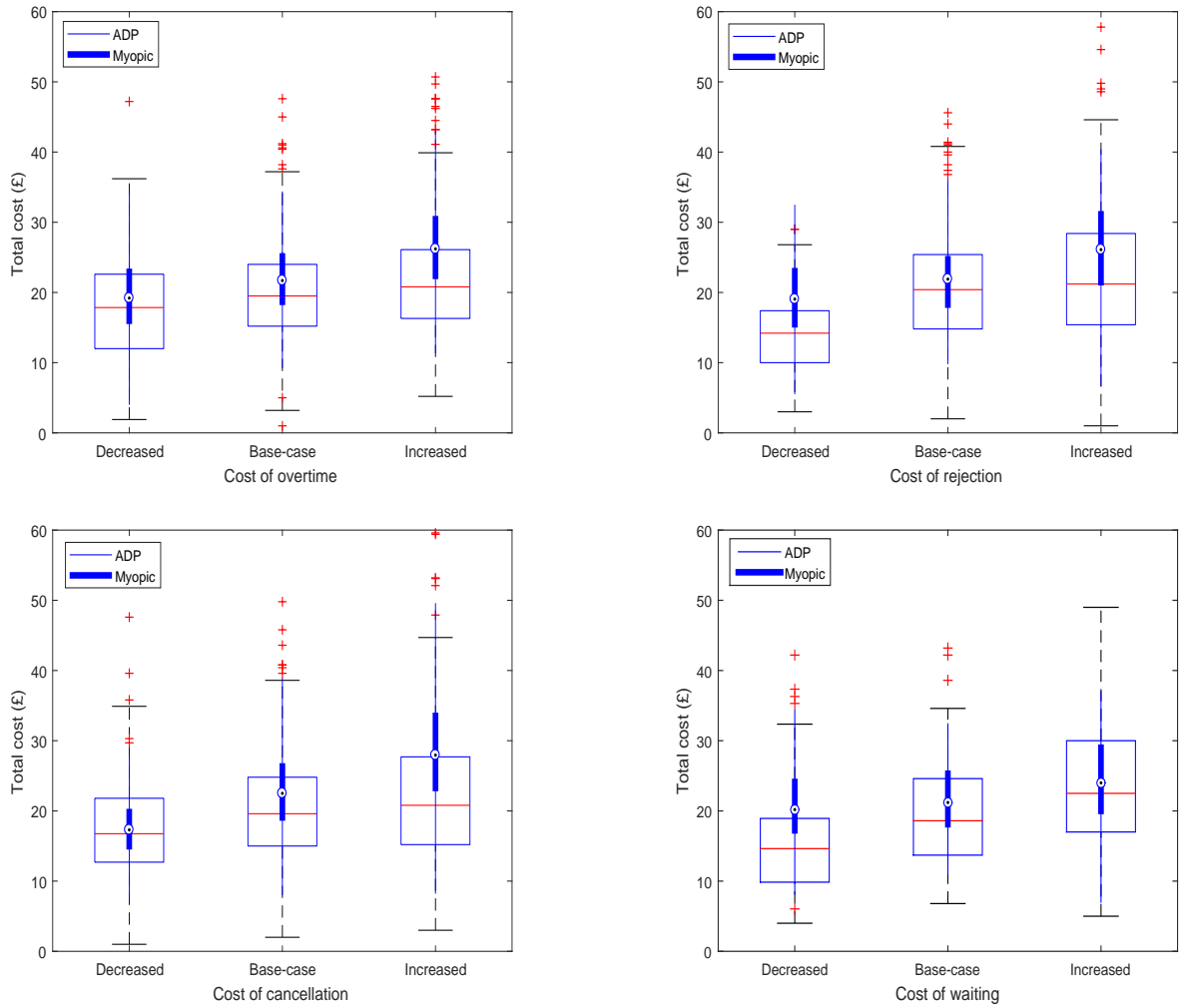
152

Figure 4.4: Box plots for the total costs obtained with the ADP policy and the myopic heuristic in different levels of cost components
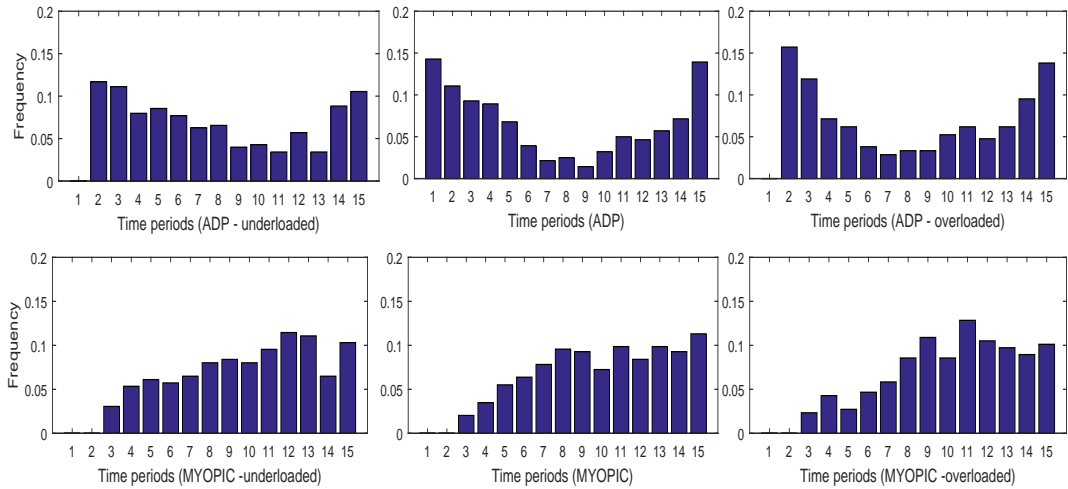
Figure 4.5: Frequency histograms for time periods with an acceptance decision by ADP and myopic policies under different arrival rates

As opposed to the ADP policy, the myopic policy rejects the non-electives arrived in the earlier periods and usually accepts them afterwards. This is because, in the earlier periods, the cost of acceptance is higher than that of rejection due to a large waiting cost. However, after a certain time period, the cost of acceptance gets lower than that of rejection, therefore, the non-electives are accepted without considering the future risk of overtime.

*Maximum completion time:* We also investigate whether the latest surgery completion time in the schedule (so-called maximum completion time) can be used in practice as a decision-making rule for the acceptance or rejection decision of non-elective patients. For this, we compare the maximum completion times when an acceptance and rejection decisions are made in the ADP and myopic policies. The statistical significance tests show that, both policies tend to accept the non-elective patient if the maximum completion time is low at his/her arrival time. However, it is not possible to identify a threshold maximum completion time level over which the non-electives are always rejected.

*Surgery duration of non-elective:* Intuitively, the acceptance or rejection

154

decision for a non-elective may depend on the type of surgery as well as the expected surgery duration. The statistical analysis shows that the average acceptance rates in both policies are the same if the non-elective patient possesses an expected surgery duration of 2 and 4 time slots. On the other hand, if the expected surgery duration is about 6 time slots, the acceptance rates drop by 42% compared to other types of surgeries.

### 4.5.3 Impact of Various Elective Scheduling Strategies

In practice, hospitals may prefer to use different elective scheduling strategies. In this section, we investigate the effect of alternative elective scheduling policies on the overall cost.

*Strategies based on separating/combining surgery types:* When there are multiple operating rooms available, the surgery suite manager may apply different strategies for assigning elective patients. For instance, the same type of surgeries may be scheduled to the same operating room using a 'divided' strategy. Alternatively, the same type of surgeries can be shuffled across the rooms, named as 'shuffled' strategy. In order to establish the possible impact of these elective scheduling strategies on the overall performance, we consider two operating rooms with different initial schedules. The initial schedules of surgery types in the first and second rooms are set as [1,1,1,1,1,1] and [2,2,2], respectively, for the divided strategy; [1,1,2,1,1] and [2,1,2,1], for the shuffled strategy. Note that the expected completion times of the rooms are fixed as 12 time slots in both strategies. Figure 4.6 shows the relative frequency of time periods in which a non-elective patient is accepted by the ADP policy in two strategies.

In the divided case, the policy initially tends to assign a non-elective into the second room, possibly because this room has a lower risk of delay due to smaller number of surgeries. Later during the day, the non-electives are assigned
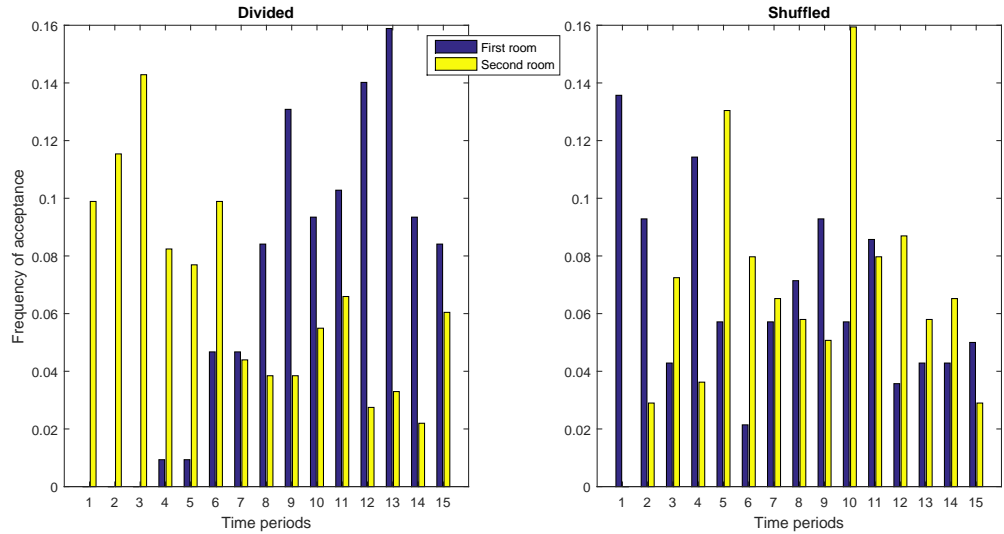
155

Figure 4.6: Frequency histogram for time periods with an acceptance decision in two rooms with different initial schedules

mostly to the first room since the schedule of the second room is already tight with the earlier non-elective additions. However, in the shuffled case, the time period of acceptance disperses evenly throughout the day. Moreover, the shuffled case results in 41% lower overall cost than the divided case.

*Leaving buffer times within the elective schedule:* An alternative elective scheduling strategy is to leave buffer times within the elective schedule to accommodate the non-elective arrivals (Van Riet and Demeulemeester, 2015). In this way, the elective schedule would be less interrupted, and, consequently, the waiting time of the electives would be reduced. In order to test the effect of such policy on the overall cost, we construct a new instance of the base data (given in Table 4.2) by considering the same initial schedule with two empty time slots (buffer) in the middle of the day. Recall that the schedule in the base data do not have any buffer. The ADP algorithm is run with two data instances (with and without buffer) by changing the level of one cost component and keeping the others at the same levels as in the base data.

Table 4.3 displays the (statistically significantly) better strategy (leaving vs. not leaving a buffer) in terms of the overall cost for varying cost component levels. If both strategies produce not (statistically significantly) different costs, then it is shown with 'ND' to represent 'no difference'. As it can be seen from Table 4.3, having no buffer provides better performance in terms of the overall cost in more cases. If the relative cost of waiting is as low as 0.01, leaving a buffer is worse than not leaving. As the cost of waiting is increased, leaving buffer is less costly.

Table 4.3: Performance comparison of both strategies with and without buffer at various levels of cost components

| Cost | Waiting | Overtime | Cancellation | Rejection |
|------|---------|----------|--------------|-----------|
| Doubled | With buffer | Without buffer | Without buffer | ND |
| Base-case | ND | ND | ND | ND |
| Half reduced | ND | ND | ND | Without buffer |

*Elective Sequencing Rules (schedule longest or shortest surgery first):* Next, we implement two elective sequencing rules, namely longest surgery first (LF) and shortest surgery first (SF) that have been widely applied in practice (Testi et al., 2007). There is no clear preference among these two strategies in practice but in the literature, it is claimed that SF produces smaller overtime and cancellations (Testi et al., 2007). We solve the operating room planning problem with LF and SF elective surgery sequencing rules using the base data except with the different initial schedules (completion times) of [6, 12, 14, 16, 18] for a LF schedule and [2, 4, 6, 12, 18], for a SF schedule.

Table 4.4 shows the better rule (LF vs. SF) that provides the smaller overall cost under different levels of cost components. When the overtime cost is doubled, the LF rule produces a lower overall cost. On the other hand, for the cases of high rejection and low waiting costs, the SF rule performs better than

the LF.

Table 4.4: Performance comparison of the LF and SF elective sequencing rules at various cost component levels

| Cost | Waiting | Overtime | Cancellation | Rejection |
|---|---|---|---|---|
| Doubled | ND | **LF** | ND | **SF** |
| Base-case | ND | ND | ND | ND |
| Half reduced | **SF** | ND | ND | ND |

*A Brief Summary of Findings:* Our findings observed from the second part of the computational results are summarized below:

- An acceptance of the non-elective patient at early and late periods during a day is less costly.

- There is no threshold latest completion time level after which any non-elective patient will be always rejected.

- Shuffling different types of surgeries among the operating rooms is less costly than assigning same type of surgeries into one room.

- The performance of the SF and LF rules for scheduling elective patients highly depends on the levels of cost components. On the other hand, there is no significant difference between different sequencing strategies if the cost levels are fixed as the original levels presented in Table 4.2.

## 4.6 Conclusions

Surgeries are the major sources of the costs and revenues in hospitals. The real-time management of a surgery suite is a dynamic problem affected by several uncertainties such as non-elective arrivals and surgery durations. It requires to

take several decisions related to acceptance/rejection of non-electives and cancellation of electives. In this chapter, we develop an MDP model for the real-time surgery management problem. However, due to 'curse of dimensionality', the real-sized instances cannot be solved to optimality. Thus, we apply a simulation-based ADP algorithm with lookup table and double-pass approaches to obtain approximately optimum policy for the real-sized instances. To obtain a better algorithm performance, different parameter settings are tested by using the optimum solution of a small instance. The numerical comparisons show that the exact and approximate policies coincide in 96% of the states.

We also consider the myopic heuristic as a benchmark for the performance comparison of the ADP approach. The computational results show that the approximate policy provides lower overall cost than the myopic policy produces. We also test the effect of different scheduling policies such as leaving buffers within the initial schedule or longest/shortest surgery first. The experiments show that leaving buffers may be beneficial if the cost of waiting is very high. The experiments also suggest that 'shortest surgery first' rule is better than 'longest surgery first' in most of the cost combinations.

The modelling framework introduced in this chapter can be extended to a longer planning horizon than one day. In this case, the decisions made on the admission of a patient today will have an impact on the next day's schedules. The problem would be larger and more complex in this case which may require further approximations to obtain a solution. Another possible extension is, instead of assigning an accepted non-elective to the first place in the schedule, the order of this non-elective in the schedule can be a decision variable. However, this would increase the size of the action space considerably.

# Chapter 5

# Conclusions

Healthcare processes are subject to several uncertainties such as patient arrivals, operation durations, test results, etc. To obtain robust solutions for healthcare management problems, the planning and management of healthcare processes should take these uncertainties into account. This thesis models and solves three healthcare decision-making problems under uncertainty. We use robust optimization, queuing theory, scenario-based modelling, and MDP to model these problems. This section concludes the thesis by summarizing the research and the main findings, mentioning several limitations encountered during the research, and finally providing some future research directions.

## 5.1   Summary of Research and Findings

In Chapter 2, we study the capacity planning problem in a network of stem-cell donation centres. The uncertainties in patient arrivals, results of blood tests, donor travel times, and the number of donors are incorporated into the model with a scenario-based approach. The advanced blood testing is modelled as a multi-server, first-come first-served queue with general interarrival and service time distributions. We consider the maximum waiting time in this queue since

the worst-case, the patient death, should be avoided as much as possible. The maximum waiting time is approximated with a robust queuing approach. The resulting non-linear integer programming model is reformulated into a linear one and solved with branch-and-bound. We design out-of-sample and in-sample experiments to investigate the real-life performance of the optimum solution. The approximate maximum waiting times calculated by the optimization model are very close to the simulated maximum waiting times, indicating accuracy of the approximation. The computational experiments show that increasing the number of stem-cell donation centres is more cost-effective. Also, the results indicate that increasing the budget more than a certain level does not affect the maximum waiting time in the network. Lastly, we analyze the service performance of the network for different budget levels and patient arrival rates.

Chapter 3 introduces the resource allocation problem in a healthcare outsourcing network. Given a fixed budget, a central healthcare authority needs to decide the capacities to outsource from available providers in several regions and allocate the patient demand in the network accordingly. Each service provider is modelled as a multi-server, first-come first-served queue where the patient arrivals and service durations are assumed to follow general distributions. The maximum waiting time in each provider is approximated with the robust queuing approach, same as in the second chapter. To solve the resulting non-linear integer programming model, we propose an alternating optimization based heuristic combined with the bisection search. In the computational experiments, we show that the heuristic performs better than the available commercial solvers especially for medium and large size instances. The sensitivity analysis provides several important managerial insights. First, the results suggest that larger regions with more providers is better than the smaller ones. Secondly, defining the budget for a larger network with more regions results in a smaller maximum waiting time.

161

These results are possibly due to the increased risk sharing with larger regions and networks. The final part of the experiments uses the real data obtained from an NHS audiology network. These experiments show that the current access times in the network can be reduced.

Chapter 4 focuses on the real-time management of a surgery schedule while non-elective patient arrivals and surgery durations are uncertain. We develop an MDP model with a single-day planning horizon where the action set consists of accepting or rejecting a non-elective arrival. The overall cost is a weighted sum of the costs of surgery cancellation, rejection, and waiting time and staff overtime as well. The real-sized instances of the model cannot be solved within a reasonable time with the exact method. Thus, we develop a backward-pass ADP algorithm with a lookup table. The comparison of the optimum and approximate solutions for a small instance shows that the optimality gap of the algorithm is less than 5%. The experiments with the generated data illustrate that the approximate policies result in significantly lower costs than the myopic policies in almost all of the cost levels considered. The analysis of the approximate policies shows that the non-elective patients arriving in earlier or later during a day are more likely to be accepted. Also, it is found that leaving buffer times within the elective schedule does not always result in less cost. Another significant observation is that the shortest-first scheduling of electives results in less cost than the longest-first in most of the cases. Finally, the results indicate that assigning different types of surgeries to a room is more cost-effective than assigning same type of surgeries.

## 5.2 Methodological and Practical Limitations of Thesis

During the research, we encountered with several limitations as summarized below:

• The stem-cell donation centres in the UK and Turkey were quite reluctant to release their data due to the sensitivity of donation operations and donor safety issues. Thus, the input data for the experiments in Chapter 2 are generated with simulation based on the average values publicly available. The impact of the generated data on the solution was then analyzed through the sensitivity analysis. Also, we had to make several assumptions to simplify the search operations and to develop tractable models.

• Real data of the budget and unit-capacity prices for the computational experiments of Chapter 3 were not available. Thus, we identified plausible ranges for these parameters and obtained the results for these ranges. We also assumed that unit-capacity prices are given as inputs to the model while in the reality, these may be subject to negotiations between the central authority and providers, thus, may be decision variables. Including these negotiations into the modelling would require a different approach such as game theory and would make the model highly complex.

• Real data of cost coefficients for the computational experiments of Chapter 4 were not available. These data can be deduced from decision-makers by applying a multi-criteria decision analysis that is out of the scope of this thesis. Instead, we obtained the approximate policies for possible combinations of cost coefficient levels that may be applicable to a range of decision-makers.

• Finally, the solutions obtained in these chapters are not implemented in the real-life. Instead, we used simulation to imitate the reality and investigate the

solution performances. However, it should be noted that the real-life performances of the solutions can be different than the simulated ones. On the other hand, implementation of the findings usually encounters with many practical challenges and was not the main aim of this thesis.

## 5.3 Future Research Directions

As suggested by the previous section, there are several future research directions for this thesis that are summarized below.

• Some of the assumptions made in Chapter 2 to simplify the model can be relaxed in the future considering the rapid development of the powerful solvers. The most significant one of these assumptions is that one patient can only make one search process. Sometimes, a patient can initiate another search process after an initial failed attempt. Thus, the search duration of these patients would consist of several cycles, instead of a single one. Secondly, the model presented in this chapter can be extended by considering the operations in the donor-side of a stem-cell donation centre which would affect the donor database level and the number of donors found for each patient.

• Chapter 3 can be extended by considering the incentives and pricing issues between providers and a central authority. Also, the network can be modelled as a network of queues instead of modelling each provider as a separate single queue. However, note that the analysis of queuing networks is more challenging. Additionally, different patient referral strategies such as patient choice-based can be investigated. This would require to extend the model with the choice issues. Moreover, different types of medical services, such as cardiology and audiology can be included in the model. Finally, we can easily extend the model to analyze different types of payment methods between providers and the central authority.

• To obtain the real cost coefficient levels for the experiments of Chapter 4, a case study can be conducted in a hospital. This would require to use a weight elicitation technique. Secondly, the multiple objectives in the model can be handled with multi-objective optimization methods such as genetic algorithm. The problem can be expanded by considering the decision-making in a week while the decisions in a day affect the future schedules in a rolling horizon fashion. Finally, instead of assigning an accepted non-elective surgery just after the current one is finished, the slot to assign the non-elective can be a decision variable which would make the model even more complex.

As a conclusion, this thesis aims to assist the decision-making for three healthcare management problems under uncertainty. For this purpose, we use various modelling and solution approaches under OR that is proven to be helpful for many decision-making problems under uncertainty (Brandeau et al., 2004). Specifically, we focus on three capacity planning and resource allocation problems in healthcare management. The models developed in the thesis are quite generic and applicable to many cases with slight modifications, if needed. Reflecting the complexity of the problems considered, the models are non-linear and hard to solve. Thus, we consider possible approximate solution methods for these models and compare them as an attempt to identify the best one.

In my opinion, the current healthcare management practices are generally not supported by analytical methods like the research presented in this thesis. Considering the highly uncertain nature of these practices, the lack of analytical support results in deterioration in patients' health and even their death, as well as the waste of resources. Our research has indicated that especially the capacity planning and resource allocation in healthcare settings can be significantly improved by OR methods such as stochastic programming. We believe that the resulting managerial insights can be very beneficial for healthcare man-

agers. However, we are also aware of the special challenges of the implementation of the results in real-life.

# Bibliography

Addis, B., Carello, G., Grosso, A. & Tànfani, E. (2016), 'Operating room scheduling and rescheduling: a rolling horizon approach', *Flexible Services and Manufacturing Journal* **28**(1), 206–232.
**URL:** *http://dx.doi.org/10.1007/s10696-015-9213-7*

Adelman, D. (2004), 'A Price-Directed Approach to Stochastic Inventory/Routing'.

Adelman, D. (2007), 'Dynamic Bid Prices in Revenue Management'.

Ahmed, S., Tawarmalani, M. & Sahinidis, N. V. (2004), 'A finite branch-and-bound algorithm for two-stage stochastic integer programs', *Mathematical Programming* **100**(2), 355–377.

Aksin, O. Z., De Vericourt, F. & Karaesmen, F. (2008), 'Call center outsourcing contract analysis and choice', *Management Science* **54**(2), 354–368.

Alfonso, E., Xie, X., Augusto, V. & Garraud, O. (2013), 'Modelling and simulation of blood collection systems: improvement of human resources allocation for better cost-effectiveness and reduction of candidate donor abandonment', *Vox sanguinis* **104**(3), 225–233.

Aloulou, M. A. & Portmann, M.-C. (2003), An efficient proactive reactive scheduling approach to hedge against shop floor disturbances, *in* 'In Pro-

ceedings of the 1 st Multidisciplinary International Conference on Scheduling:
Theory and Applications, MISTA 2003', Citeseer, p. 1.

Antony Nolan Registry (2016), 'Finding Donors'.
**URL:** *https://www.anthonynolan.org/patients-and-families/stem-cell-or-bone-marrow-transplant/finding-donor-your-stem-cell-or-bone*

April, J., Glover, F., Kelly, J. P. & Laguna, M. (2003), Simulation-based op-
timization: practical introduction to simulation optimization, *in* 'Proceed-
ings of the 35th conference on Winter simulation: driving innovation', Winter
Simulation Conference, pp. 71–78.

Asaduzzaman, M., Chaussalet, T. J. & Robertson, N. J. (2010), 'A loss network
model with overflow for capacity planning of a neonatal unit', *Annals of Op-
erations Research* **178**(1), 67–76.

Avriel, M. & Schaible, S. (1978), 'Second order characterizations of pseudocon-
vex functions', *Mathematical Programming* **14**(1), 170–185.

Bandi, C. & Bertsimas, D. (2012), *Tractable stochastic analysis in high dimen-
sions via robust optimization*, Vol. 134.

Bandi, C., Bertsimas, D. & Youssef, N. (2015), 'Robust queueing theory', *Oper-
ations Research* **63**(3), 676–700.

Barz, C. & Rajaram, K. (2015), 'Elective Patient Admission and Scheduling un-
der Multiple Resource Constraints', *Production and Operations Management*
pp. n/a–n/a.
**URL:** *http://doi.wiley.com/10.1111/poms.12395*

Batun, S., Denton, B. T., Huschka, T. R. & Schaefer, A. J. (2011), 'Operating

room pooling and parallel surgery processing under uncertainty', *INFORMS journal on Computing* **23**(2), 220–237.

Beksac, M. (2014), TURKIYE DE AKRABA DISI DOKU BILGI VE KOR-DON KANI BANKACILIGI GUNCEL DURUM, Antalya, pp. 42–44.

Beliën, J. & Demeulemeester, E. (2007), 'Building cyclic master surgery schedules with leveled resulting bed occupancy', *European Journal of Operational Research* **176**(2), 1185–1204.

Ben Abdelaziz, F. & Masmoudi, M. (2012), 'A multiobjective stochastic program for hospital bed planning', *J Oper Res Soc* **63**(4), 530–538.
**URL:** *http://dx.doi.org/10.1057/jors.2011.39*

Ben-Tal, A. & Nemirovski, A. (1999), 'Robust solutions of uncertain linear programs', *Operations Research Letters* **25**(1), 1–13.

Ben-Tal, A. & Nemirovski, A. (2000), 'Robust solutions of linear programming problems contaminated with uncertain data', *Mathematical Programming* **88**(3), 411–424.

Berman, O. & Krass, D. (2002), '11 Facility Location Problems with Stochastic Demands and Congestion', *Facility location: applications and theory* p. 329.

Bertsekas, D. P. (1999), *Nonlinear programming*, Athena scientific Belmont.

Bertsekas, D. P. & Tsitsiklis, J. N. (1995), Neuro-dynamic programming: an overview, *in* 'Decision and Control, 1995., Proceedings of the 34th IEEE Conference on', Vol. 1, IEEE, pp. 560–564.

Bertsimas, D., Brown, D. B. & Caramanis, C. (2011), 'Theory and applications of robust optimization', *SIAM review* **53**(3), 464–501.

Bezdek, J. C. & Hathaway, R. J. (2002), Some notes on alternating optimization, *in* 'Advances in Soft Computing—AFSS 2002', Springer, pp. 288–300.

Bhat, U. N. (2015), *An introduction to queueing theory: modeling and analysis in applications*, Birkhäuser.

Billaut, J.-C. & Roubellat, F. (1996), 'A new method for workshop real time scheduling', *International Journal of Production Research* **34**(6), 1555–1579.

Birge, J. R. (1985), 'Decomposition and partitioning methods for multistage stochastic linear programs', *Operations research* **33**(5), 989–1007.

Birge, J. R. & Louveaux, F. (2011), *Introduction to stochastic programming*, Springer.

Blake, J. T., Dexter, F. & Donald, J. (2002), 'Operating room managers' use of integer programming for assigning block time to surgical groups: a case study', *Anesthesia & Analgesia* **94**(1), 143–148.

Boffey, B., Galvao, R. & Espejo, L. (2007), 'A review of congestion models in the location of facilities with immobile servers', *European Journal of Operational Research* **178**(3), 643–662.

Bondareva, M. & Seidmann, A. (2012), Peaker Outsourcing for Service Systems with Time-Varying Arrival Rates, *in* 'System Science (HICSS), 2012 45th Hawaii International Conference on', IEEE, pp. 4806–4813.

Borgman, N. J. (2017), Managing urgent care in hospitals, PhD thesis, University of Twente.

Boyan, J. A. & Littman, M. L. (2000), Exact solutions to time-dependent MDPs, *in* 'NIPS', pp. 1026–1032.

Brandeau, M. L., Sainfort, F. & Pierskalla, W. P. (2004), *Operations research and health care: a handbook of methods and applications*, Vol. 70, Springer.

Bretthauer, K. M., Heese, H. S., Pun, H. & Coe, E. (2011), 'Blocking in healthcare operations: A new heuristic and an application', *Production and Operations Management* **20**(3), 375–391.

Burden, R. L. & Faires, J. D. (1993), 'Numerical analysis', *PWS, Boston* .

Campbell, D. & Arnett, J. (2015), 'A&E meltdown forces thousands of operations to be cancelled', *The Guardian* .
**URL:** *http://www.theguardian.com/society/2015/jan/10/a-and-e-doctors-warn-patient-misery-planned-surgery*

Cardoen, B., Demeulemeester, E. & Beliën, J. (2010), 'Operating room planning and scheduling: A literature review', *European Journal of Operational Research* **201**(3), 921–932.

Carey, K., Burgess, J. F. & Young, G. J. (2011), 'Hospital competition and financial performance: the effects of ambulatory surgery centers', *Health Economics* **20**(5), 571–581.
**URL:** *http://dx.doi.org/10.1002/hec.1617*

Castillo, I., Ingolfsson, A. & Sim, T. (2009), 'Social optimal location of facilities with fixed servers, stochastic demand, and congestion', *Production and Operations Management* **18**(6), 721–736.

CCG, H. & District, R. (2014), 'Operational Resilience and Capacity Planning_Harrogate and Rural District CCG'.
**URL:** *http://www.harrogateandruraldistrictccg.nhs.uk/data/uploads/governing-body-papers/7-august-2014/9.3-operational-resilience-and-capacity-plan.pdf*

CCG, M. K. (2017), 'Clinical Commissioning Cycle_ A Summary for GPs_Milton Keynes CCG'.
**URL:** *http://www.miltonkeynesccg.nhs.uk/resources/uploads/files/Commissioning Cycle for GPs.ppt*

Chao, X., Liu, L. & Zheng, S. (2003), 'Resource Allocation in Multisite Service Systems with Intersite Customer Flows', *Management Science* **49**(12), 1739.

Chen, M., Mehrotra, S. & Papp, D. (2015), 'Scenario generation for stochastic optimization problems via the sparse grid method', *Computational Optimization and Applications* **62**(3), 669–692.

(CHPI), C. f. H. & Interest, P. (2015), The contracting NHS-can the NHS handle the outsourcing of clinical services?, Technical report, London.

Cochran, J. K. & Roche, K. T. (2009), 'A multi-class queuing network analysis methodology for improving hospital emergency department performance', *Computers & Operations Research* **36**(5), 1497–1512.

Creemers, S. & Lambrecht, M. (2009), 'An advanced queueing model to analyze appointment-driven service systems', *Computers & operations research* **36**(10), 2773–2785.

Crouzeix, J.-P. & Ferland, J. A. (1982), 'Criteria for quasi-convexity and pseudo-convexity: Relationships and comparisons', *Mathematical Programming* **23**(1), 193–205.
**URL:** *http://dx.doi.org/10.1007/BF01583788*

Dantzig, G. B. (1955), 'Linear programming under uncertainty', *Management Science* **1**(3-4), 197–206.

De Angelis, V., Felici, G. & Impelluso, P. (2003), 'Integrating simulation and optimisation in health care centre management', *European Journal of Operational Research* **150**(1), 101–114.

Deb, K. (2001), *Multi-objective optimization using evolutionary algorithms*, Vol. 2012, John Wiley & Sons Chichester.

Delage, E. & Ye, Y. (2010), 'Distributionally robust optimization under moment uncertainty with application to data-driven problems', *Operations Research* **58**(3), 595–612.

Deloitte (2016), 2016 Global health care outlook, Technical report.
**URL:** *https://www2.deloitte.com/uk/en/pages/life-sciences-and-healthcare/articles/2016-global-health-care-sector-outlook.html*

Department of Health (2014), Examining new options and opportunities for providers of NHS care-The Dalton review, Technical report, London.
**URL:** *https://www.gov.uk/government/publications/dalton-review-options-for-providers-of-nhs-care*

Department of Health (2016), 'Department of Health:Referral to Treatment Statistics'.
**URL:** *https://www.gov.uk/government/publications?keywords=referral+to+treatment& publication_ filter_ option=all&topics%5B%5D=all&departments %5B%5D=all&official_ document_ status=all&world_ locations%5B%5D=all &from_ date=&to_ date=*

Diskapi Yildirim Beyazit Medical School (2015), 'DNA Typing and Trans. Lab'.
**URL:** *http://www.diskapieah.gov.tr/diskapi/index.php?option=com_ content& view=article&id=357:doku-tipleme-ve-transplantasyon-laboratuvar&catid=63:dal-poliklinii&Itemid=305*

Duma, D. & Aringhieri, R. (2015), 'An online optimization approach for the Real Time Management of operating rooms', *Operations Research for Health Care* **7**, 40–51.

Dutech, A. & Scherrer, B. (2013), 'Partially observable Markov decision processes', *Markov Decision Processes in Artificial Intelligence* pp. 185–228.

Düzgün, R., Thiele, A., Cochran, J. J., Cox, L. A., Keskinocak, P., Kharoufeh, J. P. & Smith, J. C. (2010), Dynamic Models for Robust Optimization, *in* 'Wiley Encyclopedia of Operations Research and Management Science', John Wiley & Sons, Inc.
**URL:** *http://dx.doi.org/10.1002/9780470400531.eorms0271*

Earwicker, S. C. & Whynes, D. K. (1998), 'General Practitioners' referral thresholds and choices of referral destination: an experimental study', *Health Economics* **7**(8), 711–722.

Erdelyi, A. & Topaloglu, H. (2010), 'Approximate dynamic programming for dynamic capacity allocation with multiple priority levels', *IIE Transactions* **43**(2), 129–142.

Erdem, E., Qu, X. & Shi, J. (2012), 'Rescheduling of elective patients upon the arrival of emergency patients', *Decision Support Systems* **54**(1), 551–563.

Erdoğan, G., Erkut, E., Ingolfsson, A. & Laporte, G. (2010), 'Scheduling ambulance crews for maximum coverage', *Journal of the Operational Research Society* **61**(4), 543–550.

Erdogan, S. A., Denton, B. T., Cochran, J. J., Cox, L. A., Keskinocak, P., Kharoufeh, J. P. & Smith, J. C. (2011), 'Surgery planning and scheduling'.

Ferrand, Y., Magazine, M. & Rao, U. (2010), Comparing two operating-room-allocation policies for elective and emergency surgeries, *in* 'Simulation Conference (WSC), Proceedings of the 2010 Winter', IEEE, pp. 2364–2374.

Flessa, S. (2000), 'Where efficiency saves lives: A linear programme for the optimal allocation of health care resources in developing countries', *Health Care Management Science* **3**(3), 249–267.

Fomundam, S. & Herrmann, J. W. (2007), 'A survey of queuing theory applications in healthcare'.

Fruchtman, S. (2003), 'Stem cell transplantation', *MOUNT SINAI JOURNAL OF MEDICINE* **70**(3), 166–170.

Fuloria, P. C. & Zenios, S. A. (2001), 'Outcomes-adjusted reimbursement in a health-care delivery system', *Management Science* **47**(6), 735–751.

Gabrel, V., Murat, C. & Thiele, A. (2014), 'Recent advances in robust optimization: An overview', *European Journal of Operational Research* **235**(3), 471–483.

Gallivan, S., Utley, M., Treasure, T. & Valencia, O. (2002), 'Booked inpatient admissions and hospital capacity: mathematical modelling study', *BMJ: British Medical Journal* **324**(7332), 280.

Gerchak, Y., Gupta, D. & Henig, M. (1996), 'Reservation planning for elective surgery under uncertain demand for emergency surgery', *Management Science* **42**(3), 321–334.

Gnanlet, A. & Gilland, W. G. (2009), 'Sequential and simultaneous decision making for optimizing health care resource flexibilities', *Decision Sciences* **40**(2), 295–326.

Gocgun, Y. & Puterman, M. L. (2014), 'Dynamic scheduling with due dates and time windows: an application to chemotherapy patient appointment booking', *Health care management science* **17**(1), 60–76.

Goh, J. & Sim, M. (2010), 'Distributionally robust optimization and its tractable approximations', *Operations Research* **58**(4-part-1), 902–917.

Govind, R., Chatterjee, R. & Mittal, V. (2008), 'Timely access to health care: Customer-focused resource allocation in a hospital network', *International Journal of Research in Marketing* **25**(4), 294–300.

Gratwohl, A., Baldomero, H. & Passweg, J. (2013), 'Hematopoietic stem cell transplantation activity in Europe', *Current Opinion in Hematology* **20**(6), 485–493.

Green, L. (2006), Queueing analysis in healthcare, *in* 'Patient flow: reducing delay in healthcare delivery', Springer, pp. 281–307.

Green, L. V., Kolesar, P. J. & Whitt, W. (2007), 'Coping with Time-Varying Demand When Setting Staffing Requirements for a Service System', *Production and Operations Management* **16**(1), 13–39.

Guerriero, F. & Guido, R. (2011), 'Operational research in the management of the operating theatre: a survey', *Health care management science* **14**(1), 89–114.

Güneş, E. D. & Yaman, H. (2010), 'Health network mergers and hospital re-planning', *Journal of the Operational Research Society* **61**(2), 275–283.

Gunpinar, S. (2013), 'Supply Chain Optimization of Blood Products'.

Gupta, D. & Wang, L. (2008), 'Revenue management for a primary-care clinic in the presence of patient choice', *Operations Research* **56**(3), 576–592.

Gupta, V. & Osogami, T. (2011), 'On Markov–Krein characterization of the mean waiting time in M/G/K and other queueing systems', *Queueing Systems* **68**(3-4), 339–352.

Gurvich, I. & Perry, O. (2012), 'Overflow networks: approximations and implications to call center outsourcing', *Operations Research* **60**(4), 996–1009.

Harper, P. R., Powell, N. H. & Williams, J. E. (2010), 'Modelling the size and skill-mix of hospital nursing teams', *Journal of the Operational Research Society* **61**(5), 768–779.

Haugh, M. B. & Kogan, L. (2007), 'Duality theory and approximate dynamic programming for pricing American options and portfolio optimization', *Handbooks in operations research and management science* **15**, 925–948.

Heydari, M. & Soudi, A. (2016), 'Predictive/Reactive Planning and Scheduling of a Surgical Suite with Emergency Patient Arrival', *Journal of medical systems* **40**(1), 1–9.

(HFMA), H. F. M. A. (2005), Achieving operating room efficiency through process integration, Technical report, Healthcare Financial Management Association.

Hosseini, N. (2012), Managing Elective and Non-elective Case Assignments for an Operating Room Suite, PhD thesis, Clemson University.
**URL:** *http://tigerprints.clemson.edu/cgi/viewcontent.cgi?article=2065&context=all_dissertations*

Howard, D. H., Meltzer, D., Kollman, C., Maiers, M., Logan, B., Gragert, L., Setterholm, M. & Horowitz, M. M. (2008), 'Use of cost-effectiveness analysis to determine inventory size for a national cord blood bank', *Medical Decision Making* **28**(2), 243–253.

Hulshof, P. J. H., Mes, M. R. K., Boucherie, R. J. & Hans, E. W. (2013), Tactical planning in healthcare using approximate dynamic programming, Technical report, University of Twente.

Hurley, C. K., Vina, M. F. & Setterholm, M. (2003), 'Maximizing optimal hematopoietic stem cell donor selection from registries of unrelated adult volunteers', *Tissue Antigens* **61**(6), 415–424.

Johnson, M. M. D. (2008), 'Current trends of outsourcing practice in government and business: causes, case studies and logic', *Journal of Public Procurement* **8**(2), 248.

Kakabadse, N. & Kakabadse, A. (2000), 'Critical review-outsourcing: A paradigm shift', *Journal of Management Development* **19**(8), 670–728.

Kaut, M. & Wallace, S. W. (2007), 'Evaluation of scenario-generation methods for stochastic programming', *Pacific Journal of Optimization* **3**(2), 257–271.

Kendall, D. G. (1953), 'Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded Markov chain', *The Annals of Mathematical Statistics* pp. 338–354.

KIBANK (2005), Istanbul Universitesi Tıp Fakultesi Kemik Iligi Bankası Yıllık Rapor, Technical report, Istanbul Universitesi Tıp Fakultesi Kemik Iligi Bankasi, Istanbul.
**URL:** *http://istanbultip.istanbul.edu.tr/kbank/*

Kimura, T. (1983), 'Diffusion approximation for an M/G/m queue', *Operations Research* **31**(2), 304–321.

Koçaga, Y. L., Armony, M. & Ward, A. R. (2015), 'Staffing Call Centers with

Uncertain Arrival Rates and Cosourcing', *Production and Operations Management* **24**(7), 1101–1117.

Kollman, C., Abella, E., Baitty, R. L., Beatty, P. G., Chakraborty, R., Christiansen, C. L., Hartzman, R. J., Hurley, C. K., Milford, E. & Nyman, J. A. (2004), 'Assessment of optimal size and composition of the US National Registry of hematopoietic stem cell donors', *Transplantation* **78**(1), 89–95.

Lakshmi, C. & Iyer, S. A. (2013), 'Application of queueing theory in health care: A literature review', *Operations Research for Health Care* **2**(1), 25–39.

Lamiri, M., Xie, X., Dolgui, A. & Grimaud, F. (2008), 'A stochastic model for operating room planning with elective and emergency demand for surgery', *European Journal of Operational Research* **185**(3), 1026–1037.

Lee, D. K. K. & Zenios, S. A. (2012), 'An evidence-based incentive system for Medicare's End-Stage Renal Disease program', *Management Science* **58**(6), 1092–1105.

Lee, S. J., Klar, N., Weeks, J. C. & Antin, J. H. (2000), 'Predicting costs of stem-cell transplantation', *Journal of clinical oncology* **18**(1), 64.

Lejeune, M. & Noyan, N. (2010), 'Mathematical programming approaches for generating p-efficient points', *European Journal of Operational Research* **207**(2), 590–600.

Lin, J., Muthuraman, K. & Lawley, M. (2011), 'Optimal and approximate algorithms for sequential clinical scheduling with no-shows', *IIE Transactions on Healthcare Systems Engineering* **1**(1), 20–36.

Liu, X., Cai, X., Zhao, R. & Lan, Y. (2015), 'Mutual referral policy for coordi-

nating health care systems of different scales', *International Journal of Production Research* pp. 1–23.

Lu, M. & Donaldson, C. (2000), 'Performance-Based Contracts and Provider Efficiency', *Disease Management and Health Outcomes* **7**(3), 127–137.

Ma, X.-m. & Zhang, F. (2002), 'A genetic algorithm based stochastic programming model for air quality management', *Journal of Environmental Sciences* **14**(3), 367–374.

Macario, A., Vitez, T. S., Dunn, B. & McDonald, T. (1995), 'Where are the costs in perioperative care?: Analysis of hospital costs and charges for inpatient surgical care', *Anesthesiology* **83**(6), 1138–1144.

Mahar, S., Bretthauer, K. M. & Salzarulo, P. A. (2011), 'Locating specialized service capacity in a multi-hospital network', *European Journal of Operational Research* **212**(3), 596–605.

Marianov, V. & Serra, D. (2002*a*), '4 Location Problems in the Public Sector'.

Marianov, V. & Serra, D. (2002*b*), 'Location–allocation of multiple-server service centers with constrained queues or waiting times', *Annals of Operations Research* **111**(1-4), 35–50.

Mason, J. E. (2012), Markov Decision Processes and Approximate Dynamic Programming Methods for Optimal Treatment Design., PhD thesis, North Carolina State University.

Maxwell, M. S., Restrepo, M., Henderson, S. G. & Topaloglu, H. (2010), 'Approximate dynamic programming for ambulance redeployment', *INFORMS Journal on Computing* **22**(2), 266–281.

May, J. H., Strum, D. P. & Vargas, L. G. (2000), 'Fitting the lognormal distribution to surgical procedure times', *Decision Sciences* **31**(1), 129–148.

Milliyet (2016), 'İlik naklinde TÜRKÖK umudu'.
**URL:** *http://www.milliyet.com.tr/ilik-naklinde-turkok-umudu-gundem-2312524/*

Min, D. & Yih, Y. (2014), 'Managing a patient waiting list with time-dependent priority and adverse events', *RAIRO-Operations Research* **48**(01), 53–74.

Monitor (2013), Local price setting and contracting practices for NHS services without a nationally mandated price: a research papers, Technical report, London.
**URL:** *https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/284182/LocalPricingReport23Sept13.pdf*

Monitor (2016), 'Local Variations'.
**URL:** *https://ldp.monitor-nhsft.gov.uk/Pages/LocalVariations.aspx?&&p_SortBehavior=0&p_FileLeafRef=BMI locl-vartn-temp-eto15.5.15_20150624154804817.pdf&&PageFirstRow=1&&View=%7BE7EB6E4C-CC76-47F5-BBA5-F5DA22B6D5F4%7D*

Mousazadeh, M., Torabi, S. A. & Pishvaee, M. S. (2016), Health Service Network Design Under Epistemic Uncertainty, *in* 'Fuzzy Logic in Its 50th Year', Springer, pp. 257–281.

Müller, C. R., Ehninger, G. & Goldmann, S. F. (2003), 'Gene and haplotype frequencies for the loci HLA-A, HLA-B, and HLA-DR based on over 13,000 German blood donors', *Human immunology* **64**(1), 137–151.

Naboureh, K. & Safari, E. (2016), 'A Stochastic Location-Allocation Model for

Specialized Services in a Multihospital System', *Advances in Operations Research* **2016**.

NHS (2013), 'NHS Policies'.
**URL:** *https://www.england.nhs.uk/2013/12/ccg-fund-allocs/*

NHS (2014), Understanding the new NHS, Technical report, National Health Services, London.
**URL:** *https://www.england.nhs.uk/wp-content/uploads/2014/06/simple-nhs-guide.pdf*

NHS (2016*a*), NHS Direct Access Audiology Waiting Times Data, Technical report.
**URL:** *https://www.england.nhs.uk/statistics/statistical-work-areas/direct-access-audiology/*

NHS (2017), 'NHS waiting time pledge'.
**URL:** *http://www.uclh.nhs.uk/pandv/choosingourservices/18weeks/Pages/Home.aspx*

NHS, O. D. T. (2016*b*), 'Organ and Donor Transplantation- Activity Report'.
**URL:** *http://www.odt.nhs.uk/uk-transplant-registry/annual-activity-report/*

Odejide, O. O., Salas Coronado, D. Y., Watts, C. D., Wright, A. A. & Abel, G. A. (2014), 'End-of-life care for blood cancers: A series of focus groups with hematologic oncologists', *Journal of Oncology Practice* **10**(6), e396–e403.

Ozkarahan, I. (2000), 'Allocation of surgeries to operating rooms by goal programing', *Journal of Medical Systems* **24**(6), 339–378.

Patrick, J., Puterman, M. L. & Queyranne, M. (2008), 'Dynamic multipriority patient scheduling for a diagnostic resource', *Operations Research* **56**(6), 1507–1525.

Pehlivan, C., Augusto, V., Xie, X. & Crenn-Hebert, C. (2012), Multi-period capacity planning for maternity facilities in a perinatal network: A queuing and optimization approach, *in* 'Automation Science and Engineering (CASE), 2012 IEEE International Conference on', IEEE, pp. 137–142.

Powell, W. B. (2007), *Approximate Dynamic Programming: Solving the curses of dimensionality*, Vol. 703, John Wiley & Sons.

Powell, W. B. (2009), 'What you should know about approximate dynamic programming', *Naval Research Logistics* **56**, 239–249.

Punke, H. (2013), 'Outsourcing is exploding in healthcare, will the trend last', *Becker's Hospital Review* .
**URL:** *http://www.beckershospitalreview.com/workforce-labor-management/outsourcing-is-exploding-in-healthcare-will-the-trend-last.html*

Querol, S., Rubinstein, P., Marsh, S. G. E., Goldman, J. & Madrigal, J. A. (2009), 'Cord blood banking:'providing cord blood banking for a nation'', *British journal of haematology* **147**(2), 227–235.

Registry, A. N. (2016), 'Finding Donor Your Stem-cell or Bone'.
**URL:** *https://www.anthonynolan.org/patients-and-families/stem-cell-or-bone-marrow-transplant/finding-donor-your-stem-cell-or-bone*

Riabacke, M., Danielson, M. & Ekenberg, L. (2012), 'State-of-the-art prescriptive criteria weight elicitation', *Advances in Decision Sciences* **2012**.

Roy, B. V., Bertsekas, D., Lee, Y. & Tsitsiklis, J. (1997), 'A neuro-dynamic programming approach to retailer inventory management', *Proceedings of the 36th IEEE Conference on Decision and Control* **4**.

Sahinidis, N. V. (2004), 'Optimization under uncertainty: state-of-the-art and opportunities', *Computers & Chemical Engineering* **28**(6), 971–983.

Salomon, J. A., Mathers, C. D., Murray, C. J. L. & Ferguson, B. (2001), 'Methods for life expectancy and healthy life expectancy uncertainty analysis'.

Salzarulo, P. A., Bretthauer, K. M., Côté, M. J. & Schultz, K. L. (2011), 'The impact of variability and patient information on health care system performance', *Production and Operations Management* **20**(6), 848–859.

Santibáñez, P., Bekiou, G. & Yip, K. (2009), 'Fraser Health uses mathematical programming to plan its inpatient hospital network', *Interfaces* **39**(3), 196–208.

Saunders, R. & Westerink, A. (2014), 'Using clinical outsourcing to derive value-based care', *Journal of the Healthcare Financial Management Association* .
**URL:** *http://www.hfma.org/Content.aspx?id=25786#%23*

Sauré, A., Patrick, J., Tyldesley, S. & Puterman, M. L. (2012), 'Dynamic multi-appointment patient scheduling for radiation therapy', *European Journal of Operational Research* **223**(2), 573–584.

Schmid, V. (2012), 'Solving the dynamic ambulance relocation and dispatching problem using approximate dynamic programming', *European Journal of Operational Research* **219**(3), 611–621.

Schofield, W. N., Rubin, G. L., Piza, M., Lai, Y. Y., Sindhusake, D., Fearnside, M. R. & Klineberg, P. L. (2005), 'Cancellation of operations on the day of intended surgery at a major Australian referral hospital', *Med J Aust* **182**(12), 612–615.

Schrieck, J., Akşin, Z. & Chevalier, P. (2014), 'Peakedness Based Staffing for Call Center Outsourcing', *Production and Operations Management* **23**(3), 504–524.

Schütz, H.-J. & Kolisch, R. (2011), 'Approximate dynamic programming for capacity allocation in the service industry', *European Journal of Operational Research* **218**, 239–250.
**URL:** *http://www.sciencedirect.com/science/article/pii/S0377221711008101*

Shapiro, A. (2013), Sample average approximation, *in* 'Encyclopedia of Operations Research and Management Science', Springer, pp. 1350–1355.

Simao, H. & Powell, W. (2009), 'Approximate dynamic programming for management of high-value spare parts', *Journal of Manufacturing Technology Management* **20**(2), 147–160.

Stancu-Minasian, I. M. (1984), *Stochastic programming with multiple objective functions*, Vol. 13, D Reidel Pub Co.

Stougie, L. & Van Der Vlerk, M. H. (2003), 'Approximation in stochastic integer programming'.

Strum, D. P., May, J. H. & Vargas, L. G. (1998), 'Surgical procedure times are well modeled by the lognormal distribution', *Anesthesia & Analgesia* **86**(2S), 47S–47S.

Strum, D. P., May, J. H. & Vargas, L. G. (2000), 'Modeling the Uncertainty of Surgical Procedure TimesComparison of Log-normal and Normal Models', *The Journal of the American Society of Anesthesiologists* **92**(4), 1160–1167.

Stuart, K. & Kozan, E. (2012), 'Reactive scheduling model for the operating theatre', *Flexible Services and Manufacturing Journal* **24**(4), 400–421.

Stummer, C., Doerner, K., Focke, A. & Heidenberger, K. (2004), 'Determining location and size of medical departments in a hospital network: a multiobjective decision support approach', *Health Care Management Science* **7**(1), 63–71.

Sun, Y. & Li, X. (2013), 'Response surface optimisation of surgery start times in a single operating room using designed simulation experiments', *International Journal of Healthcare Technology and Management* **14**(1), 61–72.

Syam, S. S. & Côté, M. J. (2010), 'A location allocation model for service providers with application to not-for-profit health care organizations', *Omega* **38**(3), 157–166.

Talluri, K. T. & Van Ryzin, G. J. (2006), *The theory and practice of revenue management*, Vol. 68, Springer Science & Business Media.

Testi, A., Tanfani, E. & Torre, G. (2007), 'A three-phase approach for operating theatre schedules', *Health Care Management Science* **10**(2), 163–172.

Tijms, H. C., Van Hoorn, M. H. & Federgruen, A. (1981), 'Approximations for the steady-state probabilities in the M/G/c queue', *Advances in Applied Probability* pp. 186–206.

Topaloglu, H. & Powell, W. B. (2006), 'Dynamic-programming approximations for stochastic time-staged integer multicommodity-flow problems', *INFORMS Journal on Computing* **18**(1), 31–42.

Tsitsiklis, J. N. & Van Roy, B. (2001), 'Regression methods for pricing complex American-style options', *IEEE Transactions on Neural Networks* **12**, 694–703.

UK, N. H. S. (2017), 'NHS Scotland'.
**URL:** *https://www.england.nhs.uk/ourwork/demand-and-capacity/about/*

University of Twente (2017), 'Surgery Scheduling Benchmark Set'.
   **URL:** *https://www.utwente.nl/en/choir/research/benchmarkORscheduling/*

Utley, M., Gallivan, S., Treasure, T. & Valencia, O. (2003), 'Analytical methods for calculating the capacity required to operate an effective booked admissions policy for elective inpatient services', *Health care management science* **6**(2), 97–104.

van Essen, J. T., Hurink, J. L., Hartholt, W. & van den Akker, B. J. (2012), Operating room rescheduling.
   **URL:** *http://eprints.eemcs.utwente.nl/21177/01/wp_368.pdf*

Van Riet, C. & Demeulemeester, E. (2015), 'Trade-offs in operating room planning for electives and emergencies: A review', *Operations Research for Health Care* **7**, 52–69.

Waeber, R., Frazier, P. I. & Henderson, S. G. (2013), 'Bisection search with noisy responses', *SIAM Journal on Control and Optimization* **51**(3), 2261–2279.

Wallace, S. W. & Helgason, T. (1991), 'Structural properties of the progressive hedging algorithm', *Annals of Operations Research* **31**(1), 445–455.

Whitt, W. (1993), 'Approximations for the GI/G/m queue', *Production and Operations Management* **2**(2), 114–161.

Wu, S. D., Byeon, E.-S. & Storer, R. H. (1999), 'A graph-theoretic decomposition of the job shop scheduling problem to achieve scheduling robustness', *Operations Research* **47**(1), 113–124.

Wullink, G., Van Houdenhoven, M., Hans, E. W., van Oostrum, J. M., van der

Lans, M. & Kazemier, G. (2007), 'Closing emergency operating rooms improves efficiency', *Journal of Medical Systems* **31**(6), 543–546.

Yahia, Z., Eltawil, A. B. & Harraz, N. A. (2015), 'The operating room case-mix problem under uncertainty and nurses capacity constraints', *Health care management science* pp. 1–12.

Yan, X., Diaconis, P., Rusmevichientong, P. & Roy, B. V. (2004), Solitaire: Man versus machine, *in* 'Conference on Advances in Neural Information Processing'.

Zhang, Y., Berman, O., Marcotte, P. & Verter, V. (2010), 'A bilevel model for preventive healthcare facility network design with congestion', *IIE Transactions* **42**(12), 865–880.

Zhang, Z., Xie, X. & Geng, N. (2013), 'Dynamic Surgery Assignment of Multiple Operating Rooms With Planned Surgeon Arrival Times', *Automation Science and Engineering* **PP**(99), 1–12.

Zhou, Y.-P. & Ren, Z. J. (2010), 'Service outsourcing', *Wiley encyclopedia of operations research and management science. Hoboken, NJ: Wiley. Available from: http://faculty. washington. edu/yongpin/Service% 20Outsourcing. pdf [Accessed 1 Sepetember 2011]* .

Zonderland, M. E., Boucherie, R. J., Litvak, N. & Vleggeert-Lankamp, C. L. A. M. (2010), 'Planning and scheduling of semi-urgent surgeries', *Health care management science* **13**(3), 256–267.