

A Network-Based Formulation for Scheduling Clinical Rotations

Andre A. Cire*

Department of Management, University of Toronto Scarborough and Rotman School of Management, Toronto, Ontario M1C 1A4, Canada, andre.cire@rotman.utoronto.ca

Adam Diamant

Schulich School of Business, York University, Toronto, Ontario M3J 1P3, Canada, adiamant@schulich.yorku.ca

Tallys Yunes

Miami Business School, Coral Gables, Florida 33124-8237, USA, tallys@miami.edu

Alejandro Carrasco

School of Medicine, American University of the Caribbean, Coral Gables, Florida 33134, USA, acarrasco@aucmed.edu

We investigate the scheduling practices of a medical school that must assign a cohort of students to a series of clinical rotations, while respecting both operational and quality-of-service requirements. Students become available to start clerkship progressively throughout the year and can complete rotations at hospitals in different geographic regions. Each hospital may offer a subset of the clinical rotations, with different start dates, capacities, and cost rates. We propose a novel network-flow model based on decision diagrams, a graphical structure that compresses the state space of a dynamic program, to model feasible schedules. We demonstrate that our network model has several interesting structural features, is computationally superior as compared to a classical mixed-integer linear program, and can be used to generate useful insights that can aid in managerial decision-making. Using a dataset collected from the American University of the Caribbean, we perform a counterfactual analysis which shows that had our scheduling approach been implemented, a cost reduction of approximately 19% on average could have been achieved. To understand how assignment decisions can affect future costs, we develop a discrete-event simulation of the licensing examination and clerkship scheduling process. We then compare our exact scheduling approach with current practice and achieve an average cost reduction of 25%. We also show that this cost reduction is robust with respect to estimation and forecast uncertainty, specifically, the licensing exam failure rate and the future cohort size.

Key words: rotation scheduling; medical training; network-based formulation; optimization; simulation

History: Received: April 2018; Accepted: October 2018 by Sergei Savin, after 1 revision.

1. Introduction

The assignment of personnel to shifts is one of the most fundamental problems in scheduling (Ernst et al. 2004). Of particular interest is the field of health care, where scarce resources must be appropriately utilized and budgets carefully balanced. We consider a context where medical students are assigned to a series of rotations of different lengths at hospitals in various geographic locations. This in-hospital component of medical training (i.e., clerkship) provides students with an opportunity to develop clinical skills and ensure they learn how to deliver patient-centered, evidence-based care. It also exposes students to a range of specialization options, such as general surgery, internal medicine, psychiatry, pediatrics, and obstetrics & gynecology. Finally, this type of

experiential education is mandatory; only after the successful completion of these rotations are students eligible to graduate and participate in the well-known hospital residency matching programs. In 2016, for example, 19,254 students graduated from 147 medical schools, with the size of a cohort ranging from 50 to 400 students for each school (Association of American Medical Colleges 2016).

Because clerkship is a required component of medical education, the scheduling of students to core rotations is the responsibility of the medical school. Nonetheless, several factors make this task complex. First, a student can begin their clerkship only after passing a licensing examination, which takes place continuously throughout the year, is scheduled months in advance, and is administered on specialized computers at test centers worldwide. Thus, there

is a constant (dynamic) flow of students who become eligible to begin in-hospital training. Second, although the pass rate among first-time test takers is high, some students do fail; in 2015, 4% of first-time test takers had to retake their examination at a later date (U.S. Medical Licensing Examination 2015). For these students, the start of clerkship must be postponed and the sequence of rotations may have to be rescheduled. Third, multiple hospitals may be affiliated with a single medical school's clerkship program, each offering different rotations, on varying starting dates, with distinct capacities in terms of how many students they can accept on each date. Finally, every medical student in a cohort must be assigned to all core rotations. Since the duration of a clerkship program is typically one to two years, multiple cohorts interact and the length of a medical school's planning horizon ranges from 12 to 36 months.

The assignment of students to clinical rotations has a direct impact on the students' careers, their educational experience, and the school's financial status. We address the complexities of clerkship scheduling by formulating the rotation assignment and scheduling problem (RASP). The RASP seeks a minimum-cost policy that assigns a cohort of medical students, who become eligible to start their training over a period of several months, to a series of required clinical rotations, while observing operational and quality-of-service (QoS) requirements. We first present a mixed-integer linear programming (MILP) formulation of the RASP which, although capable of finding optimal solutions, is prohibitively large for practical-sized problem instances. As a result, we propose an alternative network-based formulation (NBF) that is derived from a *decision diagram*, a graphical structure that encodes the set of all feasible clinical rotations, that is, those that satisfy a medical school's operational and QoS requirements. We prove that both formulations are equivalent and show several structural properties of the NBF that establish its theoretical attractiveness and, more importantly, its usefulness for obtaining managerial insights. We then use real data to generate a collection of problem instances to demonstrate that the NBF exhibits significantly superior computational performance over a range of parameter values when compared to the MILP model.

Although our formulation and solution approach are quite general, the study is motivated by our experience with the American University of the Caribbean (AUC) School of Medicine (AUC 2018). Each year, AUC schedules between 300 and 350 students to five core rotations at one or more of 25 affiliated hospitals. Rotations are completed one at a time, in any order, and last either 6 or 12 weeks. Hospitals are located in several US states as well as a number of countries

around the world (e.g., the UK and Australia). Each hospital may offer a subset of the five rotations, with different start dates, capacities, and cost rates. AUC commits to a certain level of service, ensuring, for example, that students complete rotations in a small number of geographic regions and that they are not idle for excessively long periods between rotations. The school's scheduling decisions are constrained, however, by national licensing requirements (Federation of State Medical Boards and National Board of Medical Examiners, 2018) that limit when a student can begin clerkship, operational restrictions on where and when a student can complete a particular rotation, and QoS constraints that ensure AUC provides its students with a positive learning experience.

Of particular concern is the current practice's demand on human resources. Several staff members are responsible for interacting with students and ensuring that the school's scheduling restrictions are obeyed. Further, assignment decisions are made on an ad hoc basis using informal rules and somewhat ambiguous institutional guidelines. Previous attempts at an analytical scheduling approach that captured the full complexity of the problem have either proved too computationally intensive, did not allow for subsequent managerial analysis, or did not take into account rescheduling decisions. As a result, it is important that any proposed method address these issues.

We apply the NBF technique to AUC's RASP by incorporating school-specific constraints. Using six years' worth of student and administrative data collected from AUC, we present two experiments to estimate the cost savings of implementing our scheduling approach in practice. First, we perform a counterfactual analysis by comparing AUC's historical scheduling decisions to those derived from our scheduling policy. We demonstrate that had our scheduling approach been implemented, the estimated savings would have been approximately 19% per year on average. This is due to the better use of capacity at hospitals where AUC has contracts that guarantee a fixed number of rotation slots per cohort. It also indicates that AUC negotiated too many contracts given the number of students they needed to place. We then develop a discrete-event simulation (DES) of the licensing examination and clerkship scheduling process to understand how AUC's strategic decisions can affect future performance. We compare our network-based solution with the school's current scheduling practice (at best, a greedy approach), demonstrating that they can expect an average cost reduction of 25% going forward. We show that this result is robust to mis-specifications in the failure rate in that the difference between the NBF and the greedy approach increases as the failure rate decreases. We also

demonstrate that our approach dominates the school's current scheduling practice for a wide range of congestion levels (i.e., the number of students versus the number of affiliated hospitals and their corresponding capacities). In some cases, a cost reduction of more than 50% is achieved. Only when the size of the cohort is exceedingly small does the greedy approach offer similar performance as compared to the NBF technique. Finally, we discuss our model's applicability as a managerial decision-making tool by interviewing staff about how it can be utilized going forward.

2. Literature Review

Scheduling in the broad context of health care is an important and active area of research. With the appropriate management of resources, the cost of delivering health services can be reduced while simultaneously providing patients with greater and more equitable access to care. Systems of such type have been developed for many areas of the health services space, in particular for scheduling care providers; see, e.g., Hall (2012), Nickel et al. (2012), Guo et al. (2014), Cappanera and Scutellà (2014), Bard et al. (2014), Smet et al. (2016). Our focus is on the scheduling of medical students to educational opportunities that enable them to learn how best to become care providers in the future.

The nurse scheduling/rostering problem has an extensive literature (e.g., Burke et al. 2004, Cheang et al. 2003, Van den Bergh et al. 2013, Warner 1976). The objective is to create work schedules over a planning horizon that ensure the appropriate number of nurses, each with a requisite skill level, are present for each shift. The flexibility of the schedule is limited by organizational constraints, such as the number of working hours before a break must be taken, financial restrictions that limit the amount of overtime that is allowed, and temporal constraints such as the number of consecutive work days. A feature that is often important in nurse scheduling are the preferences of those being scheduled (e.g., seniority, holidays, day/night shifts). From conversations with AUC staff, students do indicate which hospitals or geographic regions they prefer. Nonetheless, these are used as guidelines by the human schedulers, having lower priority than more general QoS requirements.

Several papers address the scheduling of medical residents to shifts, a stage that comes after students graduate from medical school. Cohn et al. (2009) develop a mathematical programming tool that constructs annual on-call schedules for residents specializing in psychiatry at the Boston University School of Medicine. Agarwal (2016) formulate and solve a multi-objective optimization problem for generating

resident schedules for Rochester General Hospital. Bard et al. (2016) develop several optimization-based heuristics to construct annual block schedules for family medicine residents at the University of Texas Health Science Center. Smalley and Keskinocak (2016) propose two integer programming models to assign surgical residents to shifts for the Surgeon Residency Program at Emory University School of Medicine. Finally, Lemay et al. (2017) introduce a pair of algorithms that solve a sequence of optimization problems that, in addition to ensuring appropriate patient coverage and adequate training opportunities, incorporates resident preferences for downtime. A feature that concerns resident scheduling, but not the RASP, is their working conditions (Ludmerer 2009) and the quality of patient care provided. In addition, as addressed in Smalley and Keskinocak (2016), there is a difference between *rotation* and *shift* scheduling for residents; the former being an inter-day problem, and the latter an intra-day problem. For the RASP, intra-day decisions do not play a role.

The RASP is also closely related to classical machine scheduling problems, with students, hospitals, start dates, and capacity constraints playing the role of jobs, machines, release dates, and resource constraints, respectively (Baptiste et al. 2001, Pinedo 2008). Specifically, the rotations can be perceived as tasks of a particular job subject to additional deadline constraints and assignable to any subset of machines that can perform them. Problems of this sort are notoriously difficult to solve efficiently. For example, even if the order of the rotations were fixed and students were already assigned to hospitals, the RASP would reduce to an instance of the job-shop scheduling problem, which is NP-Hard in general (Garey et al. 1976) and still challenging for realistically-sized problems.

Our work is closely related to scheduling models that take advantage of an underlying network structure. Eppen and Martin (1987) present an alternative formulation of the multi-item capacitated lot-sizing problem by reformulating it as a shortest-path problem. Beasley and Cao (1998) analyze the crew scheduling problem and use a dynamic programming approach to represent the complex MILP model as a network. Raffensperger (1999) applies this technique to the cutting stock, tank scheduling, and traveling salesperson problems. Côté et al. (2011) introduce a grammar-based approach to solving a large-scale MILP model where the cost of assigning an identical set of workers to multiple activities during a single day of operation is minimized. Restrepo et al. (2012) analyze a similar setting, but present a column generation approach where the solution of the subproblem (i.e., a feasible shift schedule) is obtained by formulating and solving a shortest-path problem with resource constraints.

In this stream of research, an optimization model is reformulated as a network flow problem on an extended variable space that yields better relaxations and solves more quickly. We add to this research by representing the problem as a decision diagram, which exploits redundancy to reduce the size of the network and is the basis for our network-based MILP model. Moreover, we demonstrate that our network can be used to derive additional managerial insights and as a decision support tool. The idea of a network-based formulation also appears in Zhang et al. (2016). For a comprehensive treatment of extended formulations, see Conforti et al. (2010) and Lancia and Serafini (2018); for decision diagrams and optimization, we refer the reader to Bergman et al. (2016).

The remainder of this study is organized as follows. We formally describe the RASP in section 3 and then introduce a classical MILP model in section 3.1. In section 3.3, we introduce multivalued decision diagrams (MDDs) and, in section 3.5, we use the RASP's MDD to create a network-based formulation (NBF) for the problem. In section 4, we apply our solution approach to AUC by incorporating their specific QoS constraints into our models. Our numerical analysis with both actual and randomly generated RASP instances appears in section 5. Finally, we provide managerial insights and concluding remarks in sections 6 and 7. All proofs are located in Appendix B.

3. The Rotation Assignment and Scheduling Problem

Consider a set of medical students $\mathcal{I} = \{1, \dots, I\}$ that are to be scheduled to a set of clinical rotations $\mathcal{J} = \{1, \dots, J\}$ at a set of hospitals $\mathcal{K} = \{1, \dots, K\}$. Let t index the weeks in the planning horizon, which ranges from $t = 1$, the earliest time any student can start in-hospital training, up to $t = T$, the latest end time of any student's final rotation. Each student $i \in \mathcal{I}$ must complete all rotations $j \in \mathcal{J}$ exactly once before week T . Let $\mathcal{I}(t) \subseteq \mathcal{I}$ be the set of students who are eligible to start clerkship on week t because they have successfully passed their licensing examination prior to week t . Further, for each $i \in \mathcal{I}$, we define t_i to be the earliest week student i can begin clerkship, that is, $t_i = \min\{t | i \in \mathcal{I}(t)\}$.

We consider a deterministic version of the problem where $\mathcal{I}(t)$ is known with certainty for all t . This is a sensible assumption for two reasons. First, the licensing examination is booked 3–6 months in advance and test scores are typically released eight weeks after the examination is taken. Medical schools are generally aware of these dates and typically schedule students by assuming these dates are firm. Second, most students pass their examination; among the 20,213 students who took the US licensing examination in 2015,

96% passed on their first attempt (U.S. Medical Licensing Examination 2015). Thus, although students do fail, this occurrence happens infrequently. Nevertheless, in section 4.2, we address the issue regarding the re-scheduling of students who do not pass their examination.

For each rotation $j \in \mathcal{J}$, let $\Delta_j \in \mathbb{Z}_+$ be its duration in weeks and let $\mathcal{K}(j) \subseteq \mathcal{K}$ be the set of hospitals that offer rotation j . Conversely, let $\mathcal{J}(k)$ be the set of rotations offered by a hospital $k \in \mathcal{K}$. In addition to the set of rotations that are offered, a hospital is defined by several other attributes:

Rotation capacity: Let C_{jkt} be the maximum number of students who can start rotation $j \in \mathcal{J}(k)$, at hospital k , on week t . If students cannot begin a rotation at hospital k , on week t , $C_{jkt} = 0$.

Fixed cost: In exchange for a fixed amount c_k^F , hospital k may sign a multi-year contract with the medical school guaranteeing a pre-determined number of rotation slots per year. The set of hospitals that impose this type of fee structure is denoted by \mathcal{K}_F ; we refer to these institutions as fixed-cost hospitals. Since each fixed-cost hospital may fixed-cost hospitals have different characteristics (e.g., rotations offered, capacities, management procedures), the terms of each contract may also differ. As a result, even if $c_k^F = c_{k'}^F$ for hospitals $k, k' \in \mathcal{K}_F$, the number of guaranteed slots may not be equal. Further, each hospital negotiates a different contract duration (in years) with the medical school. To this end, let $U_k(t)$ be a function that maps the week number t to the appropriate contract *year* for some hospital $k \in \mathcal{K}_F$; let u_k be the maximum duration of any contract that can be signed with hospital $k \in \mathcal{K}_F$. In our setting, the duration of each contract can be no less than one year while multi-year, consecutive and non-consecutive contracts can also be signed. Finally, note that fixed-cost hospitals may still charge a per student fee of c_{jk} for rotation $j \in \mathcal{J}(k)$.

Variable cost: The set of hospitals that do not negotiate fixed contracts, but instead, charge an amount c_{jk} for an assigned student to complete rotation $j \in \mathcal{J}(k)$ is denoted by \mathcal{K}_V . We refer to these institutions as variable-cost hospitals. To facilitate a more compact formulation, we define $U_k(t)$ for a hospital $k \in \mathcal{K}_V$. It is the identity function, $U_k(t) = t$, that maps the week number t to the appropriate *week* for some hospital $k \in \mathcal{K}_V$. The interpretation here is that for variable-cost hospitals, only weekly contracts for a single slot per rotation can be signed with the medical school. Thus, $u_k = T$ represents the maximum weekly contract that can be signed with hospital $k \in \mathcal{K}_V$ over the planning horizon.

In practice, if a fixed-cost hospital charges a rotation cost c_{jk} , it is much smaller than the fee charged by variable-cost hospitals. However, partitioning the set of hospitals into these cost classes allows the medical school to track when their accounts must be settled. For fixed-cost hospitals, payment for in-hospital training is remitted prior to the start of the planning horizon. For variable-cost hospitals, payment can be delayed until just before the start of the student's rotation.

Each student $i \in \mathcal{I}$ must be assigned a schedule S_i which corresponds to a sequence of J stages, one for each core rotation. In stage $s = 1, \dots, J$, the student starts rotation j_s at hospital k_s on week t_s . That is, a schedule for a student i is the sequence

$$S_i = ((j_1, k_1, t_1), (j_2, k_2, t_2), \dots, (j_J, k_J, t_J)). \quad (1)$$

A schedule is feasible if it satisfies two *core* requirements, as follows. These requirements are scheduling restrictions that are typically observed by all medical schools in North America.

- R.1 A schedule must contain all rotations in \mathcal{J} exactly once, in any order.
- R.2 Rotations must be completed one at a time (no overlapping).

A feasible solution to the RASP is an assignment of I feasible schedules to students (i.e., one schedule for every student in a cohort) satisfying the following licensing and capacity constraints.

- R.3 Student i cannot start clinical rotations before week t_i (i.e., before passing the licensing examination).
- R.4 The number of students that start rotation j , at hospital k , on week t , must not exceed C_{jkt} .

The objective of the RASP is to find a collection of feasible schedules for a given cohort of students \mathcal{I} that satisfies the licensing/capacity restrictions and minimizes the sum of fixed and variable costs.

3.1. An MILP Model

We first present a time-indexed MILP formulation for the RASP to determine S_i for every $i \in \mathcal{I}$. This type of model is attractive because different objectives and side constraints can be easily incorporated. It also provides strong dual bounds (Van den Akker et al. 2000) compared to event-based formulations, and is standard when formulating resource constraints (Artigues 2017).

Let x_{ijkt} be a binary variable that equals one if student i is assigned to start rotation j , at hospital k , on week t , and zero otherwise. Let z_{ku} be a binary variable that equals one if there is at least one student assigned to a hospital $k \in \mathcal{K}_F$ during contract year u or

a hospital $k \in \mathcal{K}_V$ during week u and zero otherwise. Then, the MILP formulation of the RASP is given by

$$Z_{RASP} = \min \left\{ \sum_{k \in \mathcal{K}_F} \sum_{u=1}^{u_k} c_k^F z_{ku} + \sum_{i=1}^I \sum_{j=1}^J \sum_{k \in \mathcal{K}(j)} \sum_{t=t_i}^T c_{jk} x_{ijkt} \right\}, \quad (\text{MP})$$

subject to:

$$\sum_{k \in \mathcal{K}(j)} \sum_{t=t_i}^T x_{ijkt} = 1, \quad \forall i, j, \quad (2)$$

$$\sum_{j=1}^J \sum_{k \in \mathcal{K}(j)} \sum_{s=t-\Delta_j+1}^t x_{ijks} \leq 1, \quad \forall i, t \geq t_i, \quad (3)$$

$$\sum_{i \in \mathcal{I}(t)} x_{ijkt} \leq C_{jkt} z_{ku}, \quad \forall j, k \in \mathcal{K}(j), t, u \in U_k(t), \quad (4)$$

$$x_{ijkt} \in \{0, 1\}, \quad \forall i, j, k \in \mathcal{K}(j), t \geq t_i, \quad (5)$$

$$z_{ku} \in \{0, 1\}, \quad \forall k, u \quad (6)$$

The objective minimizes the fixed and variable costs associated with assigning students to hospitals and clinical rotations. Constraint (2) enforces (R.1): each student is assigned to exactly one of each type of rotation and all rotations are completed. Constraint (3) enforces (R.2): no more than one rotation can be completed at any given time. Requirement (R.3) is enforced by not creating x_{ijkt} variables with $t < t_i$. Constraint (4) enforces (R.4) for fixed- and variable-cost hospitals: the maximum capacity of rotation j , at hospital k , during week t . If week t is a valid start date for rotation j at hospital k , $C_{jkt} > 0$; otherwise $C_{jkt} = 0$. It also ensures that $z_{ku} = 1$ whenever a student is assigned to a hospital $k \in \mathcal{K}_F$ during contract year u or, for a hospital $k \in \mathcal{K}_V$, week u .

3.2. Formulation Size: A Case For and Against Column Generation

Problem (MP) has $O(I \times J \times K \times T)$ variables and $O(I \times T)$ constraints. In the case of AUC, this amounts to over 4,000,000 decision variables and approximately 50,000 constraints. Further, many medical schools have QoS requirements, that is, additional scheduling restrictions that are unique to the school (see section 4 for AUC's QoS constraints). When faced

with such a large model, a typical approach is to formulate an alternative MILP model that creates one variable for each feasible schedule, as defined by Equation (1). These schedules are sometimes referred to as *patterns*. In such a model, the complexity of what constitutes a feasible pattern is embedded in a pattern-generation step that is exogenous to the MILP. As a consequence, the model itself becomes relatively straightforward: pick one pattern for each student so that costs are minimized and a small number of simple side constraints are satisfied. The issue with these formulations is the large number of variables (all distinct feasible schedules/patterns), which leads to the popular approach known as *branch-and-price*, that is, column generation embedded inside a branch-and-bound tree search (Barnhart et al. 1998).

Branch-and-price algorithms are notoriously difficult to implement and fine tune for good performance. There are typically several issues that must be addressed, such as choosing an appropriate branching scheme (Vanderbeck 2011), dual stabilization techniques to improve convergence (Clautiaux et al. 2011), primal heuristics (Joncour et al. 2010), and proper management of the pool of generated columns (Barnhart et al. 1998). As a result, we present an alternative approach that is, in many ways, equivalent to column generation but possesses several advantages. This includes a significantly more compact (but still exact) representation of the set of feasible patterns, dual bounds that are just as strong, the ability to perform sensitivity analysis (e.g., see section 6), and a much simpler implementation that takes advantage of existing solvers. As shown in the following sections, the idea is to develop network-based reformulation of the problem, extracted from a dynamic programming model, that is equivalent to (MP) in that it captures all feasible student schedules.

3.3. A Multivalued Decision Diagram for the RASP

We present a reformulation of the RASP based on decoupling the assignment and scheduling elements of the problem. This is achieved by encoding all feasible schedules into a *multivalued decision diagram* (MDD) (Bergman et al. 2016). This network-style representation carries significantly more structure than the MILP model. In particular, in certain settings, the optimal solution can be obtained by solving a minimum-cost flow problem on the MDD. Further, this MDD can help derive additional managerial insights associated with rescheduling students and dealing with undesired assignments.

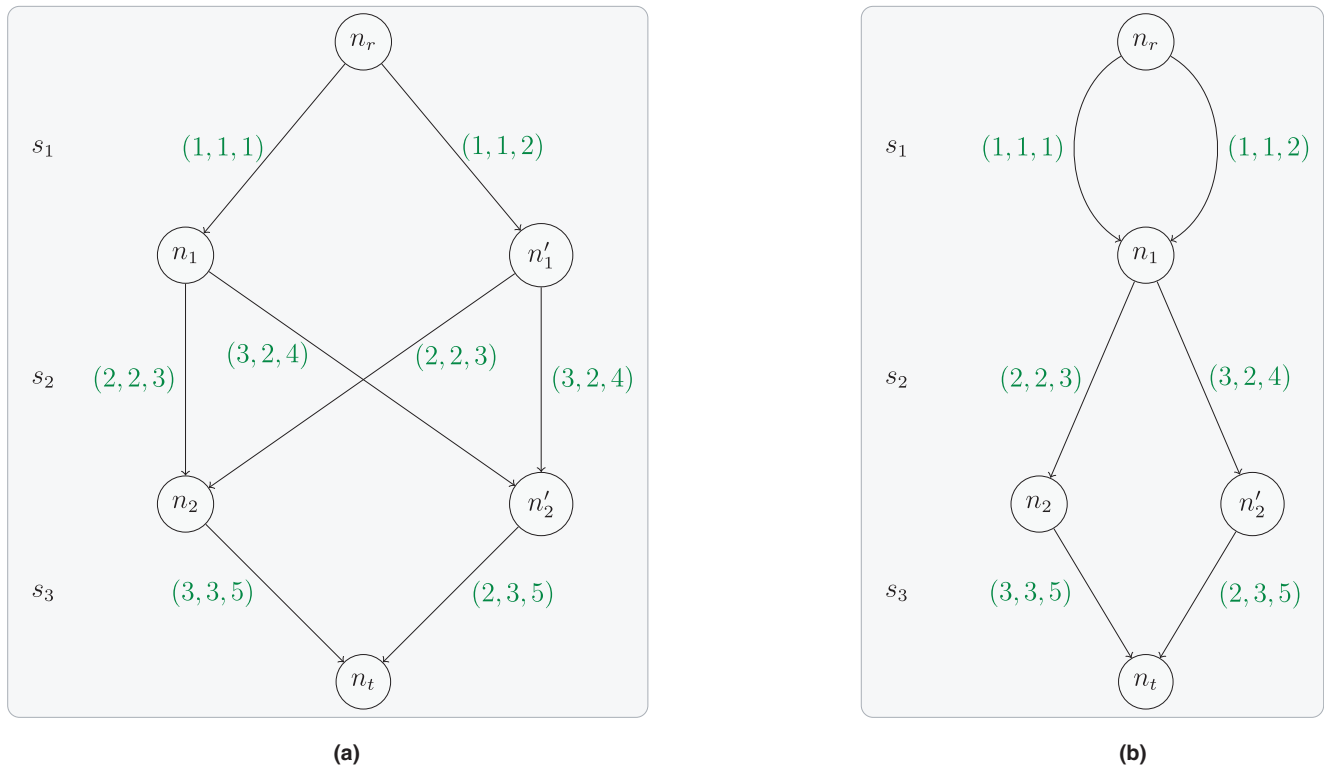
Formally, we create a directed acyclic graph $D = (N, A)$ that compactly encodes the set of feasible schedules \mathcal{S} that can be assigned to a student. The set

of nodes N is partitioned into $J + 1$ layers L_0, \dots, L_J where L_0 and L_J are singletons containing a root node n_r and a terminal node n_t , respectively. Arcs in A only connect nodes in adjacent layers of D . An arc $a = (n, n')$, with $n \in L_s$, $n' \in L_{s+1}$, is associated with the triplet $\theta(a) = (j_s, k_s, t_s)$ which represents the s -th stage in the student's schedule. Thus, a path of length J specified by the sequence of arcs (a_1, a_2, \dots, a_J) starting at n_r and terminating at n_t (also known as an n_r, n_t -path) encodes a feasible schedule $\mathcal{S} = (\theta(a_1), \theta(a_2), \dots, \theta(a_J))$. In this graphical representation, there exists a one-to-one mapping between the set of feasible schedules and all n_r, n_t -paths in D . Using the network alone, the RASP (without R.3 and R.4) now reduces to finding a set of paths of length J in D with the lowest total cost such that all medical students are assigned a schedule and all core requirements are met.

EXAMPLE 1. Consider an instance of the RASP with three hospitals ($K = 3$) and three rotations ($J = 3$), each lasting one week. The first hospital offers rotation 1 starting on weeks $t = \{1, 2\}$ (i.e., C_{111} and $C_{112} > 0$). The second hospital offers rotations 2 and 3 starting on weeks $t = 3$ and $t = 4$, respectively. The third hospital also offers rotations 2 and 3, but they both start on week $t = 5$.

Two MDDs encoding the set of feasible schedules for Example 1 are presented in Figure 1. The path (n_r, n_1, n'_2, n_t) in Figure 1a encodes a schedule where a student starts rotation 1 at the first hospital on week $t = 1$, proceeds to rotation 3 at the second hospital on week $t = 4$, and starts rotation 2 at the third hospital on week $t = 5$. Figure 1b shows the exact same set of schedules in a more compact fashion. Our methodology constructs the smallest of such graphs, known as a *reduced MDD*, for any arbitrary instance of the RASP. This is achieved by exploiting the fact that feasible schedules (i.e., paths through the network) often share equivalent subsequences, leading to redundancies that can be eliminated. For example, in Figure 1b, the two paths/schedules starting with $(1, 1, 1)$ and $(1, 1, 2)$ have the same subsequences in the following layers and, hence, do not need to be represented by two distinct nodes in their layer as they were in Figure 1a.

The MDD is constructed by developing a dynamic program (DP). Nodes encode the states of the system (i.e., partial feasible schedules) while arcs represent feasible actions (i.e., possible rotation assignments) that can be taken at each state. Hence, the DP enumerates all possible states, which are implicitly encoded as paths from n_r to n_t in the MDD. In our case, paths are constructed using a recursive model that considers both the set of clinical rotations left to complete at

Figure 1 Two Distinct Multivalued Decision Diagrams Encoding the Same Set of Feasible Rotation Schedules of Example 1 [Color figure can be viewed at wileyonlinelibrary.com]

each stage and the weeks in which the remaining rotations can begin. Thus, the MDD is a compressed version of the underlying state-transition graph of the DP, where isomorphic subtrees are merged and costs are attributed to arcs as opposed to state transitions (Hooker 2013).

We now present a dynamic programming formulation that models a feasible student schedule over a finite planning horizon. Such a model provides the theoretical basis for calculating the size of D which, for practical instances of our problem, is sufficiently small.

PROPOSITION 1. *A schedule S , as defined in Equation (1), is feasible if and only if there exists a sequence of rotations-to-go (\mathcal{G}) and earliest start time (τ) pairs $(\mathcal{G}_0, \tau_0), (\mathcal{G}_1, \tau_1), (\mathcal{G}_2, \tau_2), \dots, (\mathcal{G}_J, \tau_J)$, where $\mathcal{G}_s \subseteq \mathcal{J}$ and $\tau_s \in \mathbb{Z}_+$ for $s = 0, \dots, J$, such that $(\mathcal{G}_0, \tau_0) = (\mathcal{J}, 0)$ and*

$$\mathcal{G}_s = \mathcal{G}_{s-1} \setminus \{j_s\}, \quad (7)$$

$$\tau_s = t_s + \Delta_{j_s}, \quad (8)$$

$$(j_s, k_s, t_s) \in \mathcal{A}(\mathcal{G}_{s-1}, \tau_{s-1}), \quad (9)$$

for $s = 1, \dots, J$, where $\mathcal{A}(\mathcal{G}_s, \tau_s) = \{(j, k, t) : j \in \mathcal{G}_s, k \in \mathcal{K}(j), t \geq \tau_s\}$.

A consequence of Proposition 1 is that every feasible schedule can be obtained by enumerating all (\mathcal{G}_s, τ_s) pairs. Assuming a planning horizon of T weeks, the total number of pairs is bounded by $O(2^J \times T)$, which is small in many practical instances (e.g., if $J = 5$ and $T = 104$, $2^J \times T = 3328$). Even in large-scale implementations, the number of such pairs can be exponentially smaller than the set of feasible schedules. In fact, in section 4, we show that the number of (\mathcal{G}_s, τ_s) pairs remains tractable even when AUC's QoS constraints are incorporated into the network.

3.4. MDD Construction

The MDD is built in two phases. In the first phase, we construct a directed acyclic graph $D = (N, A)$ encoding the set of transitions of the recursive model from Proposition 1. Specifically, the set of nodes N is partitioned into layers L_0, L_1, \dots, L_J . There exists a one-to-one mapping between a node $n \in L_s$ and a pair (\mathcal{G}_s, τ_s) that is reachable at stage s by some feasible schedule. The pair associated with a node is henceforth denoted by $\zeta(n)$, e.g. $\zeta(n_r) = (\mathcal{G}_0, \tau_0)$. An arc $a = (n, n')$ connects nodes in adjacent layers, $n \in L_s, n' \in L_{s+1}$, if and only if there exists a tuple $\theta(a) = (j_s, k_s, t_s) \in \mathcal{A}(\zeta(n))$ in some feasible schedule which transitions $\zeta(n)$ to $\zeta(n')$ according to Equations (7)–(8). This procedure is summarized in Algorithm 1.

It considers each pair $\xi(n)$ (i.e., a node in the network) exactly once. Because a node can have at most $|\mathcal{A}(\xi(n))|$ outgoing arcs, the algorithm to generate the network has a (conservative) worst-case time complexity of $O(2^J \times T) \times O(J \times K \times T) = O(J \times 2^J \times K \times T^2)$.

EXAMPLE 2. Figure 1a depicts the MDD obtained after applying Algorithm 1 to the instance described in Example 1. In particular, $\xi(n_r) = (\{1, 2, 3\}, 0)$, $\xi(n_1) = (\{2, 3\}, 2)$, $\xi(n'_1) = (\{2, 3\}, 3)$, $\xi(n_2) = (\{3\}, 4)$, $\xi(n'_2) = (\{2\}, 5)$, and $\xi(n_t) = (\emptyset, 6)$.

Algorithm 1: MDD Compilation Algorithm

```

1: Create node  $n$  associated with pair  $\xi(n) = (\mathcal{G}_0, \tau_0)$ .
2: Let  $L_0 = \{n_r\}$ ,  $A = \emptyset$ .
3: For each  $s = 1, \dots, J - 1$ :
  (a) Let  $L_{s+1} = \emptyset$ .
  (b) For each node  $n \in L_s$  and  $(j_s, k_s, t_s) \in \mathcal{A}(\xi(n))$ :
    (i) Let  $n'$  be a node such that  $\xi(n') = (\mathcal{G}_{s+1}, \tau_{s+1})$  in  $L_{s+1}$  according to Equations (7)–(8). If no such node exists, add it to  $L_{s+1}$ .
    (ii) Create arc  $a = (n, n')$  associated with  $\theta(a) = (j_s, k_s, t_s)$  and add it to  $A$ .
4: Return  $D = (N, A)$ .
```

The second phase consists of merging *equivalent* nodes of D to remove redundancy. Two nodes $n, n' \in L_s$ are considered equivalent if, for any arc-specified path (a_1, \dots, a_{J-s}) starting at node n and ending at a node $\tilde{n} \in L_J$, there exists another path (a'_1, \dots, a'_{J-s}) starting at n' and ending at \tilde{n} such that $\theta(a_i) = \theta(a'_i)$ for all $i = 1, \dots, J - s$. That is, any partial schedule that reaches a state $\xi(n) = (\mathcal{G}_s, \tau_s)$ can be completed in the same way as any partial schedule that reaches $\xi(n') = (\mathcal{G}'_s, \tau'_s)$. The task of identifying and merging equivalent nodes is called *reduction* and is described in Algorithm 2. The reduction phase has a quadratic worst-case time complexity in the number of nodes of D . More sophisticated methods can be implemented with a linear-time complexity (Wegener 2000). We note that the reduction operation is a critical step in our context. In our numerical experiments, MDDs were approximately 50% smaller on average after the second phase as compared to the first.

Algorithm 2: MDD Reduction Algorithm

```

1: Merge all nodes of  $L_J$  into a single terminal node  $n_t$ .
2: For  $s = J - 1, J - 2, \dots, 1$ :
  (a) For each pair of nodes  $n, n' \in L_s, n \neq n'$ :
    If for each arc  $a = (n, \tilde{n})$  there exists an arc  $a' = (n', \tilde{n})$  with  $\theta(a) = \theta(a')$  for some  $\tilde{n} \in L_{s+1}$ , redirect incoming arcs of  $n'$  to  $n$  and remove  $n'$  from  $L_s$ .
```

EXAMPLE 3. Figure 1b depicts the result of applying Algorithm 2 to the MDD in Figure 1a. Starting in the

second-to-last layer, nodes n_2 and n'_2 cannot be merged because $\theta((n_2, n_t)) \neq \theta((n'_2, n_t))$. Going up one layer, for every arc leaving node n_1 there exists an arc leaving node n'_1 with the same endpoint (either n_2 or n'_2) and the same value of $\theta(\cdot)$. Therefore, n_1 and n'_1 can be merged.

3.5. A Network-Based Formulation

The MDD D encodes the set of all feasible schedules, that is, those that satisfy core requirements R.1 and R.2. Thus, determining feasible schedules that minimize the sum of all variable costs can be obtained by solving a minimum-cost flow problem. All feasible solutions to the RASP are feasible schedules in D , although not all feasible schedules are feasible solutions to the RASP in that they may not satisfy R.3 and R.4. Further, fixed-cost hospitals must also be included in the formulation. Thus, we present a network-based MILP formulation (NBF) of the RASP which uses the MDD's underlying graphical structure as input. The network encodes all feasible schedules while the MILP includes fixed-cost hospitals and enforces R.3 and R.4 to find the minimum-cost solution.

LEMMA 1. Given an MDD $D = (N, A)$, let $\delta^+(n)$ and $\delta^-(n)$ denote the set of outgoing and incoming arcs at a node $n \in N$, respectively. For any $m \in \mathbb{Z}_+$, there exists a one-to-one correspondence between m -sized subsets of feasible schedules and integral points of the network-flow polyhedron

$$\text{flow}(D, m) = \left\{ \mathbf{y} \in \mathbb{R}_+^{|A|} : \sum_{a \in \delta^+(n_r)} y_a = m \text{ and } \sum_{a \in \delta^-(n)} y_a - \sum_{a \in \delta^+(n)} y_a = 0, \forall n \in N \setminus \{n_r, n_t\} \right\},$$

where y_a is a non-negative decision variable that corresponds to the amount of flow on arc a .

By Lemma 1, this reformulation of the RASP consists of finding an m -unit flow over D satisfying R.3 and R.4. For each arc $a \in A$, with a slight abuse of notation, let j_a, k_a , and t_a be the rotation, hospital, and start week associated with $\theta(a)$, respectively. In addition, let $\mathcal{C}(j, k, t) \subseteq A$ be the set of arcs a with $j_a = j, k_a = k$, and $t_a = t$, and let $c_{j_a k_a}$ be the rotation cost per unit of flow on arc a . If y_a is the amount of flow on arc a , then the reformulation of the RASP is given by

$$Z_{\text{RASP}} = \min \left\{ \sum_{k \in K_F} \sum_{u=1}^{u_k} c_k^F z_{ku} + \sum_{a \in A} c_{j_a k_a} y_a \right\}, \quad (\text{NBF})$$

subject to:

$$\mathbf{y} \in \text{flow}(D, I), \quad (10)$$

$$\sum_{a \in \mathcal{C}(j,k,t)} y_a \leq C_{jkt} z_{ku}, \quad \forall j, k \in \mathcal{K}(j), t, u \in U_k(t), \quad (11)$$

$$\sum_{\{a \in \delta^+(n_r) : t_a \geq t\}} y_a = \sum_{s \geq t} |\mathcal{I}(s) \setminus \mathcal{I}(s-1)|, \quad \forall t, \quad (12)$$

$$y_a \in \{0, \dots, I\}, \quad \forall a, \quad (13)$$

$$z_{ku} \in \{0, 1\}, \quad \forall k, u. \quad (14)$$

where we assume $\mathcal{I}(0) = \emptyset$. The objective of (NBF) is equivalent to that of (MP). Constraint (10) ensures \mathbf{y} is a subset of I feasible student schedules (by Lemma 1). Constraint (11) plays the role of (4), as $\mathcal{C}(j, k, t)$ aggregates flow units (students) with the same rotation j , hospital k , and start week t . Constraint (12) ensures that students do not start rotations before they become available.

A benefit to using an MDD for the reformulation is that its network structure can be leveraged to derive efficient algorithms for special cases of NBF that are relevant in practice. For instance, AUC does not often make joint contracting and scheduling decisions. Instead, they first negotiate fixed contracts with hospitals and then assign students to rotations several months later. In this special case, $c_k^F = 0$ for all $k \in \mathcal{K}_F$ as rotation slots for fixed-cost hospitals have already been guaranteed. Thus, it is advantageous to assign as many students to these hospitals as possible. AUC may also wish to impose additional scheduling restrictions when defining a feasible assignment of students to rotations. As long as these extra requirements can be incorporated directly in the decision diagram (which may increase the size of the graph), the next proposition demonstrates that the problem can still be solved efficiently in the size of the network.

PROPOSITION 2. *Let $c_k^F = 0$ for all $k \in \mathcal{K}_F$ and let $D = (N, A)$ be constructed so that all feasible paths in the network satisfy any additional requirements imposed by the medical school. If the hospital capacities are sufficiently large, (NBF) can be solved in strongly polynomial time in the size of D .*

Moreover, in some medical schools, students are required to complete clerkship in a specific order (Wimmers et al. 2007). In this case, optimal schedules can also be derived efficiently in the size of D if no additional requirements are imposed by the school.

PROPOSITION 3. *Let $c_k^F = 0$ for all $k \in \mathcal{K}_F$ and suppose that no additional requirements are imposed by the medical school. If students must complete their rotations in a*

pre-determined sequence, the optimal solution to (NBF) can be obtained in polynomial time in the size of D .

We note that the network structure of D implicitly handles the complexity of the sequencing decisions in this reformulation. However, the general version of the RASP, which includes hospital capacities and possibly more general scheduling requirements (e.g., AUC's QoS constraints), still remains theoretically hard. We formalize this statement below.

PROPOSITION 4. *The decision version of (NBF) is strongly NP-Hard even when $c_k^F = 0$ for all $k \in \mathcal{K}_F$ and D is polynomial in the size of the problem input.*

4. Incorporating Quality-of Service Requirements: AUC Case Study

In this section, we demonstrate how to incorporate QoS requirements into our model. These requirements are unique to the individual medical school. They represent additional constraints that must be added to the RASP to ensure rotation schedules satisfy internal procedures and are well-received by prospective students. For AUC, these requirements are:

- R.5 *Student preferences:* Students may request to be placed at hospitals in specific areas of the country. Let $r_k \in \{1, \dots, R\}$ identify the geographic region where hospital k is located. For example, at AUC, some of the affiliated hospitals are located in New York City, London, England, and Miami, Florida. The medical school fulfills these requests whenever possible.
- R.6 *Relocation:* It is costly for students to resettle in different geographic regions. Thus, each student must complete all of the required rotations with no more than \bar{r} region relocations and no back-and-forth movement among them. In addition, as a measure of QoS, AUC ensures that at least $\alpha\%$ of their students complete all of their rotations within a single region.
- R.7 *Idleness:* For each student, there can be no more than ϵ weeks of separation between the end of one rotation and the beginning of the next rotation.
- R.8 *All-or-none:* Any student assigned to a fixed-cost hospital must complete all the rotations offered by that hospital (in an effort to maximize the use of the fixed contract).

The combination of the *Relocation* and *Idleness* restrictions ensure that the final schedule does not

place an undue burden on the capability of a student to complete their clerkship to the best of their ability. For example, AUC ensures that students are idle for no more than eight weeks between rotations ($\epsilon = 8$), that students can be assigned to at most two geographic regions ($\bar{r} = 1$), and that 75% of all students must be assigned to exactly one region ($\alpha = 0.75$). Note that there are more hospitals than regions, as some hospitals are assigned to the same geographic region. For example, AUC is affiliated with several hospitals in New York City which are all assigned to the same region. A hospital in Miami, Florida, however, would be considered an entirely different geographic region.

We note that (R.5)–(R.8) represent the standard QoS constraints that AUC must obey when scheduling students to rotations. Additional idiosyncratic restrictions that represent each student's personal preferences may be added to the model. In this case, the solution becomes highly specialized and provides one schedule for every student in the cohort. These restrictions, however, vary from year to year and, thus, we focus on QoS requirements that remain constant for all cohorts. The added benefit is that a feasible schedule can now be assigned to any student i provided $t_i \geq t$. More specifically, the output of the model represents flow on the network D . Any path from n_r to n_t with positive flow is a sequence of rotations that a student could follow. The value of the flow equals the number of students who will have that exact same schedule. Therefore, when students sit with their advisors to pick a schedule, they can be presented with several options (i.e., different paths in D) that satisfy their starting date t_i and are not yet assigned to other students. A path may be (subjectively) more or less desirable depending on a student's personal preferences.

4.1. MILP Representation of AUC's QoS Requirements

We first add AUC's QoS requirements to the time-indexed MILP formulation for the RASP introduced in section 3.1, that is, problem (MP). Let v_{it} be a binary variable that equals one if student i completes their last rotation at time t and zero otherwise. Let w_{ir} be a binary variable that equals one if student i is assigned to a hospital in geographic region r and zero otherwise. Finally, let \tilde{w}_i be a binary variable that equals zero if student i completes their entire training in a single geographic region, and equals one if i 's rotations take place in two different regions (recall that $\bar{r} = 1$ for AUC).

$$\sum_{t=t_i}^T v_{it} = 1, \quad \forall i, \quad (15)$$

$$\sum_{k \in \mathcal{K}(j)} x_{ijkt} - \sum_{s=t+\Delta_j}^{t+\Delta_j+\epsilon} \sum_{j' \neq j} \sum_{k \in \mathcal{K}(j')} x_{ij'ks} \leq v_{it}, \quad \forall i, j, t \geq t_i, \quad (16)$$

$$v_{it} + \sum_{s=t+1}^T \sum_{k \in \mathcal{K}(j)} x_{ijks} \leq 1, \quad \forall i, j, t \geq t_i, \quad (17)$$

$$\sum_{t=t_i}^T x_{ijkt} \leq w_{ir_k}, \quad \forall i, j, k \in \mathcal{K}(j), \quad (18)$$

$$\sum_{r=1}^R w_{ir} \leq 1 + \tilde{w}_i, \quad \forall i, \quad (19)$$

$$\sum_{i=1}^I \tilde{w}_i \leq 0.25I, \quad (20)$$

$$\begin{aligned} & \forall \{k_1, k_2, k_3\} \in \mathcal{K} | r_{k_1} \neq r_{k_2}, r_{k_2} \neq r_{k_3}, x_{ij_1 k_1 t_1} + x_{ij_2 k_2 t_2} \\ & + \sum_{s=t_2+1}^T x_{ij_3 k_3 s} \leq 2, i, t_i \leq t_1 < t_2, \text{ distinct } j_1, j_2, j_3, \\ & j_1 \in \mathcal{J}(k_1), j_2 \in \mathcal{J}(k_2), j_3 \in \mathcal{J}(k_3), \end{aligned} \quad (21)$$

$$\sum_{t=t_i}^T x_{ijkt} = \sum_{t=t_i}^T x_{ij'kt}, \quad \forall i, k \in \mathcal{K}_F, j \text{ and } j' \in \mathcal{J}(k), j < j', \quad (22)$$

$$w_{ir}, \tilde{w}_i, v_{it} \in \{0, 1\}, \quad \forall i, r, t, \quad (23)$$

Constraints (15)–(17) enforce (R.7). That is, constraint (15) ensures that each student completes all core rotations before week T , constraint (16) ensures that two consecutive rotations (when $v_{it} = 0$) cannot be separated by more than ϵ weeks, and constraint (17) prevents assignments after the last rotation is completed (when $v_{it} = 1$). Although this last constraint is redundant, it strengthens the model. Constraints (18)–(21) enforce (R.6). Constraint (18) sets w_{ir_k} to one when student i is assigned to hospital k , constraint (19) sets \tilde{w}_i to one if two w_{ir} variables with distinct r values are set to one for student i (meaning i visited two different regions), and constraint (20) ensures that at most 25% of AUC's students complete rotations in two different geographic regions, that is, have $\tilde{w}_i = 1$. Note that constraint (21) is necessary to ensure that, if a student completes their rotations in two different regions, they cannot oscillate between them, i.e. once they change to the second region, they cannot change back. Constraint (22) enforces (R.8): if a student is assigned to a fixed-cost hospital, they complete all rotations offered at that institution. Finally, notice that (R.5) is not explicitly enforced in the MILP model. Instead, once a solution is obtained,

students meet with an advisor and are presented with a few schedules from which to choose. By incorporating (15)–(23) in (MP), we obtain the AUC-specific RASP formulation.

Constraints (15)–(23) highlight the challenge of using an MILP approach for this type of problem. Requirements that are easily stated can be difficult to formalize mathematically. In our case, an additional $O(I \times R + I \times T)$ variables and $O(I \times J \times T)$ constraints are added to (MP). For practical-sized instances, this amounts to tens of thousands of additional variables, and hundreds of thousands of additional constraints, resulting in increased computational complexity.

4.2. A Network-Based Formulation of AUC's QoS Requirements

We now discuss how to encode the QoS requirements (R.6), (R.7), and (R.8) in the RASP's MDD, and consequently, in the NBF. Although the overall construction and reduction procedures, as presented in Algorithms 1 and 2, remain the same, the state information and transition functions of the DP must be augmented. That is, for every stage s , we replace the original state information pair (\mathcal{G}_s, τ_s) with $(\mathcal{G}_s, \tau_s, \eta_s, \zeta_s, \mathcal{F}_s)$, where the three new pieces of information are defined as follows:

- $\eta_s \in \mathbb{Z}_+$: the number of distinct regions to which the student has been assigned to up to stage s ;
- $\zeta_s \in \{0, 1, \dots, R\}$: the region to which the student is assigned to in stage s ;
- $\mathcal{F}_s \in [\{0\} \cup \mathcal{K}_F]^J$: an ordered tuple of length J where the number at position j , $\mathcal{F}_s(j)$, indicates whether or not rotation j was completed at a fixed-cost hospital. If $\mathcal{F}_s(j) = 0$, either rotation j has not been completed by stage s or it was completed at a hospital in $\mathcal{K} \setminus \mathcal{K}_F$. Otherwise, $\mathcal{F}_s(j)$ identifies the hospital from \mathcal{K}_F where rotation j was completed.

Given the augmented state space, feasible transitions in the MDD must abide by the QoS constraints.

- (R.7): For any s , if a rotation finishes at time τ_s , the next rotation must start no later than $\tau_s + \epsilon$.
- (R.6): Recall that \bar{r} is the maximum number of region changes a student can experience. Let $\mathbb{I}(\ell)$ be an indicator function that evaluates to 1 if condition ℓ is true and 0 otherwise. For any schedule \mathcal{S} as in (1), we require that the number of region changes be no more than \bar{r} ,

$$\sum_{s=2}^J \mathbb{I}(r_{k_{s-1}} \neq r_{k_s}) \leq \bar{r}. \quad (24)$$

- (R.8): Given a schedule \mathcal{S} , define $\mathcal{J}(\mathcal{S}, k) = \{j_s | k_s = k \text{ for } (j_s, k_s, t_s) \in \mathcal{S}, s = 1, \dots, J\}$ to be the set of rotations performed at a hospital k . For any fixed-cost hospital $k \in \mathcal{K}_F$, a feasible \mathcal{S} must be such that either all or none of the rotations in $\mathcal{J}(k)$ are performed at k ,

$$(\mathcal{J}(k) = \mathcal{J}(\mathcal{S}, k)) \text{ or } (\mathcal{J}(\mathcal{S}, k) = \emptyset). \quad (25)$$

We now present an updated version of Proposition 1 that recursively defines a feasible schedule using $(\mathcal{G}_s, \tau_s, \eta_s, \zeta_s, \mathcal{F}_s)$. To simplify our exposition when it comes to enforcing (R.8), let $\mathcal{H}(\mathcal{G}_s, \mathcal{F}_s, j)$ be the set of hospitals where a student can perform rotation j , given that rotations in \mathcal{G}_s are left to be completed, and fixed-cost hospitals in \mathcal{F}_s have been used. Such a set is defined as follows:

- If all entries in \mathcal{F}_s are zero, no fixed-contract hospitals have been assigned to the partial schedule up to stage s . Thus, we have $\mathcal{H}(\mathcal{G}_s, \mathcal{F}_s, j) = \{k : j \in \mathcal{J}(k) \text{ and } (k \in \mathcal{K} \setminus \mathcal{K}_F \text{ or } k \in \mathcal{K}_F \wedge \mathcal{J}(k) \subseteq \mathcal{G}_s)\}$. That is, rotation j can either be completed at a variable-cost hospital or at a fixed-cost hospital provided, due to the *all-or-none* requirement, that all rotations offered by that institution are in the set \mathcal{G}_s .
- Otherwise, we must verify whether a fixed-cost hospital that offers rotation j has been assigned to the partial schedule in a previous stage. This leads to two conditions:
 - If there exists some rotation $j' \neq j$ such that $\mathcal{F}_s(j') = k'$ for some fixed-cost hospital $k' \in \mathcal{K}_F$ and hospital k' offers rotation j , then $\mathcal{H}(\mathcal{G}_s, \mathcal{F}_s, j) = \{k'\}$. In this case, some rotations were already completed at fixed-cost hospital k' , and because $j \in \mathcal{J}(k')$, due to the *all-or-none* requirement, rotation j must also be completed at hospital k' .
 - Otherwise, $\mathcal{H}(\mathcal{G}_s, \mathcal{F}_s, j)$ is defined as in (a).

PROPOSITION 5. A schedule \mathcal{S} , as defined in Equation (1), is feasible if and only if there exists a sequence of tuples $(\mathcal{G}_0, \tau_0, \eta_0, \zeta_0, \mathcal{F}_0), \dots, (\mathcal{G}_J, \tau_J, \eta_J, \zeta_J, \mathcal{F}_J)$ with $\mathcal{G}_s \subseteq \mathcal{J}$, $\tau_s \in \mathbb{Z}_+$, $\eta_s \in \mathbb{Z}_+$, $\zeta_s \in \{0, 1, \dots, R\}$, $\mathcal{F}_s \in [\{0\} \cup \mathcal{K}_F]^J$ for $s = 0, \dots, J$, such that $(\mathcal{G}_0, \tau_0, \eta_0, \zeta_0, \mathcal{F}_0) = (\mathcal{J}, 0, 0, 0, (0)^J)$ and, for $s = 1, \dots, J$,

$$\mathcal{G}_s = \mathcal{G}_{s-1} \setminus \{j_s\},$$

$$\tau_s = t_s + \Delta_{j_s},$$

$$\eta_s = \eta_{s-1} + \mathbb{I}(r_{k_{s-1}} \neq r_{k_s}), \quad (26)$$

$$\zeta_s = r_{k_s}, \quad (27)$$

$$\mathcal{F}_s = \mathcal{F}_{s-1} + \mathbb{I}(k_s \in \mathcal{K}_F)k_s e_{j_s}, \quad (28)$$

$$(j_s, k_s, t_s) \in \tilde{\mathcal{A}}(\mathcal{G}_{s-1}, \tau_{s-1}, \eta_{s-1}, \zeta_{s-1}, \mathcal{F}_{s-1}), \quad (29)$$

where e_j is a vector of length J with value 1 at component j and zeros otherwise, and

$$\tilde{\mathcal{A}}(\mathcal{G}_s, \tau_s, \eta_s, \zeta_s, \mathcal{F}_s) = \{(j, k, t) : j \in \mathcal{G}_s, k \in \mathcal{H}(\mathcal{G}_s, \mathcal{F}_s, j),$$

$$\tau_s \leq t \leq \tau_s + \varepsilon, \eta_s + \mathbb{I}(\zeta_s \neq r_k) \leq \bar{r} + 1\}.$$

Proposition 5 demonstrates that by using the augmented tuple $(\mathcal{G}, \tau, \eta, \zeta, \mathcal{F})$ and the valid transition function $\tilde{\mathcal{A}}(\cdot)$ in the MDD compilation procedure, requirements (R.6), (R.7), and (R.8) can be incorporated directly into the graphical structure of D . However, it has more general implications. Specifically, it suggests that modelers of large-scale problems have additional flexibility with regard to how complex requirements may be represented in their solution approach. Some restrictions may be better expressed using a mathematical programming representation while others can be incorporated into a decision diagram via dynamic programming. Thus, the NBF approach allows the modeler to examine how best to express problem requirements by experimenting between the MDD and the network-based MILP representation and studying the effect on computational performance. Another important result is that, even with the additional AUC constraints, the MDD grows sub-exponentially in all parameters except for J .

PROPOSITION 6. *The size of the MDD increases exponentially only in the number of rotations J .*

Proposition 6 implies that our approach can accommodate large numbers of students, time periods, and hospitals. Specifically, if the number of rotations is fixed, the size of the network that underlies the MDD is pseudo-polynomial with respect to the input, as it depends on the length of the planning horizon, but is polynomial in all remaining parameters. For this problem setting, the result is especially advantageous as the number of core rotations is typically small ($J = 5$). They are standardized by the Council for Graduate Medical Education (2017)—an independent organization that sets the standards for the delivery of medical education in the US—and are unlikely to change

often. Other problem parameters, however, may fluctuate significantly from year to year.

The efficiency of the solution procedure also has implications with regards to students who fail their licensing examination and whose sequence of rotations must be rescheduled. Such students require schedules that begin after the date of their next exam. The next proposition demonstrates that the MDD network can efficiently identify an optimal set of alternative schedules with the lowest cost.

PROPOSITION 7. *Under the assumptions of Proposition 2, the b least-cost schedules can be obtained in polynomial time in the size of D for any integer $b > 0$.*

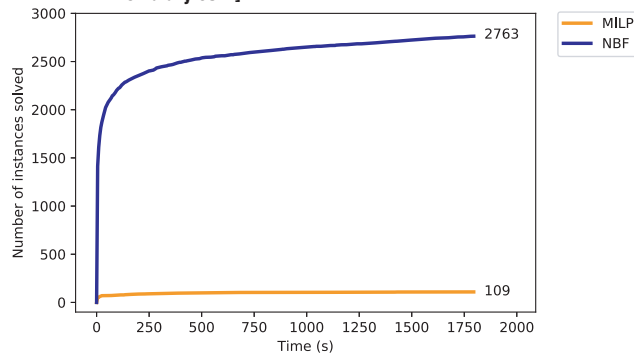
The result above implies that, although we assume a deterministic formulation for the problem, our solution methodology can efficiently account for the stochastic features of the real system.

5. Computational Experiments: An AUC Case Study

In this section, we present a comprehensive numerical study of the RASP. We first provide a performance analysis to demonstrate that the solution times, number of instances solved, and problem sizes that can be solved to optimality using the NBF approach (sections 3.3 and 3.5) are better than the MILP model (section 3.1). We then perform a counterfactual study, parametrized by six years' worth of student and administrative data, that compares AUC's historical scheduling decisions to what would have happened had the NBF framework been implemented. To understand how the optimal assignment of students to rotations can affect future costs, we develop a DES of the licensing examination and clerkship scheduling process. We use the DES to compare the long-run costs of our solution technique to a greedy heuristic. Finally, we perform a sensitivity analysis to determine how future costs may be affected by estimation uncertainty (i.e., the licensing exam failure rate) and forecast uncertainty (i.e., the future size of an AUC cohort).

For the counterfactual analysis and the DES, we compare the NBF framework to a greedy benchmark. The greedy algorithm first orders students in decreasing order of t_i . Then, one student at a time, it sequentially schedules them to a series of rotations with the minimum cost. Notice that the algorithm first schedules the student with the latest start week and progressively works backwards until it gets to the current week; this reduces the chance of generating an infeasible schedule. Interviews with staff, along with a retrospective comparison of schedules (see section 5.2), indicate that the greedy algorithm is the most coherent way to operationalize AUC's

Figure 2 Number of Problems Solved to Optimality in under 30 Minutes (out of 3240) [Color figure can be viewed at wileyonlinelibrary.com]



Notes: Test instances were generated by setting the planning horizon to $T = 52$ weeks and then varying the cohort size (10, 25, 50, 100, 200, 300), the number of hospitals (5, 10, 20, 30), the proportion of hospitals located in unique regions (20%, 40%, 60%), the percentage of fixed- versus variable-cost hospitals (25%, 50%, 75%), the number of rotations (3, 5), the rotation lengths in weeks (3, 5, 7), and the maximum idle time between rotations in weeks (3, 5, 7).

current scheduling practices. This is because scheduling decisions at AUC are made in an ad-hoc manner where staff do not attempt to solely minimize cost. Instead, staff use their experience to assign each student to a schedule (as opposed to assigning a schedule to a student) and periodically consider idiosyncratic student preferences when making assignment decisions. Nevertheless, whenever cost considerations are taken into account, placements are assigned greedily. Additionally, we did not design new heuristic approaches to use as benchmarks since they would return suboptimal solutions and the NBF approach is exact.

The MILP and NBF were solved using Gurobi 7.5.2 (Gurobi Optimization, 2016). The DES was developed in Java. Simulation experiments were performed on an Intel® Core™ i7-3770 with a 3.40GHz processor and 16 GB of RAM. For each simulation, we employed the independent-replications method with a burn-in of 50 years and trial length of 100 years.

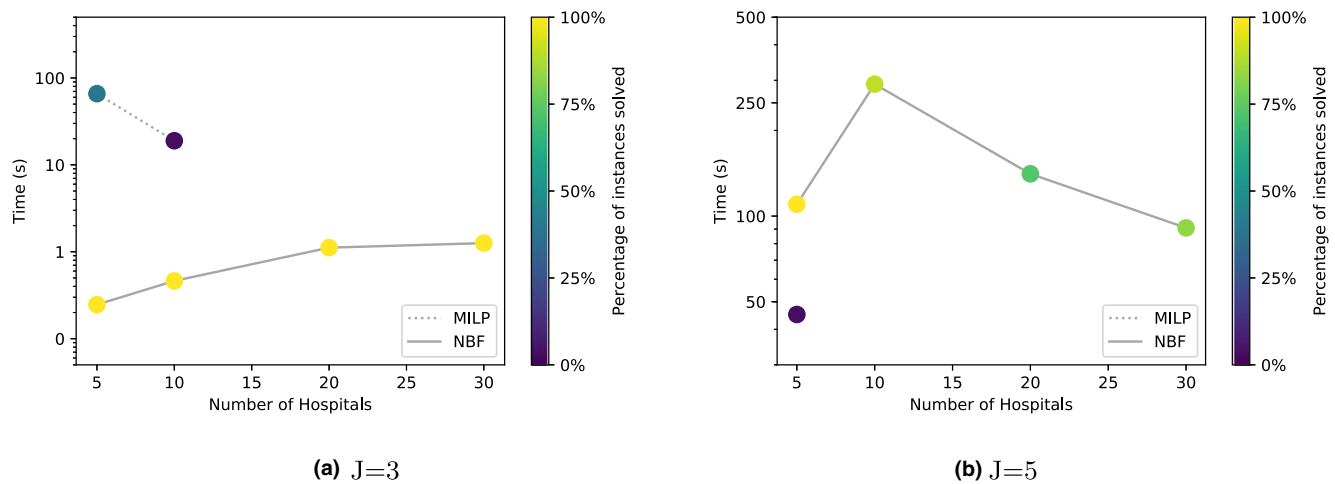
5.1. Performance Analysis

We first compare solution times and the number of instances solved between the MILP representation of AUC's RASP and the NBF approach over a wide range of problem parameters. For a planning horizon of $T = 52$ weeks, 3240 test instances were generated by varying the cohort size (10, 25, 50, 100, 200, 300), the number of hospitals (5, 10, 20, 30), the proportion of hospitals located in unique regions (20%, 40%, 60%), the percentage of fixed- versus variable-cost hospitals (25%, 50%, 75%), the number of rotations (3, 5), the rotation lengths in weeks (3, 5, 7), and the maximum idle time between rotations in weeks (3, 5, 7).

Detailed results that carefully outline the parameters selected for each experiment are presented in Appendix A. Figure 2 summarizes the data from all of the experiments and gives an overview of its implications. Specifically, the NBF approach solves 2763 (85%) of the test instances to optimality in under 30 minutes, while only 109 (4%) are solved by the MILP model. Note that all optimal MILP solutions match with the optimal NBF solutions. Further, of the problems solved by the NBF, 90% can be completed in under 5 minutes. These results suggest that for any choice of input parameters, the NBF approach significantly outperforms the MILP model by efficiently solving more problem instances in less time.

To further illustrate the differences in performance between the time-indexed MILP model and the NBF approach, we consider a collection of experiments with a planning horizon of $T = 52$ weeks. We hold fixed the proportion of hospitals located in unique regions, the percentage of fixed- versus variable-cost hospitals, and the maximum idle time between rotations at 60%, 50%, and 5 weeks, respectively. In the first set of experiments (Figure 3) we vary the number of hospitals (5, 10, 20, 30) while holding the number of rotations fixed ($J = 3$ in the left plot and $J = 5$ in the right). The shade of each data point represents the proportion of problems solved from thirty randomly generated instances; five randomly sampled rotation lengths (3, 5, 7) for each of the six different cohort sizes (10, 25, 50, 100, 200, 300). For $J = 3$, the NBF solves all problem instances to optimality within 30 minutes regardless of the number of available hospitals. The MILP model, on the other hand, does not solve a single instance when there are more than 10 hospitals. Further, when there are fewer than 10 hospitals, the solution speed of the NBF is orders of magnitude smaller than the MILP. For $J = 5$, the NBF once again exhibits much better performance. Whereas the MILP model can only solve a single instance with five hospitals, the NBF approach solves all problems instances. Further, when there are thirty hospitals, upwards of 60% of all randomly generated instances can be solved to optimality in under 2 minutes. Since the number of hospitals affiliated with AUC is 25, these results suggest that the NBF is more appropriate for realistically-sized problem instances.

In the second set of experiments (Figure 4), we vary the cohort size (10, 25, 50, 100, 200, 300) while holding the number of rotations fixed ($J = 3$ in the left plot and $J = 5$ in the right). The shade of each data point now represents the number of problems solved from 20 randomly generated instances; five randomly sampled rotation lengths (3, 5, 7) for each of the four values associated with the number of available hospitals (5, 10, 20, 30). We observe similar performance as in Figure 3 in that for $J = 3$, the NBF solves all problem

Figure 3 Time (Seconds) and the Percentage of Instances Solved as a Function of the Number of Hospitals [Color figure can be viewed at wileyonlinelibrary.com]

Notes: The proportion of hospitals located in unique regions, the percentage of fixed- versus variable-cost hospitals, and the maximum idle time between rotations are held fixed at 60%, 50%, and 5 weeks, respectively. The shade of each data point represents the number of solved instances from five randomly sampled rotation lengths (3, 5, 7) for each cohort size (10, 25, 50, 100, 200, 300).

instances to optimality within 30 minutes. The MILP model does not solve all instances even when there are fewer than 50 students per cohort. For $J = 5$, the MILP model can only solve one problem instance with 10 students. In contrast, the NBF optimally solves all instances with up to 100 students, and upwards of 50% of the instances when there are 200 or more students in a cohort. This further demonstrates that the NBF is a more robust alternative for AUC, as the cohort size over the last several years has consistently exceeded 300 students.

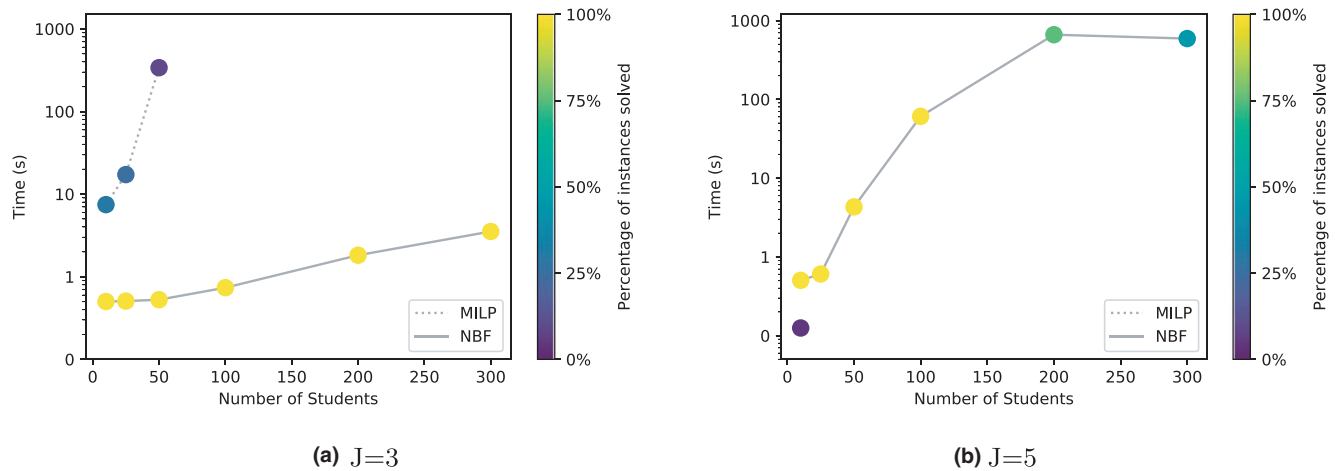
In the final set of experiments, we evaluate the impact of side constraints on MDD size. In Figure 5, we plot the number of nodes of the non-reduced MDD (x -axis) versus the number of nodes of the reduced MDD (y -axis) for all 3240 test instances. On average, the reduced MDD is typically half the size of the non-reduced MDD (6715 vs. 3691 nodes). This result remains consistent even as the number of rotations increases. For example, for $J = 3, 4$, and 5, there are approximately 1.82 times the number of nodes in the non-reduced MDD as compared to the reduced MDD. Finally, while the inclusion of QoS constraints leads to an MDD with a theoretical bound that is significantly larger than the original MDD, we find that the actual difference varies considerably. In fact, the addition of QoS constraints often results in substantial reductions in the number of nodes in D .

The performance analysis shows the effectiveness of the NBF approach at solving realistically-sized problem instances. As the number of hospitals and the size of a cohort increase, so do the number of QoS requirements that ensure students have an

excellent clerkship experience. Whereas any MILP model would naturally get larger to include more students (i.e., decision variables) and extra constraints for the additional hospitals and QoS requirements, the graphical structure associated with the NBF approach may, in some cases, shrink. This implies that as long as new constraints do not add to the MILP reformulation (i.e., section 3.5), AUC may see computational speedups for future cohorts as the scheduling requirements get more complex.

5.2. Counterfactual Study: What Could Have Been?

In this section, we compare AUC's historical scheduling decisions to what would have happened had our framework been implemented. We collected data on 1962 medical students who studied at AUC between August 2012 and January 2018. Each entry in the dataset included a unique (masked) student identifier, the rotation, the name of the hospital where the rotation was completed, and the start and end dates of the rotation. We also obtained a set of administrative documents outlining the contract details for each of the 25 affiliated hospitals. For each hospital, we examined the contract to determine the type of rotations offered, the fixed (c_k^f) and variable costs (c_{jk}) associated with training, the number of students that started a rotation on a given date (C_{jkt}), and the hospital's address (r_k). Hospitals within a 50 mile radius of each other were classified as being in the same geographic region. In accordance with the Council for Graduate Medical Education (2017), each student completed $J = 5$

Figure 4 Time (Seconds) and the Percentage of Instances Solved as a Function of the Cohort Size [Color figure can be viewed at wileyonlinelibrary.com]

Notes: The proportion of hospitals located in unique regions, the percentage of fixed- versus variable-cost hospitals, and the maximum idle time between rotations are held fixed at 60%, 50%, and 5 weeks, respectively. The shade of each data point represents the number of solved instances from five randomly sampled rotation lengths (3, 5, 7) for each number of hospitals (5, 10, 20, 30).

clinical rotations: internal medicine, general surgery, pediatrics, obstetrics & gynecology, and psychiatry. The first two rotations were twelve weeks in length, while the latter three were six weeks. We note that, for privacy reasons, neither dataset indicated whether a student failed their licensing examination nor how many times the examination was written.

The counterfactual study was performed by comparing the yearly placement costs of three scheduling approaches using past capacity and student availability data. The first approach (current practice) represents what actually transpired. The second (Greedy) approach schedules students to rotations according to the greedy algorithm. The third approach uses the NBF technique. For each scheduling model starting with the September 2014 cohort, students were assigned to rotations after passing their licensing exam. For the Greedy and NBF approaches, the placement costs of years 2015 and 2016 were obtained and compared to current practice (see Figure 6). Two important insights emerged. First, substantial savings, approximately 19% per year on average, could have been realized using the NBF technique. Although some savings are due to the better use of capacity, the majority of the cost reduction is associated with limiting the number of fixed-cost hospitals under contract. For example, in 2015, NBF needed to negotiate far fewer contracts with fixed-contract hospitals as compared to current practice. Second, we observed that the greedy approach bounds AUC's placement cost from below. Although, as can be expected, the greedy approach does not fully capture the idiosyncrasies of current scheduling practices, the result suggests that it is an appropriate benchmark.

5.3. Simulation Results: Evaluation of Future Costs

Using results from a DES, we now provide some insight as to the cost savings AUC can expect in the future. We use a discrete uniform distribution with support [300, 350] to model the number of students in a cohort. Although we do not have enough data to find the true cohort size distribution, discussions with staff indicate that this distribution and range is reasonable. We note that the mean cohort size estimated from the data set is 327. Given the realization for the number of students in a cohort for a given year, the empirical distribution models the week each student receives their examination results. Since a student may pass or fail the exam, the probability they begin clerkship is geometrically distributed with a success probability of 96% (one minus the failure rate). This is the average pass rate of first-time test takers for 2016. To ensure that AUC's QoS constraints are satisfied (see section 4), a student that begins clerkship cannot be idle for more than 8 weeks ($\epsilon = 8$). Further, at most 25% of students in a cohort can be assigned to two geographic regions ($\bar{r} = 1$) while all other students must complete in-hospital training in a single geographic region.

Since clerkship typically comprises the last 2 years of a four-year degree in medicine, we assume a rolling planning horizon of $T = 104$ weeks (i.e., 2 years with 52 weeks per year). If the week represents the beginning of a new school year, the size and examination weeks of the cohort are realized. Then, the scheduling policy (either greedy or the NBF) assigns students a schedule. On all other weeks, we realize which students can begin in-hospital training (by following their previously scheduled sequence of rotations) and

Figure 5 The Total Number of Nodes between Reduced and Non-Reduced Multivalued Decision Diagrams for All Instances [Color figure can be viewed at wileyonlinelibrary.com]

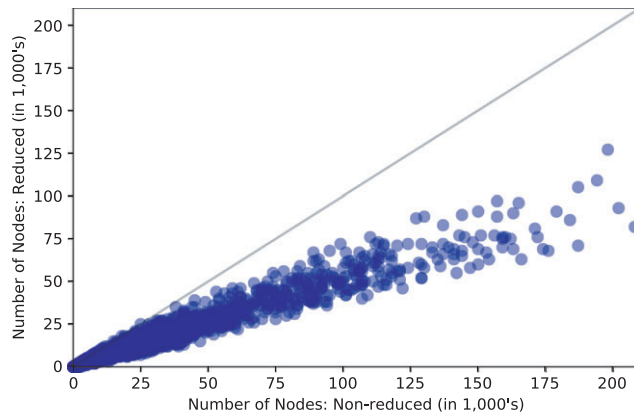
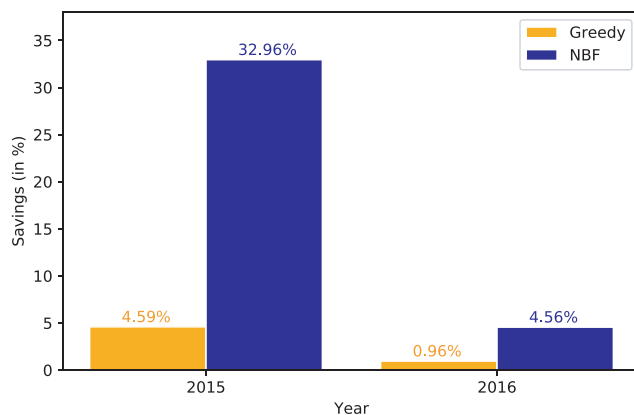


Figure 6 Comparison Against Historical Placement Costs (Current Practice) as a Function of Time [Color figure can be viewed at wileyonlinelibrary.com]



which students have failed their examination. For those that have failed, we re-assign them to a new schedule either using the greedy approach or by following Proposition 7. We also continuously track and update the capacity levels of each rotation as the current cohort affects the number of rotation slots available for next year's cohort. We note that in the DES, no time limit was imposed on the NBF although all instances were solved to optimality within 2 hours.

We find that, on average, AUC can expect around a 25% reduction in cost using the NBF approach. The translates into a real-dollar cost reduction of several hundred thousand dollars, a managerially significant quantity that validates the pursuit of a methodology that finds optimal assignment policies. Figure 7 presents a time-series plot of costs per year. In the left figure, costs associated with the greedy approach are presented, and in the right figure, the costs accrued by the NBF are given. Notice that the greedy

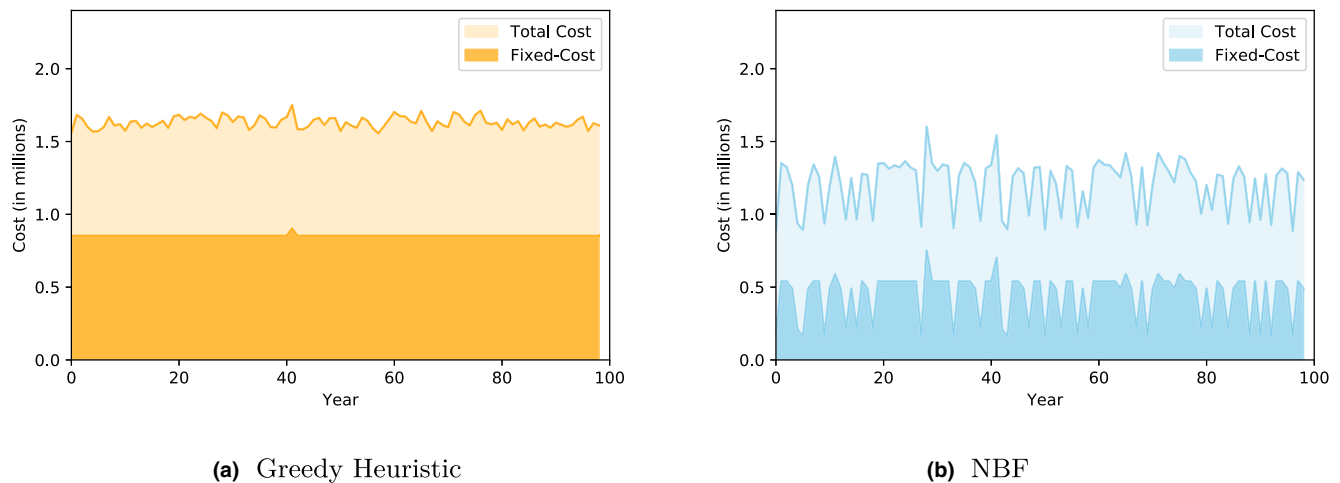
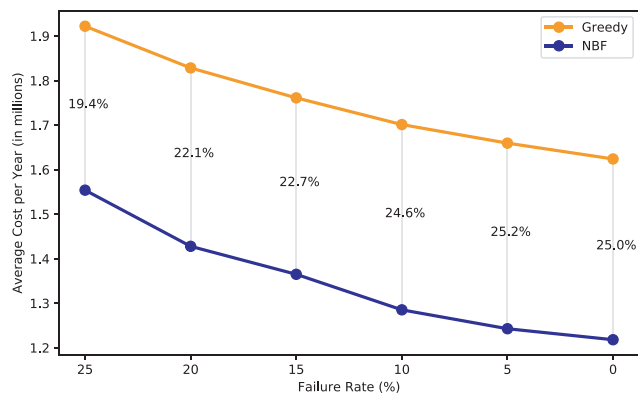
approach is more consistent than the NBF as the year-over-year costs are smaller. However, the NBF always assigns students to schedules with a smaller placement cost, and in some cases, finds assignments that are half as costly. This is because the NBF consistently opens fewer fixed-cost hospitals as compared to the greedy heuristic. This intuition is confirmed by examining the proportion of the total cost that can be attributed to fixed-cost hospitals in Figure 7. We also observe that the greedy approach tends to concentrate the use of fixed-cost hospitals early in the planning horizon, whereas the NBF technique tends to use these hospitals more evenly over time.

5.4. Simulation Results: Estimation and Forecast Uncertainty

Our results thus far assume that the licensing examination failure rate at AUC is equal to the national average and that the average number of students in a cohort is fixed as compared to the number of affiliated hospitals. However, the actual failure may be higher (or lower) and the ratio between the number of students to the number of available rotation slots may not be held fixed. Thus, we now investigate how sensitive our results are to mis-specifications in the underlying stochastic quantities. This gives AUC further insight into the value of a more systematic approach to student scheduling and also validates the choice of a deterministic model (as opposed to a stochastic program, for instance).

In Figure 8, we examine how the failure rate affects costs. As expected, we find that as the failure rate increases, the placement cost AUC incurs also increases. This is independent of the scheduling approach and highlights the fact that as uncertainty plays a larger role in the scheduling environment, including stochasticity in any model becomes compulsory. Nonetheless, since the failure rate at AUC has been approximately equal to the national average over the last several years (AUC Medical School 2018a), our approach is reasonable. We also find that the NBF reduces placement costs by approximately 20%–25% as compared to current practice, and is robust with respect to a wide range of failure rate values (see Figure 8). Thus, the NBF approach yields a substantial reduction in placement costs and, as AUC acts to increase the licensing examination pass rate (AUC Medical School 2018b), the resulting cost savings will approach the true minimum.

In Figure 9, we study how the size of the cohort, given that the number of available rotation slots are held fixed, affects placement costs. For both approaches, we find that costs increase as a function of the cohort size, as expected. However, we also find that the NBF is clearly superior when the supply of rotation slots does not substantially exceed demand.

Figure 7 Placement Cost Per Year (in Millions) as a Function of Time (Year) for a Single Trial [Color figure can be viewed at wileyonlinelibrary.com]**Figure 8** Average Placement Cost Per Year (in Millions) as a Function of the Per Student Failure Rate [Color figure can be viewed at wileyonlinelibrary.com]

Notes: Note that the probability a student can begin clerkship is geometrically distributed with success probability equal to 1 minus the failure rate. The average values over 100 trials are presented.

More specifically, for small cohorts, any scheduling approach is sufficient; there are few students and many feasible schedules to choose from. In the middle range where the cohort size is between 150 and 250 (more realistic) and the supply of rotation slots is still ample, the difference between the NBF and the greedy approach is the largest. In this region, there are many feasible schedules although choosing the optimal one is computationally difficult. This is why the NBF yields such large cost improvements as compared to the greedy approach. AUC typically operates in the upper range, where demand for rotation slots is large and supply is limited. In this region, the relative cost savings are reduced (as compared to the middle range) as the ratio of demand to supply approaches one. Nevertheless, the absolute cost savings in this

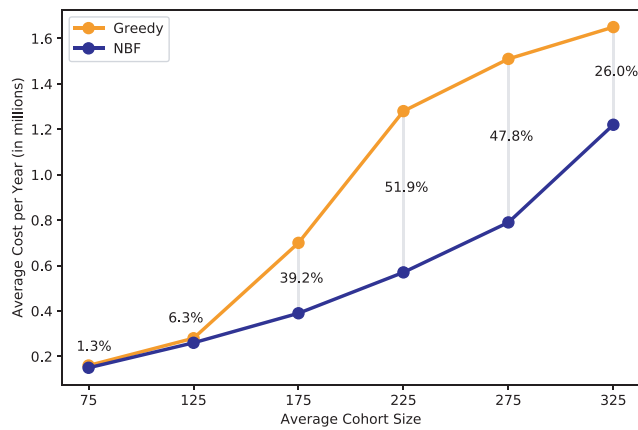
region are the most substantial due to the large number of students that must be placed. Further, because AUC typically negotiates fixed contracts with hospitals before assigning students to rotations, supply will never far outpace demand. Since it is likely that AUC will continue to operate in this region for the foreseeable future, results such as these are encouraging.

6. Managerial Insights

In addition to its computational benefits, the underlying graphical structure in our NBF approach has useful managerial implications. For example, by creating a pattern store (a collection of feasible paths in D), administrators can ensure that specific schedules are either not assigned or are more likely to be assigned to students in future cohorts. These individual patterns can be inserted into or removed from the network using MDD-based separation techniques. We refer to Cire and Hooker (2014) and Bergman et al. (2016) for MDD methodologies to include and remove paths alongside their associated computational complexities.

This flexibility allows the medical school to utilize feedback solicited from students who have finished clerkship. Negative experiences could be directly incorporated into the assignment process, that is, particularly problematic schedules would be eliminated by removing the corresponding paths in the MDD. Systematically representing these idiosyncratic restrictions as constraints in an MILP model would be more cumbersome to implement and difficult to maintain. Moreover, while, in theory, it is possible to answer some of these questions with an MILP model, we did not observe its performance to be robust on larger problem sizes. Specifically, while the number of hospitals and students of a real-life problem

Figure 9 Average Placement Cost Per Year (in Millions) as a Function of the Mean Number of Students in a Cohort [Color figure can be viewed at wileyonlinelibrary.com]



Notes: Note that the probability a student can begin clerkship on a given week is sampled from the empirical distribution. The average values over 100 trials are presented.

instance are around 25 and 200–300, respectively, the current state-of-the-art MILP model cannot solve instances with more than 10 hospitals and 50 students.

Our approach also allows a medical school to conduct analyses for capacity planning, demand management, and contract negotiations. For example, the maximum number of students that can be assigned to any hospital can be obtained by replacing the objective function of the NBF to maximize flow (instead of minimizing cost). Note that for the classical time-indexed MILP, this would require adding extra students until the model became infeasible. Further, by adjusting the cost of the arcs in the decision diagram, a medical school can determine a threshold beyond which it would be too costly to assign students to any affiliated hospital. These managerial questions are of interest to any medical school in the planning stages, particularly when determining how many students should receive an offer of admission. Although additional MILP constructs could be incorporated to address these concerns, these insights are already obtainable as a direct consequence of the NBF.

It is this versatility that AUC finds compelling about the network-based model. Interviews with scheduling staff at AUC have indicated that our approach can be used to answer other managerially relevant questions, as well as to analyze certain what-if scenarios. Consider the following examples:

1. *What is the maximum number of slots per rotation that AUC can effectively use from a hospital?* By maximizing the flow over the arcs associated with that hospital in the decision diagram, and setting the capacity of the remaining arcs

and hospitals to infinity, the number of slots can be obtained. Such a solution would also account for the underlying complex combinatorial structure of the problem imposed by the operational requirements. For example, suppose that a hospital only offers a single rotation j , but there exists some rotation $j' \neq j$ that is only available at a fixed-cost hospital which also offers rotation j . Because of the *all-or-none* requirement, the number of slots requested from the first hospital should effectively be zero.

2. *What if a certain hospital were to dramatically increase (or decrease) the number of slots available to a cohort?* This can be answered by adjusting the arc and hospital capacities accordingly. For a bound on the marginal cost associated with increasing or decreasing capacities, it is also possible to use the reduced cost of the arc variables in the NBF. Such bounds will be more accurate than the time-indexed MILP formulation, since the network flow model represents an ideal formulation of the set of student schedules (i.e., without operational requirements, the extreme points of the network flow model are integral and represent a feasible student schedule).
3. *How different are schedules that move students through clerkship the quickest versus those that keep them within a small geographic region?* If particular arcs corresponding to schedules completed at the end of the planning horizon are removed from the decision diagram, it forces students to complete clerkship earlier. The decision diagram can then be modified to only allow transitions between regions that are geographically close.
4. *How many fixed-cost hospitals should AUC secure contracts with?* In general, guaranteeing rotation slots far in advance requires significant investment as expensive contracts with fixed-cost hospitals must be signed. By taking better advantage of the capacity available at variable-cost hospitals, which offer a smaller cost per placement, AUC can reduce their overall cost. Further, if too much capacity is concentrated in a small number of fixed-cost hospitals, a large proportion of students must have their schedules restricted to the dates offered by those hospitals, limiting the universe of feasible schedules that can be created.

AUC values their current relationships with fixed-cost hospitals because they can be trusted to provide a steady number of rotation slots every year (thereby reducing uncertainty). In contrast, capacity at variable-cost hospitals is

more volatile. Therefore, it is not clear whether AUC has any interest in reducing the number of contracts they secure. Nevertheless, our model can be used by AUC, and any medical school for that matter, to better negotiate contracts by means of what-if scenarios. The model can be run several times, with different inputs for hospital capacities, costs, etc., and the resulting outputs would form the basis of a contract negotiation strategy in which the school argues for specific capacity levels, price adjustments, and so on.

Many other questions can be similarly answered by running and re-running our model after fixing certain variables (or arcs of the decision diagram) to specific values. In fact, Alejandro Carrasco, Senior Director of Business Operations, remarked that “being able to quickly answer these types of questions would significantly improve our operations, allowing us to deliver better quality schedules for our students, and with much less manual effort.”

7. Concluding Remarks

We investigate the scheduling practices of a medical school that must assign a cohort of medical students, who become eligible over a period of several months, to a series of clinical rotations at hospitals in different geographic regions. The student schedules must obey certain core restrictions, licensing requirements, and quality-of-service constraints, while keeping the school's placement costs to a minimum. To this end, we propose a network-based formulation (NBF) and show that it is significantly superior to a mixed-integer linear programming (MILP) model. Using a dataset collected from the American University of the Caribbean (AUC), we perform a counterfactual analysis and develop a DES of the licensing examination and clerkship scheduling process. We find that AUC could have saved approximately 19% per year on average in student placement costs had our scheduling approach been implemented. Using the DES, we then compare our exact scheduling approach with current practice and demonstrate that, in the future, AUC can expect an average cost reduction of 25%. We also show that this cost reduction is robust to mis-specifications in the licensing examination failure rate and the size of a cohort when the number of affiliated hospitals is held constant.

Our NBF is created by using decision diagrams to reformulate a particularly large optimization problem as a network model with side constraints. We use dynamic programming, and its ability to reduce a complex problem to a series of simpler subproblems, together with MILP modeling, which more easily represents complicated relationships between decision

variables and best communicates various problem restrictions. By carefully choosing which constraints to incorporate into the decision diagram versus which to represent in the MILP model, one can obtain significant computational speed-up and exact solutions to large-scale problems that may be difficult to implement and fine tune using other methods such as branch-and-price. Future work will aim to provide more guidance as to what classes of constraints are better represented in the state-transition graph versus those that are better suited for a mathematical programming representation.

Acknowledgments

The first two authors were supported by a Discovery Grant provided by the Natural Sciences and Engineering Research Council of Canada (NSERC). The authors are grateful to the department editor, the associate editor, and the three anonymous referees, whose comments and suggestions have considerably improved the paper. We also sincerely thank the staff from the School of Medicine at the American University of the Caribbean for the support and time dedicated to this project.

References

- Agarwal, A. 2016. Balancing medical resident education and workload while ensuring quality patient care. Ph.D. thesis, Rochester Institute of Technology.
- Artigues, C. 2017. On the strength of time-indexed formulations for the resource-constrained project scheduling problem. *Oper. Res. Lett.* **45**(2), 154–159. Available at <http://www.sciencedirect.com/science/article/pii/S016763771730072X>.
- Association of American Medical Colleges. 2016. Total enrolment by US medical school. Available at <https://www.aamc.org/download/321526/data/factstableb1-2.pdf>, <https://www.aamc.org/download/321532/data/factstableb2-2.pdf> (accessed date October 25, 2017).
- AUC. 2018. American university of the caribbean medical school. Available at <https://www.aucmed.edu> (accessed date March 29, 2018).
- AUC Medical School. 2018a. Facts & figures about AUC's medical school. Available at <https://www.aucmed.edu/about/facts-and-figures.html> (accessed date February 18, 2018).
- AUC Medical School. 2018b. Facts & figures about AUC's medical school. Available at <https://www.aucmed.edu/academics/usmle-preparation.html> (accessed date February 18, 2018).
- Baptiste, P., C. Le Pape, W. Nuijten. 2001. *Constraint-Based Scheduling*. Kluwer Academic Publishers, New York, NY.
- Bard, J. F., Y. Shao, X. Qi, A. I. Jarrah. 2014. The traveling therapist scheduling problem. *IIE Trans.* **46**(7): 683–706.
- Bard, J. F., Z. Shu, D. J. Morrice, L. K. Leykum, R. Poursani. 2016. Annual block scheduling for family medicine residency programs with continuity clinic considerations. *IIE Trans.* **48**(9): 797–811.
- Barnhart, C., E. L. Johnson, G. L. Nemhauser, M. W. P. Savelsbergh, P. H. Vance. 1998. Branch-and-price: Column generation for solving huge integer programs. *Oper. Res.* **46**(3): 316–329.
- Beasley, J. E., B. Cao. 1998. A dynamic programming based algorithm for the crew scheduling problem. *Comput. Oper. Res.* **25**(7–8): 567–582.

- Bergman, D., A. A. Cire, W.-J. van Hoeve, J. N. Hooker. 2016. *Decision Diagrams for Optimization*. Springer International Publishing, New York, NY.
- Burke, E. K., P. De Causmaecker, G. Vanden Berghe, H. Van Landeghem. 2004. The state of the art of nurse rostering. *J. Sched.* 7(6): 441–499.
- Cappanera, P., M. G. Scutellà. 2014. Joint assignment, scheduling, and routing models to home care optimization: A pattern-based approach. *Transport. Sci.* 49(4): 830–852.
- Cheang, B., H. Li, A. Lim, B. Rodrigues. 2003. Nurse rostering problems—A bibliographic survey. *Eur. J. Oper. Res.* 151(3): 447–460.
- Cire, A. A., J. N. Hooker. 2014. The separation problem on binary decision diagrams. *Proceedings of the International Symposium on Artificial Intelligence and Mathematics (ISAIM)*, Fort Lauderdale, FL, USA.
- Clautiaux, F., C. Alves, J. Valério de Carvalho, J. Rietz. 2011. New stabilization procedures for the cutting stock problem. *J. Comput.* 23(4): 530–545.
- Cohn, A. M., S. Root, C. Kymissis, J. Esses, N. Westmoreland. 2009. Scheduling medical residents at Boston University school of medicine. *Interfaces* 39(3): 186–195.
- Conforti, M., G. Cornuéjols, G. Zambelli. 2010. Extended formulations in combinatorial optimization. *4OR* 8(1): 1–48.
- Côté, M. C., B. Gendron, L. M. Rousseau. 2011. Grammar-based integer programming models for multiactivity shift scheduling. *Management Sci.* 57(1): 151–163.
- Council for Graduate Medical Education. 2017. Program requirements. Available at <http://www.acgme.org/What-We-Do/Accreditation/Common-Program-Requirements> (accessed date June 7, 2017).
- Eppen, G. D., R. K. Martin. 1987. Solving multi-item capacitated lot-sizing problems using variable redefinition. *Oper. Res.* 35(6): 832–848.
- Ernst, A. T., H. Jiang, M. Krishnamoorthy, D. Sier. 2004. Staff scheduling and rostering: A review of applications, methods and models. *Eur. J. Oper. Res.* 153(1): 3–27.
- Federation of State Medical Boards and National Board of Medical Examiners. 2018. The united states medical licensing examination. Available at <http://www.usmle.org> (accessed date March 29, 2018).
- Garey, M. R., D. S. Johnson, R. Sethi. 1976. The complexity of flowshop and jobshop scheduling. *Math. Oper. Res.* 1(2): 117–129.
- Guo, J., D. R. Morrison, S. H. Jacobson, J. A. Jokela. 2014. Complexity results for the basic residency scheduling problem. *J. Sched.* 17(3): 211–223.
- Gurobi Optimization, Inc. 2016. Gurobi optimizer reference manual. Available at <http://www.gurobi.com>.
- Hall, R., ed. 2012. *Handbook of Healthcare System Scheduling, International Series in Operations Research & Management Science*, vol. 168. Springer, New York, NY.
- Hooker, J. N. 2013. Decision diagrams and dynamic programming. C. Gomes & M. Sellmann, eds. *Proceedings of CPAIOR 2013*, vol. 7874. Springer Berlin Heidelberg, Yorktown Heights, NY, pp. 94–110.
- Joncour, C., S. Michel, R. Sadykov, D. Sverdlov, F. Vanderbeck. 2010. Column Generation based Primal Heuristics. *International Symposium on Combinatorial Optimization (ISCO)*, *Electron Notes Discrete Math.*, vol. 36, 695–702.
- Lancia, G., P. Serafini. 2018. *Compact Extended Linear Programming Models*. Springer, Cham, Switzerland.
- Lemay, B., A. Cohn, M. Epelman, S. Gorga. 2017. New methods for resolving conflicting requests with examples from medical residency scheduling. *Prod. Oper. Manag.* 26(9): 1778–1793.
- Ludmerer, K. M. 2009. Resident burnout: Working hours or working conditions? *J. Grad. Med. Educ.* 1(2): 169–171.
- Nickel, S., M. Schröder, J. Steeg. 2012. Mid-term and short-term planning support for home health care services. *Eur. J. Oper. Res.* 219(3): 574–587.
- Pinedo, M. L. 2008. *Scheduling: Theory, Algorithms, and Systems*, 3rd edn. Springer, New York, NY.
- Raffensperger, J. F. 1999. The marriage of dynamic programming and integer programming. *Proceedings of the 34th Annual Conference of the Operational Research Society of New Zealand (ORSNZ)*, Hamilton, New Zealand, pp. 49–58.
- Restrepo, M. I., L. Lozano, A. L. Medaglia. 2012. Constrained network-based column generation for the multi-activity shift scheduling problem. *Int. J. Prod. Econ.* 140(1): 466–472.
- Smalley, H. K., P. Keskinocak. 2016. Automated medical resident rotation and shift scheduling to ensure quality resident education and patient care. *Health Care Manage. Sci.* 19(1): 66–88.
- Smet, P., P. Brucker, P. De Causmaecker, G. Vanden Berghe. 2016. Polynomially solvable personnel rostering problems. *Eur. J. Oper. Res.* 249(1): 67–75.
- U.S. Medical Licensing Examination. 2015. Performance data. Available at <http://www.usmle.org/performance-data/default.aspx> (accessed date August 16, 2016).
- Van den Akker, J. M., C. A. J. Hurkens, M. W. P. Savelsbergh. 2000. Time-indexed formulations for machine scheduling problems: Column generation. *J. Comput.* 12(2): 111–124.
- Van den Bergh, J., J. Beliën, P. De Bruecker, E. Demeulemeester, L. De Boeck. 2013. Personnel scheduling: A literature review. *Eur. J. Oper. Res.* 226(3): 367–385.
- Vanderbeck, F. 2011. Branching in branch-and-price: A generic scheme. *Math. Program.* 130(2): 249–294.
- Warner, D. M. 1976. Scheduling nursing personnel according to nursing preference: A mathematical programming approach. *Oper. Res.* 24(5): 842–856.
- Wegener, I. 2000. *Branching Programs and Binary Decision Diagrams: Theory and Applications*. SIAM monographs on discrete mathematics and applications, Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Wimmers, P. F., T. A. W. Splinter, G. R. Hancock, H. G. Schmidt. 2007. Clinical competence: General ability or case-specific? *Adv. Health Sci. Educ.* 12(3): 299–314.
- Zhang, H., L. Jackson, A. Tako, J. Liu, S. Dunnett. 2016. Network based formulations for roster scheduling problems. E. K. Burke, L. Di Gaspero, E. Özcan, B. McCollum, A. Schaerf, eds. *Proceedings of the 11th International Conference on the Practice and Theory of Automated Timetabling (PATAT)*, Udine, Italy, August 23–26, 2016, pp. 403–419. Available at <http://www.patatconference.org/patat2016/proceedings/> (accessed date December 09, 2019).

Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

Appendix A: Performance Analysis Results.

Appendix B: Proofs of Lemmas and Propositions.