Innovative Applications of O.R.

# Identifying proactive ICU patient admission, transfer and diversion policies in a public-private hospital network

José Tomás Marquinez [a,*], Antoine Sauré [b], Alejandro Cataldo [c], Juan-Carlos Ferrer [a]

[a] *Departamento de Ingeniería Industrial y de Sistemas, Pontificia Universidad Católica de Chile, Vicuña Mackenna 4860, Santiago, Chile*
[b] *Telfer School of Management, University of Ottawa, 55 Laurier Ave. E., Ottawa, ON K1N 6N5, Canada*
[c] *Institute for Mathematical and Computational Engineering, School of Engineering, Pontificia Universidad Católica de Chile, Vicuña Mackenna 4860, Santiago, Chile*

## ARTICLE INFO

## ABSTRACT

Management of hospital beds is a high-impact issue for two-tier healthcare systems, due principally to their critical importance and high costs. Bed capacity in the public sector is generally insufficient to provide immediate care to all critical patients and thus a significant proportion of public expenditure is assigned to the diversion of patients for treatment in the private sector. We formulate and approximately solve a discounted infinite-horizon Markov Decision Process (MDP) that seeks to identify cost-effective policies for transferring ICU patients between hospitals or diverting them to private clinics. The solution approach employs an affine architecture for approximating the value function of the MDP model and solves the equivalent linear programming model using column generation. The approach can handle a high level of dimensionality, enabling it to consider the arriving patients' many different diagnostic groups and their corresponding lengths of stay. The decisions generated through this approach often differ from the intuitive ones produced in a typical day-by-day decision process, that does not consider the impact of the current day's decisions on the future performance of the system. In particular, the resulting policies will in many cases proactively transfer patients to a different public facility or divert them to a private one even though the hospital they first arrived at had beds available. The performance of the proposed method was evaluated by simulating a case study based on data from a hospital network in Santiago, Chile, producing savings of almost 49% due mostly to reduced usage of private services.

© 2021 Elsevier B.V. All rights reserved.

## 1. Introduction

Managing critical care network capacity is a high-impact issue, due mainly to the heavy cost of providing such services and their critical importance to the beneficiaries. The fundamental measure of the capacity of a hospital network is the number of critical care beds (Green, 2005), which typically is limited. Many patients admitted to a hospital through the casualty/emergency departments are classified as critical and in urgent need of attention, and will thus require admission to an Intensive Care Unit (ICU) as quickly as possible (Chan, Farias, & Escobar, 2017). ICUs serve a broad range of patients who have sustained an unexpected injury or illness, suffer from a pre-existing condition, or are in need of support before and after undergoing complex procedures (Canadian Institute for Health Information, 2016). Typical examples include persons with myocardial infarction or respiratory failure, or who require an abdominal aorta intervention. Such cases require high-complexity in-

patient beds or "critical" beds. Thus, an efficient management of this key resource is essential to ensure prompt access to care for the largest possible number of critical patients.

This critical bed management problem is particularly complex if it must be addressed at the level of a hospital network, by which is meant a set of two or more facilities that work together to coordinate and deliver a spectrum of services to the community (American Hospital Association, 2019). In countries with this type of arrangement such as Chile, the US or the UK —among others— patients are often transferred between hospitals, typically because the facility they first arrived at does not offer the specialty or appropriate level of care they need or is short of beds (Mahar, Bretthauer, & Salzarulo, 2011). These inter-hospital transfers are particularly common in the UK and the US. In other cases such as Hong Kong, patients are transferred between regional hospitals if there are significant waiting time differences between them (Chao, Liu, & Zheng, 2003a). However, if bed management is inefficient, the effective number of admissions at each hospital in the network may be limited. This will result in patients either waiting to be attended to (Shi, Chou, Dai, Ding, & Sim, 2016), with possible

---

* Corresponding author.
*E-mail address:* jtmarquinezv@uc.cl (J.T. Marquinez).

ARTICLE IN PRESS

JID: EOR [m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al. European Journal of Operational Research xxx (xxxx) xxx

negative impacts on their health, or sent to other hospitals with the consequent transfer costs to the system.

Typically, when patients needing a critical bed arrive at a hospital in the public network, they are first stabilized and an assessment is made of their condition. Staff then determine whether or not there is capacity to admit them (Kim, Chan, Olivares, & Escobar, 2015). If the hospital has no capacity, a patient must be transferred to another public facility that does, and if no critical beds are available anywhere in the network, there may be no other choice but to divert the patient to a costly private clinic.

Based on our experience working in this area and feedback received from physicians and ICU managers, we believe that it is possible to make admission and transfer/diversion decisions in two different ways, namely reactively and proactively. The reactive approach is the one just described above and the one most commonly used around the world. The proactive approach, which is the one considered in this paper, involves transferring some patients to other public hospitals and diverting others to private clinics even when the ICU at the hospital of first arrival has critical beds available. We believe that proactive admission and transfer/diversion decisions could significantly reduce associated costs and increase the number of critical patients who receive care in the public sector, without any detriment to patients' health. However, deciding what types of patients to transfer or divert proactively, and where to transfer or divert them, is an extremely challenging task that involves a large number of possible scenarios and alternative courses of action. Thus, there is a clear need for a systematic way of identifying proactive policies to ensure that a hospital network's available capacity is used efficiently at all times as demand requires. Such a method would reduce the number of patients treated outside the public network (i.e., diverted to private facilities) and thereby cut the fiscal burden on the public purse.

In the case of Chile, each facility in the public hospital system has its own administration (both financial and operational). When a critical patient requires ICU care and there is no critical bed available (or soon to be freed up) at the hospital of first arrival, a request is issued to an entity known as the Centralized Bed Management Unit (UGCC by its Spanish initials). This unit manages and shares information on bed availability throughout the public hospital network using a real-time online registry, enabling it to determine where the patient can be received. The destination facility may be either another public hospital or a nearby private clinic. In the latter case, the facility will be chosen from a predefined list of clinics that have contracts with the government to provide care to critical patients the public network cannot handle (Ministerio de Salud de Chile, 2018). The different types of critical care and their corresponding contracts are classified based on diagnosis-related groups (DRGs), and patients are assigned a code indicating the group their diagnosis belongs to. The clinics charge the hospitals accordingly, regardless of the patient's final length of stay (LOS). The cost structure considered in our study is based on these predetermined charges. Also, as is the case in other countries, we assume that the private sector in Chile has more than enough critical beds to absorb the excess demand for ICU care in the public sector.

Of all critical patients arriving at public hospitals in Chile, 95.42% are treated on site while the remainder are referred to other facilities, wether public or private. In practice, however, although there are 1334 critical beds in the public network, more than 70% of the referrals are directed to private clinics (Lara et al., 2016). This results in a major expense for the taxpayer given that the average cost of critical care treatment in the private sector is almost three times that for the public sector (Ministerio de Salud de Chile, 2018). It is also important to note that the Chilean public perceives the level of care they receive at private clinics to be significantly better than that they receive at public hospitals (Aravena

& Inostroza, 2015). Even so, only around 17% of the Chilean population has private health insurance, mainly because of its high cost. Chileans also perceive the quality of care they receive in the public health sector to be the same across all public hospitals, independently of the patient LOS, as medical outcomes are seen to be roughly the same across hospitals.

Under the current system in Chile, a proactive transfer and diversion policy is particularly appropriate for the country's reality. Even though critical patients must by law be taken to the closest public hospital regardless of its location, they cannot be admitted to it unless they live within its coverage area. Hence, a significant number of critical patients are transferred to other public hospitals once they are stabilized. This occurs even when the preference of the patient and his/her family is to stay at the hospital of first arrival rather than to be transferred to some other public facility. Typically, however, no objections are voiced if a patient is diverted to a private clinic due to the superior perception of private facilities mentioned above. Thus, we believe that proactive patient transfers and diversions could be successfully implemented in a public healthcare system like Chile's. Policies facilitating the proactive transfer and diversion of ICU patients should be well received by the Chilean society and would also be financially beneficial for government finances.

Hence, the present article describes a systematic approach for identifying proactive admittance, transfer and diversion policies for stable ICU patients that assist in improving the use of public hospital network capacity and thereby reduce costs. At the heart of the methodology is a Markov Decision Process (MDP) that supports decision-making on patient admissions, transfers, and diversions. It incorporates a solution approach based on an affine approximation of the MDP value function and the use of linear programming (LP) that sidesteps the "curse of dimensionality" which typically arises in real-world instances of the problem. The methodology is able to generate solutions attuned to the differences between a large number of patient categories. The resulting linear programming model is solved using a column generation algorithm for which there is a non-trivial way to find an initial set of feasible columns. Although developed for the specifics of the Chilean case, the approach should be applicable to any multi-hospital network or system that operates along similar lines.

The remainder of this paper is structured as follows. Section 2 reviews the state of the art in the literature on hospital bed management problems and proactive patient transfer policies, and the methodologies proposed for addressing them. Sections 3 and 4 describe the proposed MDP model and solution approach using approximate dynamic programming. In Section 5, the solution of a small example provides the basis for an analysis of the main characteristics of the transfer and diversion policies generated by the proposed approach under different scenarios. Also presented in this section is a comparison of the benefits obtained by the proposed policies using a simulation of the hospital decisions for an instance built with data from Chile's Greater Santiago South Health Service. Finally, Section 6 summarizes our main conclusions and briefly discusses other applications and possible extensions.

## 2. Literature review

Bed management problems have been studied extensively over the last decade but the dynamic allocation of critical beds in a multi-hospital setting involving a large number of patient categories has received limited attention. Most related to our work are the studies by Chao et al. (2003a), Chao, Liu, and Zheng (2003b), Helm, Ahmadbeygi, and Van Oyen (2011), and Hu, Chan, Zubizarreta, and Escobar (2018). Helm et al. (2011) describe an MDP model for reducing variability in hospital bed occupancy

through patient admission control. However, their analysis does not involve a multi-hospital setting. Chao et al. (2003a,b) propose an MDP model that aids in managing patient flows between hospitals in Hong Kong's public-private network based on system congestion. But to ensure the problem is tractable, they consider only two hospitals and two classes of patients who do not necessarily require immediate attention, that is, who are able to wait. Also, patient classes are defined strictly in terms of geographical location. The model's objective is to minimize wait times. Hu et al. (2018) consider a single hospital and determine a trigger point at which patients are proactively transferred from the general ward and the transitional care unit to the ICU depending on patients' risk of deterioration. They conclude that a judicious use of proactive transfers for the most severe patients could reduce mortality rates and LOS without increasing other adverse events.

The methodological challenges of managing hospital system capacity have prompted the development of a range of solution approaches to address the efficient use of insufficient system resources. The most common approach involves the design and redistribution of capacity in order to find an equilibrium between bed supply and demand. Published studies in this vein have set out to redistribute capacity across a network of hospitals (Kao & Tung, 1981), or determine each hospital's optimal capacity (Fowler et al., 2015; Mathews & Long, 2015).

A different line of research resembling the one taken in the present paper, attempts to improve the management of bed capacity and make better use of the resources already available (e.g., Ben Bachouch, Guinet, & Hajri-Gabouj, 2012). Some studies have adopted this approach with a view to capturing the importance of critical bed availability as a way of reducing mortality rates (Lara et al., 2016) or developing policies for intra-hospital transfers between different care units (González, Ferrer, Cataldo, & Rojas, 2018). Numerous other articles have demonstrated the importance of the approach we adopt here, which is the sharing of resources among multiple hospitals (Fowler et al., 2015; Olafson et al., 2015). Litvak, van Rijsbergen, Boucherie, and van Houdenhoven (2008), for example, present a mathematical model to determine how many regional beds are required as a reserve to reduce the number of rejections of incoming ICU patients to a network of hospitals in The Netherlands. However, they do not study policies for intra-network transfers, nor do they consider different LOS for different types of patients. Other studies and proposed approaches in the area of health capacity management have been reviewed by Lakshmi and Iyer (2013) and Bhattacharjee and Ray (2014).

Another issue that has not received much attention in the literature is the relationship between public and private systems, especially as regards patient transfers between them. The problem of dealing with excess demand for beds, or lack of bed capacity, has generally been treated in terms of policies for channelling waiting patients within a single facility rather than transferring some of them to other sites (e.g., Koizumi, Kuno, & Smith, 2005). There are, however, some examples of the transfer approach. One of these is proposed by Chao et al. (2003a), as already mentioned above, and another is suggested by Mahar et al. (2011), who consider inter-hospital transfers of both urgent and non-urgent patients using a model that penalizes unmet demand on site. The focus of the latter, however, is on analysing strategies for increasing and distributing capacity over a network of hospitals so as to minimize costs while maintaining a given patient service level, defined as the proportion of demand satisfied. Their approach does not, therefore, improve the efficiency in the use of a hospital network's existing resources, which is the aim of the present article.

Finally, the approximate dynamic programming approach that will be used here was originally introduced by Schweitzer and Seidmann (1985) and has been successfully applied since then in various studies (Adelman, 2004; 2007; Adelman & Klabjan, 2012; Adelman & Mersereau, 2008; de Farias & Van Roy, 2004; 2006), with many applications in healthcare settings (Bikker, Mes, Sauré, & Boucherie, 2018; Patrick, Puterman, & Queyranne, 2008; Sauré, Begen, & Patrick, 2020; Sauré, Patrick, Tyldesley, & Puterman, 2012; Sauré & Puterman, 2017). The paper by Adelman (2007) in particular takes an approach similar to the one to be presented here, but in addition to its healthcare focus our problem involves a more complex state and action structure that complicates both its formulation and solution approach. Furthermore, whereas Adelman (2007) assumes resources cannot be used more than once, our approach allows resources to be reused over an infinite planning horizon. Using the same approach in a healthcare context, Patrick et al. (2008) solve the problem of allocating the available capacity of a single resource to different types of patients. In an extension of that work, Sauré et al. (2012) pose a problem that most resembles the one to be studied here. The solution they propose features an Approximate Dynamic Programming (ADP) approach that schedules appointments at a cancer treatment centre's radiotherapy units where patients receive treatment over deterministically defined multiple days but with multiple identical treatment units of a single resource. However, none of these papers address the type of dynamic multi-resource allocation problem studied here, nor the stochastic nature of the LOS of the patients.

Summing up, our approach is developed around an approximate MDP that identifies admission, transfer and diversion policies in a combined public-private hospital network, which therefore includes the option of externalizing patients to other facilities. Our methodology also allows patients to be finely differentiated according to the standard DRG classification, which consists of a large number of diagnostic classes or groups (Zapata, 2018). This last aspect makes a direct solution of the problem intractable for real-world situations with traditional methodologies. Thus, to the best of our knowledge, the present article is the first paper on dynamic bed allocation policies to simultaneously consider the following three factors: (i) multiple hospitals; (ii) multiple types of patients; and (iii) collaboration between public and private hospitals through the externalization of care services.

Table 1 compares the characteristics of the articles most closely related to the present study. For each article, the table indicates the type of hospital network (single or multiple facility), whether or not patients need immediate attention, differentiation of patients by classes, whether patients are externalized (transferred or diverted) if capacity is lacking, whether proactive policies are considered, decision level, handling of randomness, objective function, methodology used, and intended application.

## 3. Markov decision process model

We consider a public network consisting of $H$ public hospitals (hereafter simply "hospitals") that provide services to patients classified into $G$ diagnosis-related groups (hereafter simply "groups"), each of which has potentially different arrival and length of stay (LOS) distributions. Hospital $h$ has a number of critical beds $K_h$, which also defines its patient admission capacity at any moment given the obvious fact that one bed can only accommodate one person. There also exists a set of $P$ private hospitals or clinics (hereafter simply "clinics") that together can receive a practically unlimited number of patients from any patient group.

The sequence of decisions involved in a patient's admission or possible move to a different facility ("transfer" if to a hospital, "diversion" if to a clinic) is illustrated schematically in Fig. 1. Decisions are made every so many hours by a decision-making entity at the public network level (in the Chilean case, the aforementioned UGCC). It is precisely this sequential nature of the problem as well as its size that motivate the application of approximate dynamic programming to solve it.
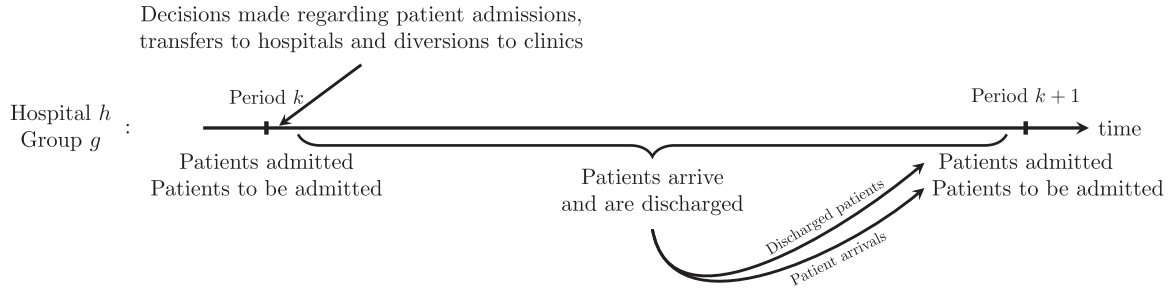
ARTICLE IN PRESS

JID: EOR [m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al. European Journal of Operational Research xxx (xxxx) xxx

**Fig. 1.** Decision-making timeline, assuming decisions are made every few hours.

We assume that a patient occupies ICU capacity only once he/she is definitively admitted to a hospital; that once so admitted, a patient continues to occupy a bed in the same hospital until discharged or removed/dismissed from the ICU (hereafter simply "discharged"); and that setup times for beds that have just been vacated and transfer/diversion times for patients are negligible given that they are much shorter than the length of a decision-making period. It should also be noted that the central administrator can only make decisions regarding patient transfers and diversions, thus leaving medical decisions such as patient discharges to individual hospitals.

### 3.1. Decision epochs

Admittance, transfer and diversion decisions are made at the start of each decision period over an infinite time horizon given that we seek to determine stationary decision policies rather than actions to be taken over a specific fixed time horizon. The time interval for a decision epoch could be a couple of hours or more without affecting the proposed solution approach.

### 3.2. State space

A state $\vec{s} \in S$ of the public hospital network can be defined as $\vec{s} = (\vec{u}, \vec{d}) = \left( u_{111}, \ldots, u_{HGL}; \; d_{11}, \ldots, d_{HG} \right)$, where $u_{hgl} \in \mathbb{Z}_0^+$ is the number of group $g$ patients currently occupying critical beds in hospital $h$ whose LOS is already $l$ time periods, and $d_{hg} \in \mathbb{Z}_0^+$ is the number of group $g$ patients who arrived at hospital $h$ during the previous time period requiring immediate occupancy of a critical bed. States where $l = L$ are defined as the number of patients whose LOS is at least $L$. Note that $u_{hgl}$ is bounded above by the number of critical beds $K_h$ that are available in hospital $h$, while $d_{hg}$ is assumed to be bounded above by a number much smaller than the population the hospital system serves. Thus, $S$ includes a finite number of states.

Given the nature of the problem, not all potential states are feasible. Since the hospitals have a limited number of critical beds, a valid state of the system must satisfy $\sum_{g=1}^{G} \sum_{l=1}^{L} u_{hgl} \leq K_h, \forall h$.

### 3.3. Action sets

At the start of each period, the decision-maker must decide on admissions and patient flows between hospitals and from hospitals to clinics. Each possible action is given by a vector $\vec{a} = (\vec{y}, \vec{z}) = (y_{111}, \ldots, y_{HGP}; z_{111}, \ldots, z_{HHG})$, where $y_{hgp} \in \mathbb{Z}_0^+$ is the number of group $g$ patients diverted from hospital $h$ to clinic $p$, $z_{hig} \in \mathbb{Z}_0^+$ is the number of group $g$ patients transferred from hospital $h$ to hospital $i$ (where $i \neq h$), and $z_{hhg} \in \mathbb{Z}_0^+$ is the number of group $g$ patients admitted to the same hospital $h$ they first arrived at (i.e., direct admissions, not transfers). A feasible action in state $(\vec{u}, \vec{d}) \in S$ must satisfy the condition that all patients requiring a critical bed must be either admitted directly at their hospital of first arrival,

transferred to another hospital, or diverted to a clinic. This implies that

$$d_{hg} = \sum_{i=1}^{H} z_{hig} + \sum_{p=1}^{P} y_{hgp} \qquad \forall h, g. \tag{1}$$

The actions available to the decision-maker must also satisfy constraint (2), which indicates that for each hospital the number of previous admissions not yet discharged plus the number of patients who will be admitted directly and those who are to be transferred from other hospitals must be less than or equal to the hospital's capacity. In formal terms,

$$\sum_{g=1}^{G} \left[ \sum_{l=1}^{L} u_{hgl} + \sum_{i=1}^{H} z_{ihg} \right] \leq K_h \qquad \forall h. \tag{2}$$

Another factor conditioning actions is that for various reasons, certain admissions and/or transfers might not be possible. For example, some hospitals may not have the resources (specialists or equipment) to handle certain diagnosis-related groups. Also, transfers may be prohibited in the case of hospitals located at too great a distance. These restrictions are reflected in the following constraint:

$$z_{hig} = 0 \qquad \forall (h, i, g) \in \mathcal{M}, \tag{3}$$

where $\mathcal{M}$ is the set of index triples for prohibited admissions or transfers.

The set of feasible actions in state $(\vec{u}, \vec{d}) \in S$ is denoted $A_{(\vec{u}, \vec{d})}$.

### 3.4. Transition probabilities

If, in network state $\vec{s} = (\vec{u}, \vec{d}) \in S$ at the start of a given period, action $\vec{a} = (\vec{y}, \vec{z}) \in A_{(\vec{u}, \vec{d})}$ is taken, the transition to the next state of the network will depend on two groups of random variables: $q_{hg}$, the number of group $g$ patients who arrive directly at hospital $h$ in that period, and $r_{hgl}$, the number of group $g$ patients whose LOS is $l$ and leave the system from hospital $h$'s ICU in that same period. The latter variables are defined in this manner because the discharge probabilities for the different patient groups generally depend on their current LOS and follow Binomial($\hat{u}_{hgl}, \rho_{hgl}$) probability distributions, where $\hat{u}_{hg1} = \sum_{i=1}^{H} z_{ihg}$, $\hat{u}_{hgl} = u_{h,g,l-1} \; \forall l : 1 < l < L$, $\hat{u}_{hgL} = u_{h,g,L-1} + u_{hgL}$, and $\rho_{hgl}$ is the probability of a group $g$ patient leaving hospital $h$ after $l$ periods since arrival, which can be obtained from the LOS distribution. For an instance of the model to be realistic, the value of $L$ should be large enough so that for any group or hospital, the probability of a patient currently in an ICU bed leaving the system does not change if he or she stays in the system for $L$ or more days. Thus, the proposed approach can be used to deal with general patient arrival and service time distributions. Even though the LOS for a specific group of patients can vary across hospitals, we assume hospitals provide the same quality of care and therefore have the same patient outcomes. Differences in

**Table 1**
Comparison of articles most closely related to the present study.

| Article | Level of application | Need imm. attention? | Diff. patient types | Externalization | Proactivity | Decision level | Handling of randomness | Objective function | Methodology | Application |
|---|---|---|---|---|---|---|---|---|---|---|
| Ben Bachouch et al. (2012) | H | ✓[a] | ✓ | ✓[a] | — | Operational | - | Total cost - Admission delays and refusals | MILP | Hospital bed use planning |
| Mathews and Long (2015) | H | — | ✓[b] | — | — | Strategic | Simulation | Average wait time | Queueing theory | Hospital bed scheduling |
| Helm et al. (2011) | H | ✓[a] | — | — | ✓ | Operational | Arrival and discharge rates | Cancellation costs | MDP | Admission control |
| Chao et al. (2003b) | MH[f] | — | — | ✓ | ✓ | Operational and Strategic | Admission and discharge rates | Average wait time | DP | Resource allocation and demand allocation |
| Litvak et al. (2008) | MH | ✓ | ✓[d] | ✓ | — | Operational | Queuing models | Acceptance rate | Equivalent Random Models | Capacity allocation |
| Mahar et al. (2011) | MH | ✓[a] | ✓[c] | ✓ | — | Strategic | Mean, St. Dev., Normal Distrib | Total expected cost - capacity increase & non-attention penalty | MINLP | Locating specialized services and quantity |
| The present article (2021) | MH | ✓ | ✓[e] | ✓ | ✓ | Operational | MDP | Total expected cost - Transfers and diversions | ADP | Admission, transfer and diversion policies |

(a) It also considers patients who can wait to receive care. (b) It considers four classes of urgency. (c) Patients are categorized as urgent or non-urgent. (d) It does not consider different length of stay per type of patient. (e) It considers multiple patient classes according to international codes. (f) It considers two hospitals in the model formulation and solution approach. H = 1 hospital; MH = Multi-hospital.

quality of care could, however, be considered in the model by penalizing transfers accordingly.

The transition probabilities that determine the state of the system in the following period, denoted $\vec{s'} = (\vec{u'}, \vec{d'})$, are then given by

$$p(\vec{s'}|\vec{s}, \vec{a}) = \begin{cases} \prod_{h=1}^{H}\prod_{g=1}^{G}\prod_{l=1}^{L} \mathbb{P}(q_{hg})\mathbb{P}(r_{hgl}|\vec{s}, \vec{a}), & \text{if } \vec{s'} \text{ satisfies constraints (4)–(7);} \\ 0, & \text{otherwise.} \end{cases}$$

$$u'_{h,g,1} = d_{hg} + \sum_{\substack{i=1 \\ i \neq h}}^{H} z_{ihg} - \sum_{\substack{i=1 \\ i \neq h}}^{H} z_{hig} - \sum_{p=1}^{P} y_{hgp} - r_{h,g,0} \qquad \forall h, g. \qquad (4)$$

$$u'_{h,g,l} = u_{h,g,l-1} - r_{h,g,l-1} \qquad \forall h, g, l : 2 \leq l \leq L - 1. \qquad (5)$$

$$u'_{h,g,L} = u_{h,g,L-1} - r_{h,g,L-1} + u_{h,g,L} - r_{h,g,L} \qquad \forall h, g. \qquad (6)$$

$$d'_{hg} = q_{hg} \qquad \forall h, g. \qquad (7)$$

Eqs. (4)–(6) above determine the number of patients in each hospital depending on their current LOS, while Eq. (7) defines the number of patients who arrived during the previous period and who have to be distributed among the different facilities.

### 3.5. Immediate costs

Since admission, transfer and diversion decisions generate a cost to the system, the total cost of taking action $\vec{a} = (\vec{y}, \vec{z}) \in A_{(\vec{u}, \vec{d})}$ in state $\vec{s} = (\vec{u}, \vec{d}) \in S$ is given by

$$c(\vec{s}, \vec{a}) = \sum_{h=1}^{H}\sum_{g=1}^{G}\sum_{p=1}^{P} k_{hgp}y_{hgp} + \sum_{h=1}^{H}\sum_{i=1}^{H}\sum_{g=1}^{G} f_{hig}z_{hig} \quad \forall (\vec{s}, \vec{a}) : \vec{a} \in A_{\vec{s}}, \vec{s} \in S,$$

where $f_{hig}$ is the transfer (e.g., transport by ambulance) plus treatment cost for a group $g$ patient sent from hospital $h$ to hospital $i$, and $k_{hgp}$ is the transport plus treatment cost for a group $g$ patient diverted from hospital $h$ to clinic $p$. These costs could also reflect any potential loss in quality of care resulting from being treated at a different hospital.

Note, however, that since it is reasonable to assume that the cost of treating a patient across all hospitals is roughly the same and that all arriving patients must be treated, this cost component is a constant. Indeed, we can assign it a value of zero, that is $f_{hhg} = 0$, so the cost of an inter-hospital transfer $f_{hig} \ \forall i \neq h$ is just the transport cost. In the case of a diversion to a clinic $k_{hgp}$, the corresponding cost can then be considered as the additional amount that measures how much more expensive it is to send a patient for treatment at a private facility.

### 3.6. Optimality equations

The value function of the MDP defines the minimum expected total discounted cost over an infinite horizon and is denoted $v(\vec{s})$ for state $\vec{s}$. To identify the optimal stationary policy we must solve the following equations, known as the Bellman equations (Puterman, 2005):

$$v(\vec{s}) = \min_{\vec{a} \in A_{\vec{s}}} \left\{ c(\vec{s}, \vec{a}) + \gamma \sum_{\vec{s'} \in S} p(\vec{s'}|\vec{s}, \vec{a})v(\vec{s'}) \right\} \qquad \forall \vec{s} \in S, \qquad (8)$$

where $\gamma$ is the periodic discount factor. These equations allow us to determine the value of each state by finding the feasible action that minimizes the immediate and future expected discounted cost, thus capturing the future impact of current decisions.

ARTICLE IN PRESS

JID: EOR [m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al. European Journal of Operational Research xxx (xxxx) xxx

Solving the system of equations in (8) using traditional dynamic programming methods can be impractical due to the curse of dimensionality. For example, the very modestly-sized example we will analyse in Section 5.1 has approximately $10^{61}$ possible states and $10^{41}$ possible actions. For this reason, we propose the alternative solution approach set out in the next section. Details about the computation of the number of states and actions for a generic instance of the problem can be found in Section F of the online supplement.

## 4. Solution approach - approximate dynamic programming

Our solution approach is based on rewriting the optimality equations as the following equivalent linear programming model (Puterman, 2005),

$$\max \quad \sum_{\vec{s} \in S} \alpha(\vec{s}) \nu(\vec{s}) \tag{9}$$

$$\text{s.t.} \quad c(\vec{s}, \vec{a}) + \gamma \sum_{\vec{s'} \in S} p(\vec{s'}|\vec{s}, \vec{a}) \nu(\vec{s'}) \geq \nu(\vec{s}) \qquad \forall \vec{s} \in S, \vec{a} \in A_{\vec{s}}$$

where $\alpha(\vec{s}) > 0$ is the weight of state $\vec{s} \in S$ in the objective function. We assume that $\sum_{\vec{s} \in S} \alpha(\vec{s}) = 1$, so $\alpha(\vec{s})$ can be interpreted as the probability distribution of the initial system state and the objective function as the weighted value.

Note, however, that since (9) has a variable for each state and a constraint for each feasible state-action pair, the model in this form still has a dimensionality problem. To solve it, we resort to an affine approximation of $\nu(\vec{u}, \vec{d})$ in the state variables. As noted earlier, this solution approach was originally introduced by Schweitzer and Seidmann (1985). Here, we use the following affine approximation in the state variables:

$$\nu(\vec{u}, \vec{d}) = \beta + \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{l=1}^{L} U_{hgl} u_{hgl} + \sum_{h=1}^{H} \sum_{g=1}^{G} D_{hg} d_{hg},$$

$$\forall \vec{s} = (\vec{u}, \vec{d}) \in S, \quad \vec{U}, \vec{D} \geq 0, \beta \in \mathbb{R}. \tag{10}$$

The approximation coefficients $U_{hgl}$ can be thought of as the marginal expected discounted cost of having an additional bed occupied by a group $g$ patient who will spend the $l$-th night of his or her stay in hospital $h$, and $D_{hg}$ as the marginal expected discounted cost of having an additional group $g$ patient waiting at hospital $h$ to be admitted, transferred or diverted. $\beta$ is a free constant offset that makes the function in (10) affine.

To estimate these parameters, which define the approximate optimal policy, we replace the value function in (9) with the approximation in (10) and then formulate the following approximate equivalent linear programming model:

$$\max_{\substack{\vec{U}, \vec{D} \geq \vec{0} \\ \beta \in \mathbb{R}}} \quad \beta + \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{l=1}^{L} \mathbb{E}_\alpha[u_{hgl}] U_{hgl} + \sum_{h=1}^{H} \sum_{g=1}^{G} \mathbb{E}_\alpha[d_{hg}] D_{hg} \tag{11}$$

$$\text{s.t.} \quad (1-\gamma)\beta + \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{l=1}^{L} \mu_{hgl}(\vec{s}, \vec{a}) U_{hgl}$$

$$+ \sum_{h=1}^{H} \sum_{g=1}^{G} \delta_{hg}(\vec{s}, \vec{a}) D_{hg} \leq c(\vec{s}, \vec{a}) \quad \forall \vec{s} \in S, \vec{a} \in A_{\vec{s}},$$

where

$$\mathbb{E}_\alpha[u_{hgl}] = \sum_{\vec{s} \in S} \alpha(\vec{s}) u_{hgl}(\vec{s}) \quad \forall h, g, l,$$

$$\mu_{hgl}(\vec{s}, \vec{a}) = u_{hgl}(\vec{s}) - \gamma \sum_{\vec{s'} \in S} p(\vec{s'}|\vec{s}, \vec{a}) u'_{hgl}(\vec{s}, \vec{a}) \quad \forall h, g, l,$$

$$\mathbb{E}_\alpha[d_{hg}] = \sum_{\vec{s} \in S} \alpha(\vec{s}) d_{hg}(\vec{s}) \quad \forall h, g,$$

$$\delta_{hg}(\vec{s}, \vec{a}) = d_{hg}(\vec{s}) - \gamma \sum_{\vec{s'} \in S} p(\vec{s'}|\vec{s}, \vec{a}) d'_{hg}(\vec{s}, \vec{a}) \quad \forall h, g.$$

The $HGL + HG + 1$ variables in (11) pose no real problem but the enormous number of constraints —one for each state-action pair— does. We therefore solve the dual of (11), shown in (12), using column generation, where $\pi(\vec{s}, \vec{a})$ represents the dual variable associated with the state-action pair $(\vec{s}, \vec{a})$. The corresponding pseudo-code is provided below in Algorithm 1.

---

**Algorithm 1** Column generation.

---

1: Obtain an initial set of feasible state-action pairs $\mathcal{P}$ using the proposed phase 1 approach described in Algorithm 2.
2: Define and set the variable $PV = \infty$.
3: **while** $PV > 10^{-5}$ **do**
4:     Solve the restricted version of (2), that considers only the current set of columns $\mathcal{P}$.
5:     Determine $\beta$, $\vec{U}$ and $\vec{D}$ as the shadow prices of the corresponding dual constraints.
6:     Solve (13) using the current shadow prices and store the optimum objective value in $PV$.
7:     **if** $PV > 0$ **then**
8:         Select the state-action pair that generates $PV$.
9:         Compute $\vec{\mu}(\vec{s}, \vec{a}), \vec{\delta}(\vec{s}, \vec{a})$ and $c(\vec{s}, \vec{a})$, and consider it a column.
10:        Update $\mathcal{P}$ by adding the column just obtained.
11:     **end if**
12: **end while**

---

$$\min_{\pi \geq \vec{0}} \quad \sum_{\vec{s} \in S} \sum_{\vec{a} \in A_{\vec{s}}} c(\vec{s}, \vec{a}) \pi(\vec{s}, \vec{a}) \tag{12}$$

$$\text{s.t.} \quad (1-\gamma) \sum_{\vec{s} \in S} \sum_{\vec{a} \in A_{\vec{s}}} \pi(\vec{s}, \vec{a}) = 1$$

$$\sum_{\vec{s} \in S} \sum_{\vec{a} \in A_{\vec{s}}} \mu_{hgl}(\vec{s}, \vec{a}) \pi(\vec{s}, \vec{a}) \geq \mathbb{E}_\alpha[u_{hgl}] \qquad \forall h, g, l$$

$$\sum_{\vec{s} \in S} \sum_{\vec{a} \in A_{\vec{s}}} \delta_{hg}(\vec{s}, \vec{a}) \pi(\vec{s}, \vec{a}) \geq \mathbb{E}_\alpha[d_{hg}] \qquad \forall h, g.$$

To determine the next state-action pair to include in the formulation of (12) in each iteration of the column generation algorithm, the following linear programming model is solved.

$$\max_{\substack{(\vec{s}, \vec{a}) \in \\ S \times A_{\vec{s}}}} \quad (1-\gamma)\beta + \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{l=1}^{L} \mu_{hgl}(\vec{s}, \vec{a}) U_{hgl}$$

$$+ \sum_{h=1}^{H} \sum_{g=1}^{G} \delta_{hg}(\vec{s}, \vec{a}) D_{hg} - c(\vec{s}, \vec{a}) \tag{13}$$

The first step of the column generation algorithm is to obtain an initial set of feasible state-action pairs $\mathcal{P}$ using the phase 1 approach described in Algorithm 2. This approach involves a similar iterative process, but with models (14) and (15) instead, until a feasible solution to (12) is found. In other words, the phase 1 algorithm iterates until all the slack variables $a^\beta$, $a^U_{hgl}$ and $a^D_{hg}$ in (14) are equal to zero.

$$\min_{\pi \geq \vec{0}, a \geq \vec{0}} \quad a^\beta + \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{l=1}^{L} a^U_{hgl} + \sum_{h=1}^{H} \sum_{g=1}^{G} a^D_{hg} \tag{14}$$

$$\text{s.t.} \quad (1-\gamma) \sum_{p \in \mathcal{P}} \pi_p + a^\beta = 1$$

$$\sum_{p \in \mathcal{P}} \mu_{phgl} \pi_p + a^U_{hgl} \geq \mathbb{E}_\alpha[u_{hgl}] \qquad \forall h, g, l$$

$$\sum_{p \in \mathcal{P}} \delta_{phg} \pi_p + a^D_{hg} \geq \mathbb{E}_\alpha[d_{hg}] \qquad \forall h, g.$$

ARTICLE IN PRESS

JID: EOR [m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al. European Journal of Operational Research xxx (xxxx) xxx

**Algorithm 2** Phase 1 approach for obtaining a small set of initial columns.

1: Start with the set of state-action pairs as $\mathcal{P} = \emptyset$.
2: Solve (14) above and store the optimum objective value in $z_{master}$.
3: **while** $z_{master} > 0$ **do**
4:     Determine $\beta$, $\vec{U}$ and $\vec{D}$ as the shadow prices of the corresponding dual constraints.
5:     Solve (15) using the current shadow prices.
6:     Select the state-action pair obtained as the optimal solution.
7:     Compute $\vec{\mu}(\vec{s}, \vec{a})$, $\vec{\delta}(\vec{s}, \vec{a})$ and $c(\vec{s}, \vec{a})$, and consider it a column.
8:     Update $\mathcal{P}$ by adding the column just obtained.
9:     Solve (14), and store the optimum value of that solved problem in $z_{master}$.
10: **end while**

$$\max_{\substack{(\vec{s}, \vec{a}) \in \\ S \times A_{\vec{s}}}} (1 - \gamma)\beta + \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{l=1}^{L} \mu_{hgl}(\vec{s}, \vec{a}) U_{hgl}$$
$$+ \sum_{h=1}^{H} \sum_{g=1}^{G} \delta_{hg}(\vec{s}, \vec{a}) D_{hg} \qquad (15)$$

Note that the structure of both (13) and (15) is that of a minimum-cost network flow problem (see Section E of the online supplement for a proof). Hence, since bed capacities are integer numbers, the linear relaxation of both problems have the integrality property. We take advantage of this structural property to solve the pricing problem more rapidly.

The column generation algorithm determines the optimal values $\vec{U}^*$ and $\vec{D}^*$. If we substitute (10) defined by $\vec{U}^*$ and $\vec{D}^*$ into the right-hand side of the Bellman equation and ignore the constant terms, we obtain the following integer programming model that defines the *approximate optimal policy (AOP)*:

$$d^*(\vec{s}) \in \arg\min_{\vec{a} \in A_{\vec{s}}} \left\{ \sum_{h=1}^{H} \sum_{g=1}^{G} \sum_{p=1}^{P} Y_{hgp} y_{hgp} + \sum_{h=1}^{H} \sum_{i=1}^{H} \sum_{g=1}^{G} Z_{hig} z_{hig} \right\}, \qquad (16)$$

where $Y_{hgp} = k_{hgp} - \gamma U_{hg1}^*, \forall h, g, p,$ and $Z_{hig} = f_{hig} + \gamma U_{ig1}^* - \gamma U_{hg1}^*, \forall h, i, g.$

The coefficients in (16), which define the AOP, can be intuitively interpreted in the following manner. A decision $y_{hgp}$ to divert a patient from hospital $h$ to a clinic involves the transport and treatment cost $k_{hgp}$ of the diversion and the benefit of not subsequently incurring the marginal expected discounted cost $\gamma U_{hg1}^*$ associated with admitting the patient directly to hospital $h$. A decision $z_{hig}$ to transfer a patient from hospital $h$ to hospital $i$ involves the net result of the cost of the transfer $f_{hig}$ plus the subsequent marginal expected discounted cost $\gamma U_{ig1}^*$ associated with having to admit the patient to hospital $i$, and the benefit to hospital $h$ of not having to incur the marginal expected discounted cost $\gamma U_{hg1}^*$ of a direct admission. Thus, on a periodic basis the action to be taken is determined by solving (16) instead of storing an approximate optimal decision for each possible state of the system, which would in practical terms be impossible due to the number of states to be considered. Implementation of (16) is simple and solutions are found in a matter of seconds. The column generation algorithm was written in GAMS 23.8.1 and solved using CPLEX 12.4, while the model defining the AOP was written in Python 3.7 and solved with Gurobi 8.1.0.

## 5. Analysis and decision-making insights

In this section we investigate the main characteristics of the approximate optimal policies derived using the proposed approach. We begin by examining a small instance of the problem, analysing the actions suggested by the AOP based on the values of $Y_{hgp}$ (for diversions to clinics) and $Z_{hig}$ (for transfers to other hospitals) and how the policies change for different cost values and levels of system congestion. Then, to evaluate the performance of the approach for a real-world case, we apply it to a full-scale instance based on a public hospital system in Santiago, Chile.

### 5.1. Approximate optimal policy analysis

In this subsection we illustrate the main properties and the performance of the approximate optimal policy (hereafter "the policy") for a base case. We assume that patient arrivals follow a discrete probability distribution with mean $\lambda_{hg}, \forall h, g$, and that patient discharges follow a Bernoulli process with a success probability $\rho_{hgl}, \forall h, g, l$.

In our base case there are $H = 4$ hospitals with 8, 10, 12 and 15 critical beds, respectively, and $P = 1$ clinic with unlimited critical bed capacity. There are $G = 2$ diagnosis-related groups, and patients whose LOS is $L = 36$ periods or more (at least 3 times the longest average LOS) accumulate in the same state component. The patient arrival rates at the hospitals are such that two of them ($h = 2$ and $h = 3$) are highly congested, one ($h = 1$) has an expected occupancy approximately equal to its capacity, and the remaining one ($h = 4$) has capacity to spare. The complete set of characteristics is set out in Table 2. The cost values are based on actual data from the public hospital system in Santiago, Chile, where diverting a patient for treatment at a clinic costs $k_{hgp} = 8400, \forall h, g, p$, while the cost of transferring a patient by ambulance between hospitals is $f_{hig} = 150, \forall h, i, g$ $(h \neq i)$. Finally, as noted earlier, we assume $f_{hhg} = 0, \forall h, g$.

All the results in this paper were obtained on a computer with an AMD A8-3510MX (1.80 gigahertz) processor and 6 gigabyte of RAM running the 64-bit version of Windows 10 Professional. The approximate optimal policy for the base case was obtained in less than 6 minutes.

#### 5.1.1. Interpretation of the policy for the base case

The values obtained for the $Y_{hgp}$ and $Z_{hig}$ coefficients in the base case are set out in Table 3. Given (16), the policy suggests that a group $g$ patient arriving at hospital $h$ must be treated at the hospital $\hat{i}$ or clinic $\hat{p}$ associated with the lowest value between $Z_{h\hat{i}g}$ and $Y_{hg\hat{p}}$, where $\hat{i} \in \arg\min_i \{Z_{hig} : \text{hospital } i \text{ has beds available}\}$ and $\hat{p} \in \arg\min_p \{Y_{hgp}\}$. Even when the hospital of first arrival has capacity, if this result indicates a different facility, the AOP will indicate that it is preferable for the hospital network to make the corresponding transfer or diversion. If, on the other hand, the minimization result is unique and zero, the default decision is to admit the patient directly.

Now consider the case in Table 3 highlighted in light grey, which reflects a regular and intuitive policy (e.g., the one currently followed in Santiago, Chile). When a group $g = 2$ patient arrives at hospital $h = 3$, the central administrator must evaluate whether to admit him or her directly (zero cost vs. positive costs for all possible transfers or diversions). But if hospital $h = 3$ has no bed available, the central administrator must attempt to transfer the patient to hospital $i = 4$, the index of the facility associated with the second lowest coefficient value ($Z_{342} = 297.00$). If that hospital has no beds, then the central administrator must check hospital $i = 1$ ($Z_{312} = 417.83$). But if that one also has no capacity, it must turn to hospital $i = 2$ ($Z_{322} = 477.09$) and if that one, too, has no bed available, the patient must be diverted to a clinic ($Y_{321} = 4,006.35$).

**Table 2**
Base case arrival rates, average lengths of stay (ALOSs) and implied utilization for AOP analysis.

|  | Group | |  |  | Group | |  |  | Implied util. |
|---|---|---|---|---|---|---|---|---|---|
|  | $g=1$ | $g=2$ |  |  | $g=1$ | $g=2$ |  |  |  |
| $h=1$ | 0.50 | 0.30 |  | $h=1$ | 12.44 | 6.20 |  | $h=1$ | 101.00% |
| $h=2$ | 0.70 | 0.40 |  | $h=2$ | 11.90 | 6.04 |  | $h=2$ | 107.46% |
| $h=3$ | 0.85 | 0.55 |  | $h=3$ | 12.31 | 5.54 |  | $h=3$ | 112.59% |
| $h=4$ | 0.80 | 0.65 |  | $h=4$ | 12.28 | 5.93 |  | $h=4$ | 91.19% |
|  |  |  |  | Average | 12 | 6 |  | Global | 102.26% |

(a) Arrival rate in patients/period ($\lambda_{hg}$).

(b) ALOSs associated with the values of $\rho_{hgl}$, in periods.

(c) Implied utilization in percent[1].

[1] If $\mu_g$ is the average length of stay of group $g$ patients, $\lambda_g$ is the arrival rate of group $g$ patients, and $K$ is the bed capacity, then the implied utilization ($IU$) according to Little's Law is given by $IU = \frac{\sum_g \lambda_g \mu_g}{K}$. Note that the result may be greater than 100%, which would indicate that demand for resources is greater than supply.

**Table 3**
Coefficient values obtained for the base case.

|  |  | Group $g=1$ | | | | Group $g=2$ | | | |  |  | Group | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | Hospital transferred to | | | | Hospital transferred to | | | |  |  | $g=1$ | $g=2$ |
|  |  | $i=1$ | $i=2$ | $i=3$ | $i=4$ | $i=1$ | $i=2$ | $i=3$ | $i=4$ |  |  |  |  |
| Hosp. of arrival | $h=1$ | 0 | 150.00 | 347.82 | 150.00 | 0 | 209.26 | −117.83 | 29.17 |  | $h=1$ | 315.00 | 3,738.52 |
|  | $h=2$ | 150.00 | 0 | 347.82 | 150.00 | 90.74 | 0 | −177.09 | −30.09 |  | $h=2$ | 315.00 | 3,679.26 |
|  | $h=3$ | −47.82 | −47.82 | 0 | −47.82 | 417.83 | 477.09 | 0 | 297.00 |  | $h=3$ | 117.18 | 4,006.35 |
|  | $h=4$ | 150.00 | 150.00 | 347.82 | 0 | 270.83 | 330.09 | 3.00 | 0 |  | $h=4$ | 315.00 | 3,859.35 |

(a) Values of coefficient $Z_{hig}$ for transfers to other hospitals.

(b) Values of coefficient $Y_{hg1}$ for diversions to a clinic.

Note finally that if there were more than one clinic, the patient would be sent to the one with the lowest coefficient.

Next we consider the case highlighted in Table 3 in dark grey, which exemplifies a novel and proactive approach. When a group $g=2$ patient arrives at hospital $h=2$, without checking whether hospital $h=2$ has beds available the central administrator must attempt to transfer him or her to hospital $i=3$ ($Z_{232} = -177.09 < 0 = Z_{222}$). If that hospital has no beds, it must then check hospital $i=4$ ($Z_{242} = -30.09$). If it also has no capacity, the central administrator must try hospital $i=2$, the hospital of first arrival, but if it, too, has no room, hospital $i=1$ ($Z_{212} = 90.74$) is checked before turning as a last resort to the clinic ($Y_{221} = 3,679.26$).

The AOP has a number of properties that make patient transfers and diversions more cost-efficient. The following observations regarding these features are worthy of note:

- Given that the ambulance cost $f_{hig}$ is the same for transfers of patients from any group $g$ between any pair of hospitals $h \neq i$, a centralized decision-maker considering only this cost would have no preferences as to which hospital a patient should be transferred to. Yet the AOP indicates a preference ordering that reflects more than just the ambulance cost, taking into account each hospital's bed capacity and the patient arrival and length of stay distributions for each group. In other words, in deciding patient transfers the policy takes into account not only the ambulance cost but also hospital expected occupancies.
- It may seem counter-intuitive to transfer patients to another hospital when there are still beds available at the hospital of first arrival (e.g., $h=2, g=2$ in Table 3). But such a policy is perfectly logical when one considers that certain groups (e.g., $g=1$) may require longer stays and that patient arrival rates and bed availability vary from one hospital to the next. Thus, making proactive transfers may reduce the number of transfers that will need to be made reactively in the future. A simple ex-

ample of this property is detailed in Section A.1 of the online supplement.

- Table 4 presents the values obtained for the $Y_{hgp}$ and $Z_{hig}$ coefficients when $f_{hig} = 100, \forall h \neq i$. This case illustrates how the AOP can make even more counter-intuitive suggestions. For example, it may recommend diverting patients to clinics proactively, even when this seems to be more expensive than admitting patients directly (e.g., $h=3, g=1$ in Table 4, where $Y_{311} = -27.07 < 0$). Once again, given hospital congestion levels, different arrival rates and longer LOSs for some groups, diverting certain patients to a clinic may work out to be more cost-efficient than admitting them directly to a hospital that is already highly congested. The policy recognizes such a possibility, seeing in the latter action a potential increase in the expected number of *reactive* diversions to be made in the future and therefore a higher expected discounted cost. A simple example of this property is set out in Section A.2 of the online supplement.

One of the most interesting results from the base case (Table 3) is the policy's suggestion that only short-stay patients ($g=2$) from three of the four hospitals ($h=1,2,3$) should be admitted (or rather transferred, in the case of hospitals $h=1,2$) to the hospital with the highest implied utilization ($h=3$). In other words, this result implies that certain hospitals should specialize, as is evident in the fact that it calls for long-stay patients ($g=1$) to be proactively transferred to other hospitals ($Z_{3i1} < 0, \forall i \neq 3$) in exchange for short-stay patients ($Z_{132}, Z_{232}, Z_{332} \leq 0$). Such a specialization should come with not only a better provision of care (Eastaugh, 1992) but also with shorter LOSs (Capkun, Messner, & Rissbacher, 2012). The only hospital that under this policy would not proactively transfer short-stay patients to hospital $i=3$ is hospital $h=4$. However, hospital $i=3$ is the first transfer option if hospital $h=4$ has no beds available. This is so because hospital $h=3$ is the hos-

ARTICLE IN PRESS

JID: EOR
[m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al.
European Journal of Operational Research xxx (xxxx) xxx

**Table 4**
Coefficient values obtained for base case if $f_{hig} = 100, \forall h \neq i$.

| | | Group $g = 1$ | | | | Group $g = 2$ | | | | | | Group | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Hospital transferred to | | | | Hospital transferred to | | | | | | $g = 1$ | $g = 2$ |
| | | $i = 1$ | $i = 2$ | $i = 3$ | $i = 4$ | $i = 1$ | $i = 2$ | $i = 3$ | $i = 4$ | | $h = 1$ | 266.00 | 3,710.7 |
| Hosp. of arrival | $h = 1$ | 0 | 100.00 | 393.07 | 100.00 | 0 | 159.62 | −119.56 | −21.56 | Hospital | $h = 2$ | 266.00 | 3,650.65 |
| | $h = 2$ | 100.00 | 0 | 393.07 | 100.00 | 40.38 | 0 | −179.18 | −81.18 | | $h = 3$ | −27.07 | 3,929.83 |
| | $h = 3$ | −193.07 | −193.07 | 0 | −193.07 | 319.56 | 379.18 | 0 | 198.00 | | $h = 4$ | 266.00 | 3,831.83 |
| | $h = 4$ | 100.00 | 100.00 | 393.07 | 0 | 221.56 | 281.18 | 2.00 | 0 | | | | |

(a) Values of coefficient $Z_{hig}$ for transfers to other hospitals.

(b) Values of coefficient $Y_{hg1}$ for diversions a clinic.

**Table 5**
Comparison of simulated occupancy rates for the base case (AOP and myopic policy).

| | Myopic | AOP |
| --- | --- | --- |
| Occupancy $h = 1$ | 92.44% | 93.47% |
| Occupancy $h = 2$ | 91.83% | 92.77% |
| Occupancy $h = 3$ | 91.81% | 78.57% |
| Occupancy $h = 4$ | 90.22% | 94.14% |
| Global occupancy | 91.40% | 89.56% |

pital that discharges group $g = 2$ patients most rapidly, prompting the policy to have that hospital admit as many patients from that group as possible. In order to ensure sufficient bed availability for that group, the policy transfers the other patient groups occupying more resources (i.e., more bed-days) to other hospitals.

We refer the reader to Section B of the online supplement for additional examples showing that the policy does not only suggest a preference order for transfers and diversions attempts (proactive or not) when there is specialization of some hospitals, but even when hospitals do not specialize and all have the same average length of stay (ALOS) for any given group.

We also conducted an analysis in which all system parameters remained the same as for the base case with the exception of the hospital capacities, which were all set at 10 critical beds. The resulting policy still suggested proactive transfers (hospital $h = 3$ transferring group $g = 1$ patients to any hospital and receiving group $g = 2$ patients from hospital $h = 1$) and diversions (for hospital $h = 3$ and group $g = 1$). Of course, results would differ for other arrival rates and ALOS configurations. We refer the reader to Section C of the online supplement for further details, where it is possible to see that, even when hospitals have the same capacity, it is convenient to make proactive decisions.

*5.1.2. Performance of the policy for the base case*

The global occupancy rates of the hospital system achieved by simulating the AOP and a myopic policy are compared in Table 5. The myopic policy considers only the immediate cost, ignoring any impact today's decisions might have on the system's future performance. In other words, the myopic policy solves (16) but using $Y_{hgp} = k_{hgp}, \forall h, g, p$, and $Z_{hig} = f_{hig}, \forall h, i, g$ instead of the AOP coefficient values. It is also very representative of the UGCC procedure currently followed in Santiago, Chile, which is reactive rather than proactive. As regards the distributions of the random variables, McManus, Long, Cooper, and Litvak (2004) has shown that Poisson arrivals and exponential lengths of stay are reflective of reality. Indeed, since many studies in the literature assume the same distributions (e.g., Chao et al., 2003b; Koizumi et al., 2005; Shi et al., 2016), we use these distributions for our simulations. A total of 100 simulation runs were conducted for each policy, using Python-Gurobi for a three-year period including one transient year. Execution times for the 100 simulations were less than 6 minutes for both the AOP and the myopic policy.

The decline in the average global occupancy rates brought about by the AOP even though hospital admissions (direct admissions plus transfers) increase and fewer patients are diverted to clinics can be appreciated in Tables 5 and 6 (all results are statistically significant with $p$-value $< 0.0001$). By decongesting the public system through proactive transfers and diversions and reducing the global occupancy, the approximate optimal policy results in a more efficient use of available beds. It is worth noting that a relatively low occupancy rate might be understood as an underuse of resources. However, it allows for better management of demand and provides a greater margin for responding to extreme situations that may arise. The AOP maintain an occupancy rate that manages this trade-off well. Furthermore, even though the number of diversions decreases, so does the global occupancy. This is attributable to the decisions made by the policy such as diverting patients in groups with longer ALOS rather than those with shorter ones. It is also noteworthy that with the AOP, the average discounted cost is reduced by 9.58% thanks to the improved admission, transfer and diversion decisions, a result that is also statistically significant.

*5.2. Sensitivity analysis of the policy*

*5.2.1. Cost analysis*

To better characterize the AOP and its performance, in this section we analyse how the approximate optimal actions change with variations in the cost of transfers and diversions. In every case, we assume that the changes affect $k_{hgp}, \forall h, g, p$, or $f_{hig}, \forall h, i, g$, equally while all other parameters remain constant.
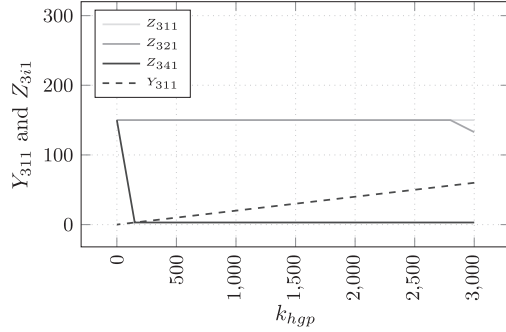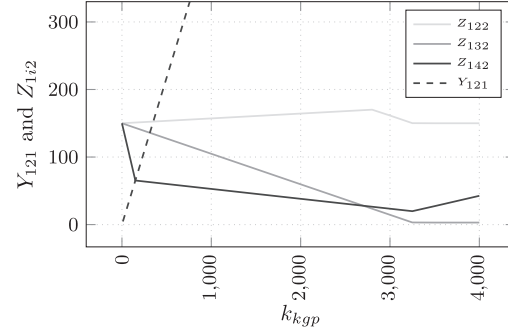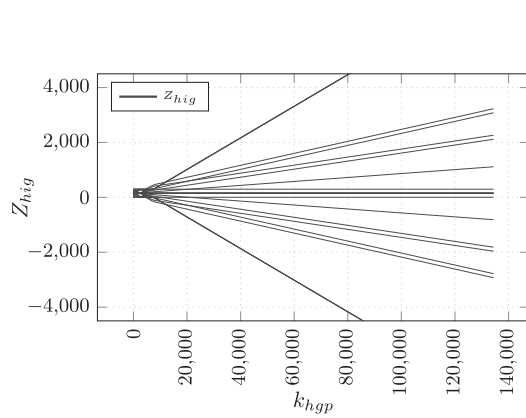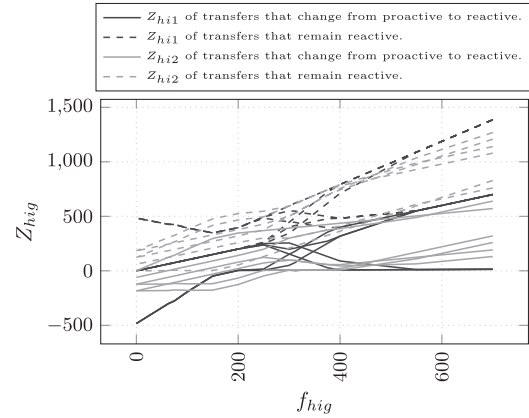
Some of the coefficient values for the base case at diversion cost levels $k_{hgp}$ lower than the 8400 originally assumed are shown in Fig. 2, illustrating what the policy would call for under different diversion cost values. Note in particular that it is not always the case that $Z_{hig} < Y_{hg1}$, which implies that diverting patients to a clinic is not always relegated to the option of last resort. Obviously, if the diversion cost was less than or equal to that for an inter-hospital transfer (i.e., $k_{hgp} \in [0, 150]$), the resulting policy would suggest sending patients to a clinic, but even as it increases above 150, the diversion coefficients are such that the clinic option ($Y_{hgp}$) remains competitive with some inter-hospital transfers ($Z_{hig}$). For example, in Fig. 2a the dashed line representing $Y_{311}$ is below the two continuous lines representing $Z_{311}$ and $Z_{321}$ over the entire range of $k_{hgp}$ in the graph. Thus, over that range the approach will result in policies that suggest it will be preferable to divert patients to a clinic rather than keeping them in some hospitals within the system.

Fig. 2b shows how the preference order for inter-hospital transfers (if there are no beds available at the hospital of first arrival) is not governed solely by the transfer cost $f_{hig}$ but also by the penalty for reactively diverting a patient to a clinic $k_{hgp}$. In this case, when $k_{hgp} = 2,000$ the preference order for group $g = 2$ patients who arrived at hospital $h = 1$ is $4 - 3 - 2$. However, when $k_{hgp} = 3,000$ the preference order changes to $3 - 4 - 2$, as reflected in the figure where $Z_{132}$ falls below $Z_{142}$ before the indicated cost level is

**Table 6**
Comparison of various performance measures for the base case (AOP and myopic policy). The margin of error was computed assuming a 95% confidence level.

|  | Myopic Policy | | AOP | | Δ% |
|---|---|---|---|---|---|
| Direct admissions | 2,218.5 | ± 8.3 | 1,621.2 | ± 7.1 | −24.80% |
| Transfers | 691.8 | ± 6.6 | 1,327.2 | ± 8.2 | +82.85% |
| Diversions | 358.4 | ± 9.9 | 324.8 | ± 8.5 | −12.33% |
| Discounted cost | 221,506 | ± 15,133 | 209,001 | ± 13,552 | −9.58% |
| Average daily cost | 4265 | ± 114 | 4010 | ± 98 | −9.73% |



(a) Analysis for hospital $h = 3$ and group $g = 1$.



(b) Analysis for hospital $h = 1$ and group $g = 2$.

**Fig. 2.** Changes in the values of coefficients $Y_{hgp}$ and $Z_{hig}$ resulting from changes in the clinic transport and treatment costs $k_{hgp}$. Note that $Z_{hhg} = 0, \forall h, g$..



(a) $Z_{hig}$ values for different $k_{hgp}$ values.



(b) $Z_{hig}$ values for different $f_{hig}$ values.

**Fig. 3.** Changes in the values of coefficient $Z_{hig}$ resulting from changes in the transport and treatment cost $k_{hgp} \in [0, 140,000]$ and the transfer cost $f_{hig} \in [0, 700]$. Note that $k_{hgp} = 8,400$ and $f_{hig} = 150$ in the base case. Note that $Z_{hhg} = 0, \forall h, g$.

reached. This change is explained by the fact that the net effect of the occupancy rates, arrival rates and ALOSs is such that it is now preferable to take advantage of the specialization of hospital $h = 3$ in order to divert fewer patients reactively.

Comparing the two figures, we also observe that as the diversion cost increases, diverting patients in groups whose ALOSs are relatively long (group $g = 1$) becomes preferable to diverting those in groups with shorter ALOSs (group $g = 2$). This is so because the former option brings with it larger reductions in system congestion. By thus following the AOP suggestions regarding what patient groups to divert, less resort is had to clinics with the consequent savings in costs.

A sensitivity analysis identifying the changes in the values of coefficient $Z_{hig}$ resulting from changes in the two cost factors, $f_{hig}$ and $k_{hgp}$, is presented graphically in Fig. 3. As can be seen, the policy to be obtained depends globally on both the inter-hospital

transfer cost and the clinic diversion cost since some of the coefficient values change sign. The following effects are worthy of note:

- While increasing the clinic transport and treatment cost $k_{hgp}$, once this cost factor reaches a value of about 10,000, the coefficient values simply follow either an increasing or a decreasing trend (see Fig. 3a). For values of $Z_{hig}$ that are negative, this trend indicates that proactive transfers become even more cost-attractive. Thus, the expected cost of a transfer (immediate cost + expected future discounted cost) to any hospital also depends implicitly on the cost of diversion to a clinic.
- While increasing the inter-hospital transfer cost $f_{hig}$, some negative $Z_{hig}$ coefficients become positive, implying that some suggested transfers change from proactive to reactive. This can be seen in Fig. 3b, where the majority of the $Z_{hig}$ coefficients change sign as the inter-hospital transfer cost increases. The

ARTICLE IN PRESS

JID: EOR [m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al. European Journal of Operational Research xxx (xxxx) xxx

**Table 7**
Coefficient values obtained for the base case if $f_{hig} = 0, \forall h, i, g$.

| | | Group $g = 1$ | | | | Group $g = 2$ | | | | | | | Group | |
| | | Hospital transferred to | | | | Hospital transferred to | | | | | | | $g = 1$ | $g = 2$ |
| | | $i = 1$ | $i = 2$ | $i = 3$ | $i = 4$ | $i = 1$ | $i = 2$ | $i = 3$ | $i = 4$ | | | $h = 1$ | 168.00 | 3,653.76 |
| Hosp. of arrival | $h = 1$ | 0 | 0 | 483.57 | 0 | 0 | 60.33 | −123.03 | −123.03 | | Hospital | $h = 2$ | 168.00 | 3,593.43 |
| | $h = 2$ | 0 | 0 | 483.57 | 0 | −60.33 | 0 | −183.36 | −183.36 | | | $h = 3$ | −315.57 | 3,776.79 |
| | $h = 3$ | −483.57 | −483.57 | 0 | −483.57 | 123.03 | 183.36 | 0 | 0 | | | $h = 4$ | 168.00 | 3,776.79 |
| | $h = 4$ | 0 | 0 | 483.57 | 0 | 123.03 | 183.36 | 0 | 0 | | | | | |

(a) Values of coefficients $Z_{hig}$ for transfers to other hospitals.

(b) Values of coefficients $Y_{hg1}$ for diversions to clinics.

following points suggest how this observation can be interpreted:

- If the cost of transferring patients was zero with our assumed configuration of arrival and discharge rates, all interhospital transfers made with a view to take advantage of unoccupied beds and hospital specialization would be proactive. Table 7 shows the coefficients for this particular scenario, in which the symmetry between $Z_{hig}$ values is very much in evidence. Where necessary, the resulting policy will also dictate that patients be diverted to clinics in order to decongest the hospitals. In our example, several $Z_{hig}$ ($h \neq i$) values are zero, which means that the policy is indifferent as between directly admitting a patient and transferring him or her to another hospital. The policy in this case suggests a common hospital preference order for each patient group independent of the hospital of first arrival. For group $g = 1$, the admission is to either hospitals $h = 1, 2, 4$, or, if no beds are available, the clinic, but never to hospital $h = 3$. The order of attempted admission for group $g = 2$, on the other hand, is first to hospitals $h = 3$ or $h = 4$ indifferently, then hospital $h = 1$, followed by hospital $h = 2$, and finally, the clinic. In addition, to take advantage of the specialization of hospitals $h = 2$ and $h = 3$ in patient groups $g = 1$ and $g = 2$, respectively, hospital $h = 2$ transfers all of its expected short-stay arrivals ($g = 2$) to other hospitals while hospital $h = 3$ does the same with all of its expected long-stay arrivals ($g = 1$).
- If the cost of moving patients between hospitals is non-zero but still low ($f_{hig} \in (0, 50]$), most of the proactive transfers would still be recommended (see Fig. 3b). Since the cost of inter-hospital transfers is low, the resulting policy would still suggest that proactive transfers should be made.
- Upon comparing, for example, the facility preference order for group $g = 2$ arrivals at hospital $h = 2$ of Table 3 ($3 − 4 − 2 − 1$–Clinic) with that of the zero transfer cost case of Table 7 ($3 − 4 − 1 − 2$–Clinic), we see that the resulting policy reorders the hospitals and eliminates the indifference. Thus, as happened when the diversion cost $k_{hgp}$ was changed, a change in $f_{hig}$ may suggest a different facility preference order. This occurs because the preference for hospitals with greater expected available capacity and shorter average stays depends on the magnitude of the "penalty" associated with having to transfer/divert a patient reactively. Hence, if transfer costs increase, the spare bed capacity obtained by carrying out all the proactive transfers and diversions indicated in Table 7 would not be sufficient to reduce costs more than those associated with reactive actions.

Finally, it is important to note that the types of decisions suggested by the approximate optimal policy depend completely on the characteristics of the hospital system such as the number of hospitals, their capacities, arrival rates, ALOSs and specializations.

A similar analysis for $U^*_{hg1}$, the marginal expected discounted cost of admitting patients to different hospitals, is presented in Section D of the online supplement.

### 5.2.2. System congestion analysis

We now examine how the resulting policy changes with variations in system congestion. We consider several scenarios in which hospital utilization levels are modified, either maintaining or increasing the differences in utilization level among hospitals. We believe the results capture well the impact of overall system congestion changes. In this analysis, we assume that every hospital has the same level of specialization for group $g = 1$ patients (i.e., the same ALOS) while the rest of the parameters remain as stated earlier for the base case.

A selection of the values obtained for the $Z_{hig}$ and $Y_{hgp}$ coefficients for different levels of global system congestion resulting from a range of group $g = 1$ ALOSs is set out in Table 8. In particular, we observe the following:

- If the ALOS is less than or equal to 10 periods, implying that the global congestion level (i.e., global implied utilization) is less than 88.17%, then $Z_{hig} = f_{hig} \geq 0$ and $Y_{hgp} = k_{hgp} > 0$. This means that $U^*_{hg1} = 0, \forall h, g$. As a result, no proactive transfers or diversions are made and the resulting policy is the myopic one. This is so because at such a relatively low congestion level, transfers and diversions would only generate a higher discounted cost.
- If the ALOS is increased to 10.5 periods, global system congestion rises to 91.34% and some $Z_{hig}$ coefficients fall close to zero depending on the congestion level at the hospital of first arrival and the patient group. Again, a zero value implies that the policy is indifferent between admitting a patient directly or transferring him/her to another hospital since the associated cost is the same, even when there is a positive transfer cost.
- As the ALOS continues to increase, so does the overall congestion level and some $Z_{hig}$ coefficients turn negative, indicating that it is now advisable to make proactive transfers. More specifically, if the ALOS is for example 11.45, it would be advantageous for hospital $h = 1$ to transfer group $g = 2$ patients to hospital $i = 3$ ($Z_{132} < 0$).
- As the ALOS increases further still, the resulting policy suggests that group $g = 1$ patients arriving at hospital $h = 3$ be transferred to another hospital ($Z_{3i1} < 0, \forall i \neq 3$) even when there is an immediate transfer cost ($f_{3i1} = 150, \forall i \neq 3$). This is the case because hospital $h = 3$ is the one most specialized in group $g = 2$ patients yet group $g = 1$ patients occupy their beds for more periods. It is therefore preferable to free up some capacity for group $g = 2$ patients.

Again, it is important to note that the resulting policy depends heavily not only on the implied utilization level of each hospital but also on the ALOS of each group at each hospital and, consequently, on the opportunities for complementary combinations of the various hospitals' specializations within the network.

ARTICLE IN PRESS

JID: EOR [m5G;March 23, 2021;3:39]

J.T. Marquinez, A. Sauré, A. Cataldo et al. European Journal of Operational Research xxx (xxxx) xxx

**Table 8**
Values of coefficients $Z_{hig}$ and $Y_{hgp}$ for different ALOSs for group $g = 1$ patients. The third row indicates the global implied utilization for each ALOS. Note that $Z_{hhg} = 0, \forall h, g$.

| | | ALOS for group $g = 1$ patients $(\forall h)$ | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | $\leq 10$ | 10.5 | $\cdots$ | 11.4 | 11.45 | 11.5 | $\cdots$ | 12 |
| Impl. Util. | | $\leq 88.17\%$ | 91.34% | $\cdots$ | 97.04% | 97.36% | 97.67% | $\cdots$ | 100.84% |
| $Z_{hig}$ | $Z_{121}$ | 150.00 | 297.00 | | 241.47 | 150.00 | 150.00 | | 150.00 |
| | $Z_{131}$ | 150.00 | 297.00 | | 297.00 | 297.00 | 297.00 | | 337.27 |
| | $Z_{141}$ | 150.00 | 150.00 | | 150.00 | 150.00 | 150.00 | | 150.00 |
| | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| | $Z_{311}$ | 150.00 | 3.00 | | 3.00 | 3.00 | 3.00 | | $-37.27$ |
| | $Z_{321}$ | 150.00 | 150.00 | | 94.47 | 3.00 | 3.00 | | $-37.27$ |
| | $Z_{341}$ | 150.00 | 3.00 | | 3.00 | 3.00 | 3.00 | | $-37.27$ |
| | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| | $Z_{122}$ | 150.00 | 244.54 | | 150.00 | 77.43 | 57.26 | | 49.51 |
| | $Z_{132}$ | 150.00 | 238.13 | | 3.00 | $-69.57$ | $-153.81$ | | $-166.91$ |
| | $Z_{142}$ | 150.00 | 150.00 | | 56.68 | 27.31 | $-6.81$ | | $-19.91$ |
| | $Z_{212}$ | 150.00 | 55.46 | | 150.00 | 222.57 | 242.74 | | 250.49 |
| | $Z_{232}$ | 150.00 | 143.60 | | 3.00 | 3.00 | $-61.06$ | | $-66.42$ |
| | $Z_{242}$ | 150.00 | 55.46 | | 56.68 | 99.87 | 85.94 | | 80.58 |
| | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | |
| $Y_{hgp}$ | $Y_{311}$ | 8,400.00 | 8,253.00 | | 3,983.56 | 2,620.99 | 1,030.91 | | 127.72 |
| | $Y_{321}$ | 8,400.00 | 8,311.87 | | 5,915.45 | 5,159.84 | 4,282.14 | | 3,925.78 |

We refer the reader to Section B.3 of the online supplement to see what the AOP suggests when all hospitals are similar in terms of the ALOS for a specific patient group, and which are the changes when only one hospital increases or decreases specialization for one patient group. The results show that the characteristics of the AOP not only depend on the level of specialization of the hospital network for a given patient group but also on the specialization level for the other patient groups, arrival rates, and number of critical beds.

### 5.3. A real-world example

In what follows we evaluate the potential benefits of the proposed approach by simulating its performance for a practical instance based on a real case and comparing the results with those obtained using a myopic policy. The arrival rates, ALOSs and costs in this instance are those observed for critical beds in the adult intensive care units of the three hospitals making up the Greater Santiago South hospital system in Chile: Barros Luco (BL), El Pino (EP) and San Bernardo (SB). Decisions are made on a daily basis. Information not available about this system was adapted from data on a set of Canadian hospitals. We assumed Poisson arrivals and discrete LOS probability distributions. The diagnosis-related groups used here were the 20 major diagnostic categories (MDC's) of the classification developed at Yale University based on body system and illness aetiology (Yale University School of Public Health, 1981), set forth in Table 9 with the complete data set used for solving and simulating the real case application.

The data for this application are summarized in Table 10. Transport costs were deemed to be proportional to the distance travelled while treatment costs at a private clinic were defined in relation to the Chilean averages derived from the cost weights of each MDC as given in the American Federal Register (1983).

The simulation of each policy consisted of 100 three-year replications with statistics collected after one transient year. The simulation model was implemented using Python-Gurobi and required less than 20 minutes to generate the results for each policy. The algorithm took about 10 hours to determine the AOP.

Performance measures aggregated across the three hospitals were obtained and compared for the AOP and the myopic pol-

icy as shown in Table 11. As can be seen, the AOP admits fewer patients to the hospital of first arrival, transferring and diverting more patients to other facilities. This is a direct consequence of the proactive transfer and diversion decisions suggested by the AOP obtained for this setting. These types of decisions are more than justified by the consequent cost savings, which come to almost 49% with respect to the myopic policy, which can be considered a good representation of the current practice (a value equivalent to some US\$13,981 daily). All of these results were statistically significant ($p$-value $< 0.0001$). The AOP achieves this result by decongesting the system and better handling the variability via proactive diversions of those patients for whom private clinic costs are the lowest (see Fig. 4), thereby allowing the public network to handle patients for whom the diversion costs are the highest.

Furthermore, as apparent from Table 12, the AOP obtains lower occupancy rates at the three hospitals and a lower global occupancy, thus leaving the system better prepared to handle demand variability. The global occupancy rate of 79.85% achieved by the AOP indicates that 24-hour critical beds were unoccupied for an average of only 4.5 hours. The decline in occupation relative to the myopic policy is attributable to the AOP's decisions on which patients to treat in which hospitals and which to divert to clinics.

## 6. Conclusions and future research

This article presents an approach built around the use of approximate stochastic dynamic programming to address the problem of admitting patients to their hospital of first arrival, transferring them between public hospitals or diverting them to private clinics in a public-private network. Typically, such decisions are based exclusively on the immediate cost of transfers and diversions (transport + treatment) without taking into account their impact on the future performance of the system. In other words, transfers and diversions are made only in a reactive manner. By using linear programming to formulate a discounted infinite-horizon Markov decision process and an affine structure for approximating the value function of the equivalent linear programming model, the proposed methodology offers a systematic way of determining admission, transfer and diversion policies for this type of problems. The policies generated by this approach, although perhaps

**Table 9**

Diagnosis-related groups used for the practical application, including corresponding arrival rates in patients/day ($\lambda_{hg}$) and ALOSs in periods associated with the values of $\rho_{hgl}$.

| ID | Name | Arrival rate | | | ALOS | | |
|---|---|---|---|---|---|---|---|
| | | BL | EP | SB | BL | EP | SB |
| g01 | Diseases & Disorders of the Nervous System | 0.20 | 0.10 | 0.19 | 7.73 | 6.97 | 6.26 |
| g02 | Diseases & Disorders of the Eye | – | 0.06 | 0.03 | – | 4.90 | 1.09 |
| g03 | Diseases & Disorders of Ear, Nose, Mouth & Throat | 0.16 | 0.07 | 0.03 | 12.22 | 4.87 | 3.84 |
| g04 | Diseases & Disorders of the Respiratory System | 0.24 | 0.28 | 0.08 | 8.75 | 8.10 | 9.81 |
| g05 | Diseases & Disorders of the Circulatory System | 10.46 | 0.40 | 0.15 | 2.36 | 4.12 | 4.48 |
| g06 | Diseases & Disorders of the Digestive System | 0.20 | 0.15 | 0.07 | 3.04 | 6.23 | 5.23 |
| g07 | Diseases & Disorders of the Hepatobiliary System & Pancreas | 0.21 | 0.09 | 0.04 | 1.44 | 6.70 | 5.48 |
| g08 | Diseases & Disorders of the Musculoskeletal System & Connective Tissue | 0.17 | 0.07 | 0.05 | 1.10 | 5.76 | 4.72 |
| g09 | Diseases & Disorders of the Skin, Subcutaneous Tissue & Breast | 0.13 | 0.06 | 0.04 | 3.36 | 5.77 | 5.26 |
| g10 | Diseases & Disorders of the Endocrine System, Nutrition & Metabolism | 0.21 | 0.08 | 0.04 | 1.65 | 4.16 | 3.62 |
| g11 | Diseases & Disorders of the Kidney, Urinary Tract & Male Reproductive System | 0.20 | 0.10 | 0.04 | 4.34 | 5.22 | 4.56 |
| g12 | Diseases & Disorders of the Female Reproductive System | – | 0.07 | 0.03 | – | 3.45 | 8.51 |
| g13 | Pregnancy & Childbirth | 0.18 | 0.07 | 0.04 | 2.91 | 2.71 | 2.34 |
| g14 | Diseases & Disorders of the Blood & Lymphatic System | 0.20 | 0.10 | 0.04 | 2.08 | 7.67 | 7.20 |
| g15 | Multisystemic or Unspecified Site Infections | 0.18 | 0.21 | 0.07 | 4.63 | 7.82 | 7.48 |
| g16 | Mental Diseases & Disorders | 0.22 | 0.06 | 0.04 | 1.17 | 5.12 | 4.77 |
| g17 | Significant Trauma, Injury, Poisoning & Toxic Effects of Drugs | 0.28 | 0.14 | 0.16 | 3.04 | 5.37 | 6.49 |
| g18 | Burns | – | 0.07 | 0.03 | – | 3.16 | 4.91 |
| g19 | Miscellaneous CMG & Ungroupable Data | – | – | 0.03 | – | – | 0.74 |
| g20 | Other Reasons for Hospitalization | 0.17 | 0.08 | 0.04 | 1.11 | 5.20 | 3.78 |

**Table 10**

Summary of the data used in the real-world application.

| | | Arrival rate (pat/day) | ALOS (days) | Capacity (ICU beds) | Implied utilization |
|---|---|---|---|---|---|
| Hospital | BL | 13.41 | 3.81 | 31 | 116.41% |
| | EP | 2.26 | 5.44 | 12 | 107.59% |
| | SB | 1.26 | 5.03 | 6 | 115.60% |
| | Global | 16.92 | 3.84 | 49 | 114.15% |

**Table 11**

Comparison of the simulation results for the real-world application (AOP and myopic policy). The margin of error was computed for a 95% confidence level.

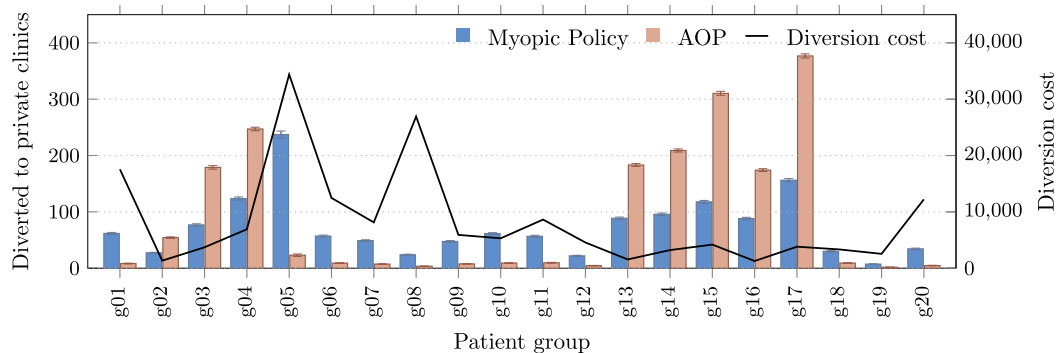| | Myopic Policy | | | AOP | | | Difference (%) |
|---|---|---|---|---|---|---|---|
| Direct admissions | 7,689.0 | ± | 16.5 | 7,209.9 | ± | 13.4 | −6.23% |
| Transfers | 1,084.7 | ± | 7.9 | 1,282.6 | ± | 8.5 | +18.24% |
| Diversions | 1,462.2 | ± | 19.3 | 1,831.2 | ± | 10.41 | +25,24% |
| Discounted cost | 1,093,394 | ± | 59,346 | 553,732 | ± | 23,045 | −49.36% |
| Average daily cost | 21,555 | ± | 407 | 11,069 | ± | 155 | −48.65% |



**Fig. 4.** Comparison of the mean number of diversions by patient group and policy, with diversion costs per group. The margin of error was computed for a 95% confidence level.

sub-optimal, can bring about a significant reduction in critical bed management costs. They allow the decision-maker to identify which types of patients should be transferred or diverted, and to which facility. Furthermore, the proposed approach is able to handle large instances and thus tailor its suggested policies to a large number of patient categories, thereby overcoming the "curse of dimensionality" that typically arises in real-world applications. At a practical level, these capabilities should constitute a major advance in the definition of bed allocation guidelines for hospital networks.

The results obtained with our approach show how better use can be made of critical beds across a hospital network through judicious application of proactive transfers between public hospitals or diversions to private clinics based on facility congestion levels and specialization of individual hospitals. Both for a base case and

**Table 12**
Simulated occupancy rates obtained by the AOP and the myopic policy for the practical application.

|                  | Myopic  | AOP    |
|------------------|---------|--------|
| BL Occupancy     | 95.76%  | 84.98% |
| EP Occupancy     | 95.42%  | 69.03% |
| SB Occupancy     | 93.83%  | 74.96% |
| Global Occupancy | 95.44%  | 79.85% |

a real-world application, a major reduction in the average total discounted cost was attained, producing for the real case of a Santiago, Chile hospital system simulated savings of around 49% in the daily expense level. Under the policies determined by the proposed method, demand is anticipated so that the system can be brought to a desired state rather than leaving it at the mercy of patient arrivals. Thus, it is shown that the best policy is not necessarily to wait until a hospital is full to capacity before transferring or diverting arriving patients to other facilities. On the contrary, there will be times when it is advantageous to make such transfers or diversions proactively even though the hospital a patient first arrived at still has beds available. Furthermore, we have shown that, depending on the characteristics of the problem, the approach can come up with intuitive and not-so-intuitive policies. These results could also be beneficial in providing hospital managers with strategic insights as regards which hospitals should specialize in which types of treatments, which should lead a more efficient allocation of limited resources.

Such fluid and proactive patient transfers and diversions, coordinated across a multi-facility hospital system, may improve efficiency and lower global costs while also achieving better patient throughput and management of variability in critical bed demand. However, in order to implement this approach and achieve similar benefits in other countries, the public hospitals in the system must share information in real time and either coordinate their transfers and diversions with each other or turn these decisions over to a centralized decision-maker like the UGCC in the Chilean system. The importance of coordination among hospitals in order to achieve a higher number of attended patients within the network is also highlighted by Litvak et al. (2008).

The proposed solution approach can be applied to any hospital system where patients can be transferred or diverted between facilities in order to make better use of bed capacity and each facility's specialization. It can also be applied to other services such as 24-hour product home-delivery chains where dispatch to a customer can be made from any branch that is part of a centralized system. In such a case, the resources are the delivery personnel who can be dispatched with an order from a branch other than the one where the order was received. In the public sector, another possible application is a prison system in which the assignment of inmates to correctional facilities considers proactively (i.e., not only reactively) each one's capacity as well as the prisoner arrival rates and sentence lengths.

A number of extensions to the proposed approach are being considered. For example, one could adapt the methodology to consider patient transfers from and to different care units (intensive, intermediate, etc.) within the same or at some other facility. The methodology could also be used to determine which hospitals require additional beds, or take into account capacity limitations at private clinics. Also, variable or LOS-dependant treatment costs at both public hospitals and private clinics could be incorporated. In addition, changes could be made to the column generation procedure to reduce solution times.

In conclusion, for centralized decision-makers in particular, the approach presented here can be very useful in defining transfer and treatment policies that will result in significant savings on health service expenditures while also boosting the number of patients that can be handled by existing hospital capacity. The approach would be especially beneficial in cases like Chile, Britain or the US where the capacity of the public hospital system is such that it must frequently call upon the resources of more costly private health facilities.

## Acknowledgments

## Supplementary material

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.ejor.2021.02.045.

## References

Adelman, D. (2004). A price-directed approach to stochastic inventory/routing. *Operations Research, 52*(4), 499–514. https://doi.org/10.2307/30036604.

Adelman, D. (2007). Dynamic bid prices in revenue management. *Operations Research, 55*(4), 647–661. https://doi.org/10.1287/opre.1060.0368.

Adelman, D., & Klabjan, D. (2012). Computing near-optimal policies in generalized joint replenishment. *INFORMS Journal on Computing, 24*(1), 148–164. https://doi.org/10.1287/ijoc.1100.0433.

Adelman, D., & Mersereau, A. J. (2008). Relaxations of weakly coupled stochastic dynamic programs. *Operations Research, 56*(3), 712–727. https://doi.org/10.1287/opre.1070.0445.

American Federal Register (1983). Table 5: List of diagnosis-related groups, relative weighting factors, mean length of stay, and length of stay outlier cutoff points used in the prospective payment system.

American Hospital Association (2019). Fast facts on US hospitals. http://www.aha.org/research/rc/stat-studies/fast-facts.shtml.

Aravena, P., & Inostroza, M. (2015). Salud Pública o Privada? Los factores más importantes al evaluar el sistema de salud en Chile. *Revista Médica de Chile, 143*, 244–24451. https://doi.org/10.4067/S0034-98872015000200012.

Ben Bachouch, R., Guinet, A., & Hajri-Gabouj, S. (2012). An integer linear model for hospital bed planning. *International Journal of Production Economics, 140*(2), 833–843. https://doi.org/10.1016/j.ijpe.2012.07.023.

Bhattacharjee, P., & Ray, P. (2014). Patient flow modelling and performance analysis of healthcare delivery processes in hospitals: A review and reflections. *Computers and Industrial Engineering, 78*, 299–312. https://doi.org/10.1016/j.cie.2014.04.016.

Bikker, I. A., Mes, M. R., Sauré, A., & Boucherie, R. J. (2018). Online capacity planning for rehabilitation treatments: An approximate dynamic programming approach. *Probability in the Engineering and Informational Sciences*, 1–25. https://doi.org/10.1017/s0269964818000402.

Canadian Institute for Health Information (2016). Care in Canadian ICUs. Technical Report, August. Ottawa, Ontario.

Capkun, V., Messner, M., & Rissbacher, C. (2012). Service specialization and operational performance in hospitals. *International Journal of Operations and Production Management, 32*(4), 468–495. https://doi.org/10.1108/01443571211223103.

Chan, C. W., Farias, V. F., & Escobar, G. J. (2017). The impact of delays on service times in the intensive care unit. *Management Science, 63*(July 2017), 2049–2072.

Chao, X., Liu, L., & Zheng, S. (2003a). Resource allocation in multisite service systems with intersite customer flows. *Management Science, 49*(12), 1739–1752. https://doi.org/10.1287/mnsc.49.12.1739.25110.

Chao, X., Liu, L., & Zheng, S. (2003b). Stochastic network models and optimization of a hospital system. In *Stochastic modeling and optimization* (pp. 363–393). New York, NY: Springer. Chapter 12

Eastaugh, S. R. (1992). Hospital specialization and cost efficiency: Benefits of trimming product lines. *Hospital and Health Services Administration, 37*(2), 223–235.

de Farias, D., & Van Roy, B. (2004). On constraint sampling in the linear programming approach to approximate dynamic programming. *Mathematics of Operations Research, 29*(3), 462–478. https://doi.org/10.1287/moor.1040.0094.

de Farias, D., & Van Roy, B. (2006). A cost-shaping linear program for average-cost approximate dynamic programming with performance guarantees. *Mathematics of Operations Research, 31*(3), 597–620. https://doi.org/10.1287/moor.1060.0208.

Fowler, R. A., Abdelmalik, P., Wood, G., Foster, D., Gibney, N., Bandrauk, N., & Jouvet, P. (2015). Critical care capacity in Canada: Results of a national cross-sectional study. *Critical Care, 19*(1), 133. https://doi.org/10.1186/s13054-015-0852-6.

González, J., Ferrer, J.-C., Cataldo, A., & Rojas, L. (2018). A proactive transfer policy for critical patient flow management. *Health Care Management Science*, 1–17. https://doi.org/10.1007/s10729-018-9437-7.

Green, L. V. (2005). Capacity planning and management in hospitals. In M. L. Brandeau, F. Sainfort, & W. P. Pierskalla (Eds.), *Operations research and health care: A handbook of methods and applications* (pp. 15–41). Boston, MA: Springer US. https://doi.org/10.1007/1-4020-8066-2_2, 10.1007/1-4020-8066-2_2, 10.1007/1-4020-8066-2_2.

Helm, J. E., Ahmadbeygi, S., & Van Oyen, M. P. (2011). Design and analysis of hospital admission control for operational effectiveness. *Production and Operations Management, 20*(3), 359–374. https://doi.org/10.1111/j.1937-5956.2011.01231.x.

Hu, W., Chan, C., Zubizarreta, J., & Escobar, G. (2018). An examination of early transfers to the ICU based on a physiologic risk score. *Manufacturing & Service Operations Management, 20*(3), 531–549. https://doi.org/10.1287/msom.2017.0658.

Kao, E. P. C., & Tung, G. G. (1981). Bed allocation in a public health care delivery system. *Management Science, 27*(5), 507–520. https://doi.org/10.1287/mnsc.27.5.507.

Kim, S.-H., Chan, C., Olivares, M., & Escobar, G. (2015). ICU admission control: An empirical study of capacity allocation and its implication on patient outcomes. *Management Science, 61*, 19–38. https://doi.org/10.2139/ssrn.2062518.

Koizumi, N., Kuno, E., & Smith, T. E. (2005). Modeling patients flows using a queueing network with blocking. *Health Care Management Science, 8*, 49–60. https://doi.org/10.1007/s10729-005-5216-3.

Lakshmi, C., & Iyer, S. A. (2013). Application of queueing theory in health care: A literature review. *Operations Research for Health Care, 2*(1-2), 25–39. https://doi.org/10.1016/j.orhc.2013.03.002.

Lara, B. A., Cataldo, A., Castro, R., Aguilera, P. R., Ruiz, C., & Andresen, M. (2016). Medicina de urgencia y unidades de cuidados intensivos: Una alianza necesaria en busca de la mejoría de la atención de pacientes críticos. *Revista Médica de Chile, 144*(7), 917–924.

Litvak, N., van Rijsbergen, M., Boucherie, R. J., & van Houdenhoven, M. (2008). Managing the overflow of intensive care patients. *European Journal of Operational Research, 185*(3), 998–1010. https://doi.org/10.1016/j.ejor.2006.08.021.

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. https://doi.org/10.1016/j.ejor.2011.03.008.

Mathews, K. S., & Long, E. F. (2015). A conceptual framework for improving critical care patient flow and bed use. *Annals of the American Thoracic Society, 12*(6), 886–894. https://doi.org/10.1513/AnnalsATS.201409-419OC.

McManus, M. L., Long, M. C., Cooper, A., & Litvak, E. (2004). Queuing theory accurately models the need for critical care resources. *Anesthesiology, 100*(5), 1271–1276. https://doi.org/10.1097/00000542-200405000-00032.

Ministerio de Salud de Chile (2018). Unidad de Gestión Centralizada de Camas, UGCC. http://www.minsal.cl/wp-content/uploads/2018/03/Informe-UGCC-2014-2018.pdf.

Olafson, K., Ramsey, C., Yogendran, M., Fransoo, R., Chrusch, C., Forget, E., & Garland, A. (2015). Surge capacity: Analysis of census fluctuations to estimate the number of intensive care unit beds needed. *Health Services Research, 50*(1), 237–252. https://doi.org/10.1111/1475-6773.12209.

Patrick, J., Puterman, M. L., & Queyranne, M. (2008). Dynamic multipriority patient scheduling for a diagnostic resource. *Operations Research, 56*(6), 1507–1525. https://doi.org/10.1287/opre.1080.0590.

Puterman, M. L. (2005). *Markov decision processes: Discrete stochastic dynamic programming*. USA: John Wiley & Sons, Inc.. https://doi.org/10.1080/00401706.1995.10484354.

Sauré, A., Begen, M. A., & Patrick, J. (2020). Dynamic multi-priority, multi-class patient scheduling with stochastic service times. *European Journal of Operational Research, 280*, 254–265. https://doi.org/10.1016/j.ejor.2019.06.040.

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. https://doi.org/10.1016/j.ejor.2012.06.046.

Sauré, A., & Puterman, M. L. (2017). Advance patient appointment scheduling. In R. J. Boucherie, & N. M. van Dijk (Eds.), *Markov decision processes in practice* (pp. 245–268). Cham: Springer International Publishing. https://doi.org/10.1007/978-3-319-47766-4_8, 10.1007/978-3-319-47766-4.

Schweitzer, P. J., & Seidmann, A. (1985). Generalized polynomial approximations in Markovian decision processes. *Journal of Mathematical Analysis and Applications, 110*(2), 568–582. https://doi.org/10.1016/0022-247X(85)90317-8.

Shi, P., Chou, M. C., Dai, J. G., Ding, D., & Sim, J. (2016). Models and insights for hospital inpatient operations: Time-dependent ED boarding time. *Management Science, 62*(January 2016), 1–28.

Yale University School of Public Health (1981). The New ICD-9-CM Diagnosis Related Groups Classification Scheme: User Manual.

Zapata, M. (2018). Importancia del Sistema GRD para Alcanzar la Eficiencia Hospitalaria. *Revista Médica Clínica Las Condes, 29*(3), 347–352. https://doi.org/10.1016/j.rmclc.2018.04.010.