# Controlling understaffing with conditional Value-at-Risk constraint for an integrated nurse scheduling problem under patient demand uncertainty

Fang He[a,*], Thierry Chaussalet[a], Rong Qu[b]

[a] School of Computer Science and Engineering, University of Westminster, W1W 6UW, UK
[b] The Automated Scheduling, Optimisation and Planning (ASAP) Group, School of Computer Science, The University of Nottingham, Nottingham, NG8 1BB, UK

## ABSTRACT

Nursing workforce management is a challenging decision-making task in hospitals. The decisions are made across different timescales and levels from strategic long-term staffing budget to mid-term scheduling. These decisions are interconnected and impact each other, therefore are best taken by considering staffing and scheduling together. Moreover, this decision-making needs to be made in a stochastic setting to meet uncertain patient demand. A sufficient and cost-efficient staffing level with desirable schedule is essential to provide good working conditions for nurses and consequently good quality of care. On the other hand, understaffing can severely deteriorate the quality of care thus should be strictly controlled.

To help with the decision making, based on our previous research we formulate in this paper an integrated nurse staffing and scheduling model under patient demand uncertainty into a two-stage stochastic programming model with an emphasis on understaffing risk control. Conditional Value-at-Risk (CVaR), a risk control measure primarily used in the financial domain, is integrated in the stochastic programming model to control understaffing risk. The IBM ILOG CPLEX solver is applied to solve the stochastic model. The model and solution approaches are tested using a case study in a real-world environment setting. We have evaluated the performance of the stochastic model and the benefit of CVaR in terms of impact on schedule quality.

## 1. Introduction

The increasing patient demand in healthcare raises challenges for hospitals from many perspectives. As nurse labour costs typically represent a large share of the total hospital budget [3], hospitals need to manage and deploy their human resources efficiently. Overtime workload, undesired work patterns, and low satisfaction are widely known issues among nurses [27]. Managing personnel cost, reducing overtime workload and undesired work patterns and improving work satisfaction efficiently all have positive impact on the quality of care provided for patients and the cost efficiency of hospitals [21]. One potential way to address these issues is to develop and analyse models and decision support systems to gain insight into the outcomes and consequences of various nurse workforce management strategies.

The management of nurse workforce is extremely challenging due to the fact that it is typically made across different time horizons and different organisational levels [16]. It is a multi-phase planning and control process that consists of staffing, shift scheduling and allocation phases [16]. Staffing is a strategic long-term planning decision that determines the mix of nursing resources. Shift scheduling focuses on the

assignment of available nurses to shifts and then constructs a mid-term roster. The roster needs to strictly meet regulations and policies and satisfy staff's personal preferences as much as possible. The regulations and policies restrict the acceptable scheduling patterns that nurses can work on, and consequently a different mix of nurse resources may be required. The linkage between these two phases suggests a more integrated approach, which motivates the model proposed in this paper.

Decision-making on nurse staffing and scheduling becomes more challenging when uncertainty is considered, which almost always presents in realistic scenarios. Many hospitals are subject to regulations to guarantee a certain level of nurses to ensure the quality of care provided [23,26,27]. The patient demand determines the number of nurses required for each shift on each day of the week. This demand fluctuation has direct impact on the number of nurses required, i.e. on labour cost (e.g. by hiring agent nurses or overtime nurses when understaffing occurs) as well as on the attractiveness of the schedule (e.g. by increasing overtime shifts etc.). Understaffing needs to be addressed as it can severely deteriorate the quality of care [23,26].

Patient demand uncertainty needs to be taken into consideration to ensure an efficient and flexible schedule. To account for these issues, we

* Corresponding author.
  *E-mail addresses:* hef@westminster.ac.uk (F. He), chausst@westminster.ac.uk (T. Chaussalet), Rong.Qu@nottingham.ac.uk (R. Qu).

propose a two-stage stochastic programming (SP) model for the integrated nurse staffing and scheduling problem under patient demand uncertainty. In the first stage, the initial staffing level and schedule is decided to minimise the labour cost, overtime workload and unattractive shift patterns. The second stage then adjusts the schedule under different demand scenarios by introducing /removing nurse shifts from/into a centre nurse pool. We propose to apply Conditional Value-at-Risk (CVaR) to control the risk of understaffing, i.e. to keep the under-staff number within a certain confidence level to ensure an adequate number of nurses.

CVaR is primarily used in finance as a risk measure to control the loss within certain confidence levels [22]. In has been introduced in healthcare operational research in the recent literature. Najjarbashi and Lim [17] use CVaR to reduce the variability of operating room scheduling under uncertainty by reducing the worst-case outcomes of an operating room schedule. Their MILP model formulation is based on a finite set of scenarios generated using the Monte Carlo sampling method. Kishimoto and Yamashita [15] apply an LP approach with CVaR type constraints for intensity modulated radiotherapy treatment (IMRT) optimisation. A key clinical criterion that measures the quality of an IMRT plan is to satisfy dose-volume constraints (DVCs), which is an NP-hard problem. The CVaR type constraints, which always satisfy the DVCs, can be described as linear constraints; therefore, the optimisation problem is transferred to an LP problem, which is much easer to solve.

One of the aims in this research is to investigate the nurse scheduling problem under patient demand uncertainty. One way to react to this uncertainty is to build some robustness in the integrated planning and scheduling phases. As summarised in Jonas Ingels and Maenhout [11], several studies in the literature propose a reactive decision support model. This model adopts options, such as allocating overtime shifts, schedule changes and allocating cross-trained nurses, to match supply and demand. In contrast, a **proactive** approach is to build some mechanisms in the scheduling phase such that a robust roster is constructed. A common proactive approach is to include buffers, such as time buffers or capacity buffers. Time buffers, e.g. flexible shift length, have been applied in personnel scheduling. Capacity buffers, e.g. reserving duties, have mostly been studied in the airline industry. Another proactive approach is two-stage stochastic approach. The first stage constructs the baseline schedule by minimising the cost while satisfying a minimum staffing requirement. The second stage takes recourse actions to adjust the shifts to meet the requirements from different scenarios. Our proposed approach falls into this category. To have more detailed information on measure of robustness and cost of robustness, we refer to Jonas Ingels and Maenhout [11], Tam et al.'s [25] work.

Two stochastic programming models are proposed in the paper. The first Stochastic Demand Model (SDM) models the demand profiles (i.e. the required number of nurses for each shift on each day) as scenarios. The second is an SDM with an additional CVaR constraint (SDM-CVaR) to control the understaffing risk. Both models aim to optimise the labour cost, to improve work satisfaction, as well as to reduce overtime workload and undesired work patterns. A practical yet very efficient solution procedure with CPLEX solver is applied to solve these models.

The main contribution of the paper can be summarised as follows: 1) An integrated staffing and scheduling model under patient demand uncertainty is proposed to produce a more flexible schedule, which accounts for reducing labour cost and overtime workload, and unattractive work patterns. 2) Applying CVaR as a risk control measure for understaffing, aiming at sufficient staff level within desired confidence level.

The remainder of this paper is organised as follows. Section 2 presents a literature review on the problem and related research. In Section 3, we formulate the integrated models under stochastic demand. Section 4 presents the solution method. Section 5 presents a case study to evaluate and compare the performance of the models. Finally,

we draw our conclusions and present future work in Section 6.

## 2. Literature review

Various models have been proposed in the literature on nurse scheduling problems [27]. Early papers focussed on problem constraints. Van den Bergh et al. [27] classify the constraints into different categories such as coverage, time-related, fairness and balance constraints. These constraints can be treated as hard constraints (which must be satisfied) or soft constraints (which can be violated but usually associated with a penalty) to achieve flexibility of problem modelling and solving. Nowadays, the quality of a nurse roster is increasingly measured in terms of personal satisfaction [27]. Overtime workload, work patterns and job satisfaction are the key factors that are investigated to achieve a satisfactory roster.

To simplify the modelling of the problem, nursing workforce management has been divided into a multi-phase sequential process [9,16,28] with different time horizons and different management levels. Early research focussed on phase-specific problem modelling and solving methodologies [9,10,16,28].

Then some researchers realised that workforce management should not be considered in isolated phases because of the inter-relationship of staffing-size and scheduling, as well as the conflicting multiple objectives of minimising cost and maximising costumer service [14,16,19]. This line of research can be generally concluded as a two-step approach: it first determines the staffing levels required to meet the desired performance at low cost, and then generates the minimum cost shift schedules to meet these requirements. Dantzig's set covering formulation [7], dated back to the '50s, is still highly relevant and used frequently in this approach. In the first step, the staffing level requirement is interpreted as a strict constraint to be met in Dantzig's model. The constraints introduced in the second step are commonly related to working regulation and employee preferences. The two-step approach is appealing because it evades the difficulty of stochastic performance constraints in the mathematical models. With this approach, the performance constraints are taken care of in the staffing stage, so that shift scheduling becomes a deterministic problem. However, the two-step approach may lead to sub-optimal shift schedules [12]. Therefore, recent literature has started to focus on more integrated approaches.

Maenhout and Vanhoucke [16] propose a more compact integrated staffing and scheduling model for a long-term nurse management problem over multiple departments. It is a single aggregated model compared with the two-step approach described above. It shows that staffing multiple departments simultaneously and integrating nurse characteristics into the staffing decision can lead to substantial improvements in schedule quality. Wright and Mahar [30] tackle the staffing and scheduling problem and achieve reduced cost and improved nurse satisfaction by scheduling cross-trained nurses, which come from multiple departments in a centre nurse pool. Wright and Bretthauer [29] present coordinated decision-making models to coordinate nurses inside the hospital, and agent nurses outside the hospital to reduce labour cost, as well as overtime workload. The results show how centralised scheduling can be used to reduce cost and improve nurse satisfaction. However, all these studies assume a deterministic setting.

Healthcare systems, like many other service systems, are featured with non-stationary and uncertain demands: the number of patients/ customers fluctuates over time in a stochastic manner. Defraeye and Nieuwenhuyse [8] provide a state-of-the-art literature review on staffing and scheduling approaches that account for non-stationary demand, mainly focusing on applications in call centres and emergency departments. In healthcare systems, patient demand uncertainty is prominent. Most hospitals enforce a patient-to-nurse ratio. Therefore, uncertainty should be taken into account in the decision making to produce a flexible schedule. Zinouri [32] addresses staff scheduling problems through a demand prediction and scenario-based approach. In

his work, based on historical data, a time series forecasting method is applied to predict daily surgical case volume. Based on the prediction, a scenario set is generated for the staff-scheduling problem.

Stochastic programming is a well-developed method to model decision-making under uncertainty in a flexible way, which imposes real-world constraints relatively easily. Bard and Purnomo [5] consider the problem of short-term nurse rescheduling for daily fluctuation in patient demand, where a given mid-term schedule is revised to cover shortage. In Bagheri et al.'s [4] stochastic model, in addition to fluctuation in patient demand, uncertainty in patient stay period over time is also considered. Zhu and Sherali [31] present a two-stage stochastic workforce-planning model in which the second stage decisions assign continuous workload to each worker. Kim and Mehrotra [14] propose a two-stage stochastic integer programming model to the integrated staffing and scheduling problem, where the second stage decision variables are integer. They assume that all acceptable working schedule patterns are pre-generated; then a modified multi-cut aggregation in an integer L-shaped algorithm with a priority branching strategy is proposed to solve the model. In Bagheri et al. [4], a sample average approximation method is applied to obtain an optimal schedule with the minimum regular and overtime assignment cost.

The model proposed in this paper seeks to efficiently schedule nurses in a ward facing uncertain demand while simultaneously optimising the number of nurses assigned to the ward based on an initial staffing number. We formulate the problem as a two-stage stochastic program. Patient demand i.e. the required number of nurses for each day is modelled using scenarios that vary over time. Thus, overstaffing and understaffing may occur. In order to keep understaffing under certain level, a CVaR constraint, which is often used to measure uncertainty in finance, is utilised in the model.

## 3. Problem description and modelling

### 3.1. Problem description

In nursing workforce management systems at most hospitals, the staffing level in each department needs to be decided, based on which a schedule for the corresponding staffing level over a period of 4 weeks (usually) can be constructed under stochastic patient demand. This usually starts with an initial base line number of nurses available for each department within the hospital's budget. The scheduling policy, which is defined in terms of practices rules, needs to be tackled to construct a satisfactory schedule, therefore staffing and scheduling need to be simultaneously considered in the process.

Previous work [16,29,30] showed that coordination of staffing across different departments can improve the quality of decision-making. Nurses are typically assigned according to a fixed or cross-utilisation policy. The former policy states that a nurse is permanently assigned to a specific ward. The latter implies that a nurse who is a member of a centre pool can be referred to a different unit. The hospital in our study applies a mixed policy. That is, an initial base number of nurses within the department's budget are assigned to the department. However, a centre pool of nurses is maintained from where extra nurse shifts can be transferred from the pool to cover the shortage in certain department, or redundant nurse shifts can be transferred into the pool.

### 3.1.1. Objective function

The quality of a nurse staffing plan and scheduling should be measured from multiple perspectives as stated above for both hospital and nurses. In Maenhout and Vanhoucke [16], the quality of a nurse staffing and shift scheduling plan is measured using three dimensions representing the hospital's and nurses' objectives, i.e., the effectiveness in providing nursing care, the efficiency of a nursing unit and the job satisfaction among nursing staff. We adapt similar measurements in our objective function, explained as follows:

(1) Personnel cost: The personnel cost consists of regular payment and overtime payment. In practice, the salary scale of nurses varies according to their experience, length of employment and other factors. The regular payment and overtime payment are represented by their corresponding parameters. In this work overtime payment is 1.5 times of regular payment, and nurses' pay is doubled on bank holidays.

(2) The quality of a nurse roster in modern working environment is increasingly measured using personnel job satisfaction [27] including violations of balanced workload and individual preferences. This is captured in our model.

(3) The recourse cost is the cost of over-staffing and understaffing in the second-stage SP model.

The overall objective function is thus an integrated function of all the above costs.

### 3.1.2. Constraints

Nurse scheduling in hospitals involves many constraints including working regulations, legal requirements, and nurses' preferences, etc. The constraints concerned in this paper are derived from real-life scenarios in hospital wards and are mostly tested in benchmark problems in the literature. Rules and regulations have been directly taken from real-world cases and preserved with essential characteristics. The problem can have several variants with respect to the number of nurses, the number of shifts and the length of the scheduling period.

### 3.2. Problem modelling

### 3.2.1. Background on two-stage stochastic program

Stochastic programming is a well-developed optimisation method under uncertainty. Shapiro and Philpott [24] provide a very good introduction to the topic. The classical two-stage linear stochastic programming problems can be formulated as

$$\min_{x \in X} \{g(x) := c^T x + \mathbb{E}[Q(x, \delta(\omega))]\}$$

where $Q(x, \delta)$ is the optimal value of the second-stage problem

$$\min_y q^T y$$

Subject to $Tx + Wy \leq h$

Here $x \in \mathbb{R}^n$ is the first-stage decision vector, X is a polyhedral set, defined by a finite number of linear constraints, $y \in \mathbb{R}^m$ is the second-stage decision vector, and $\delta = (q, T, W, h)$ contains the data of the second-stage problem.

The first stage variables $x$ must be decided before the realisations of the random variable $\omega$, and the second stage or recourse variables $y$ are taken, as corrective actions after the value of random variables become known. That is, the recourse actions are a compensation for any infeasibility from the first stage decisions; the objective is to minimise the sum of the first stage cost and the expected value of recourse costs.

### 3.2.2. Stochastic demand model (SDM) for the integrated nurse scheduling problem

We formulate the Stochastic Demand Model (SDM) as a two-stage integer stochastic program. In the first stage, before a realisation of patient demand is known, staffing decisions are made, i.e. the assignment of shifts to nurses based on the available nurses and (estimated) baseline requirement. In the second stage, the patient demand, i.e. the real required number of nurses is realised, and adjustment needs to be made to meet the requirement. The recourse actions are adding additional nurse shifts to cover understaffing or cancelling surplus shifts when overstaffing happens. The expected value of shortfall and surplus of shifts will be minimised. In Table 1 we present the notations used in the model.

**Table 1**
Notations.

| The first-stage problem: | |
| --- | --- |
| *Parameters* | |
| $I$ | The set of nurses (index $i$) |
| $J$ | The set of days during the planning period (index $j$) |
| $W$ | The set of weeks in the planning period (index $w$) |
| $K$ | The set of shift types, for example, {$E$ (Early), $D$ (Day), $L$ (Late), $N$ (Night)} (index $k$) |
| $K^U$ | The set of unwanted shift patterns, for example, {$DE$, $LE$, $LD$, $EN$} (index $k'$) |
| $n_1$ | Maximum number of working shifts a nurse can take in the period |
| $n_2$ | Maximum number of night shifts a nurse can take in the period |
| $n_3$ | Minimum number of regular shifts a nurse need to take in the period |
| $n_4$ | Minimum number of weekends off a nurse should take in the period |
| $M$ | A big constant number |
| $R_{jk}$ | Baseline required number of nurses on day $j$ with shift $k$ |
| $c_1$ | Regular wage rate per shift |
| $c_2$ | Overtime wage rate per shift |
| $c_3$, $c_4$ | Penalty for violating the corresponding soft constraint |
| *Decision variables* | |
| $sr_{ijk}$ | Binary, takes value 1 if nurse $i$ on day $j$ takes shift $k$ with regular pay, 0 otherwise. |
| $so_{ijk}$ | Binary, takes value 1 if nurse $i$ on day $j$ takes shift $k$ with overtime pay, 0 otherwise. |
| $SR_i$ | Binary, takes value 1 if nurse $i$ works regular shifts, 0 otherwise. |
| $SO_i$ | Binary, takes value 1 if nurse $i$ works overtime shifts, 0 otherwise. |
| $dev1, dev2$ | Integer, the amount of deviation when modelling the corresponding soft constraints |
| The second-stage problem: | |
| *Parameters* | |
| $\Omega$ | The set of all scenarios (index $\omega$) |
| $p^\omega$ | The probability of scenario $\omega$ |
| $R_{jk}^\omega$ | The required number of nurses under scenario $\omega$ on day $j$ with shift $k$ |
| $q^+$ | The cost of adding a shift |
| $q^-$ | The cost of cancelling a shift |
| *Decision variables* | |
| $\alpha_{jk}^\omega$ | Integer, the additional number of nurse shift need to be added on day $j$ with shift $k$ for scenario $\omega$ |
| $\beta_{jk}^\omega$ | Integer, the excess number of nurse shift need to be cancelled on day $j$ with shift $k$ for scenario $\omega$ |

The stochastic demand model (SDM) for the integrated nurse scheduling problem can be formulated as follows:

$$\min c_1 \sum_i \sum_j \sum_k sr_{ijk} + c_2 \sum_i \sum_j \sum_k so_{ijk} + c_3 \sum_i \sum_j dev1_{ij}$$
$$+ c_4 \sum_i \sum_j \sum_k dev2_{ijk} + \sum_\omega p^\omega \sum_j \sum_k \left( q^+ \alpha_{jk}^\omega + q^- \beta_{jk}^\omega \right)$$

s.t.

$$\sum_k sr_{ijk} + so_{ijk} \leq 1, \forall i, j \tag{1}$$

$$\sum_j \sum_k sr_{ijk} \leq MSR_i, \forall i \tag{2}$$

$$\sum_j \sum_k so_{ijk} \leq MSO_i, \forall i \tag{3}$$

$$SO_i \leq SR_i, \forall i \tag{4}$$

$$\sum_i (sr_{ijk} + so_{ijk}) \geq R_{jk}, \forall j, k \tag{5}$$

$$\sum_j \sum_k (sr_{ijk} + so_{ijk}) \leq n_1, \forall i \tag{6}$$

$$\sum_j (sr_{ijN} + so_{ijN}) \leq n_2, \forall i \tag{7}$$

$$\sum_j \sum_k sr_{ijk} \geq n_3, \forall i \tag{8}$$

$$\sum_w \sum_k (sr_{iw(sat)k} + sr_{iw(sun)k} + so_{iw(sat)k} + so_{iw(sun)k}) \leq 2|W| - 2n_4, \forall i \tag{9}$$

$$sr_{i(j-1)N} - sr_{ijN} + sr_{i(j+1)N} \geq 0, \forall i, j \in \{2, |J| - 1\} \tag{10}$$

$$sr_{i(j-1)N} - \sum_{k \in \{E,L,D\}} sr_{ijk} + \sum_{k \in \{E,L,D\}} sr_{i(j+1)k} \leq 1, \forall i, j \in \{2, |J| - 1\} \tag{11}$$

$$sr_{i(j-1)N} + \sum_{k \in \{E,L,D\}} sr_{ijk} - \sum_{k \in \{E,L,D\}} sr_{i(j+1)k} \leq 1, \forall i, j \in \{2, |J| - 1\} \tag{12}$$

$$sr_{i(j-1)N} + \sum_{k \in \{E,L,D\}} sr_{ijk} + \sum_{k \in \{E,L,D\}} sr_{i(j+1)k} \leq 2, \forall i, j \in \{2, |J| - 1\} \tag{13}$$

$$\sum_k (sr_{i(j-1)k} - sr_{ijk} + sr_{i(j+1)k}) + dev1_{ij} \geq 0, \forall i, j \in \{2, |J| - 1\} \tag{14}$$

$$sr_{ijk_1} + sr_{ijk_2} - dev2_{ijk'} \leq 1, \forall i, j \in \{1, |J| - 1, (k_1, k_2) \in K^U\} \tag{15}$$

$$\sum_i (sr_{ijk} + so_{ijk}) + \alpha_{jk}^\omega - \beta_{jk}^\omega \geq R_{jk}^\omega, \forall \omega, j, k \tag{16}$$

$$sr_{ijk}, so_{ijk}, SR_i, SO_i \in \{0, 1\}, \forall i, j, k \tag{17}$$

$$\alpha_{jk}^\omega, \beta_{jk}^\omega, \text{ Integer}, \forall j, k, \omega \tag{18}$$

The objective function minimises the aggregated cost which consists of regular time wage and overtime wage, the penalty from violations of soft constraints and unwanted shift patterns, as well as the expected penalty costs occurred from nurse shortage and surplus when the second stage patient demand is realised.

Constraint (1) states that each nurse can only start one shift each day. It also serves as the exclusive constraint stating that each nurse on a single day cannot take a regular shift and an overtime shift at the same time. Constraint (2) imposes a relation constraint using a large constant value $M$ between the indicator variable $SR_i$ and assignment variables $sr_{ijk}$. It states that if nurse $i$ takes a regular shift $sr_{ijk}$, then the indicator variable $SR_i = 1$, while $SR_i = 0$ means nurse $i$ is not assigned. The same rule applies to overtime shift, defined by constraint (3). Constraint (4) ensures that when a nurse works additional time over a regular shift, he or she also works through a required shift. i.e. if $SO_i = 1$, then $SR_i = 1$. Constraint (5) ensures a sufficient number of required nurses are assigned over the planning period based on the baseline requirement before the realisation of patient demand. Constraint (6) limits the maximum number of working shifts during the planning period. Constraints (7) limits the maximum number of night shifts that a nurse should take. Constraint (8) states that a nurse should work a certain minimum number of regular shifts during the planning period. Constraint (9) states a nurse must receive a certain minimum number of complete weekends off during the planning period, where $|W|$ denotes the number of weeks in the planning period. Constraint (10) states there should be no stand-alone night shift, i.e. no night shift between two non-night shifts. Constraints (11, 12), and (13) impose that there must be at least two days off after a night shift, i.e. no sequence of "NOW" "NWO" "NWW", where N, O and W denote a night shift, a day off and a regular working shift, respectively. Constraint (14) penalises a stand-alone regular shift, i.e. there should be only one working day between two off days. Constraint (15) penalises unwanted regular shift patterns such as Day Early, Late Early, Late Day, and Early Night. Constraint (16) is the adjusted coverage constraint after realisation of patient demand. It states that on each day, the assigned number of nurse shifts at the first stage, after cancelling excess nurse shifts and adding additional nurse shifts, should meet the demand for each shift in each scenario.

### 3.2.3. Stochastic demand model with CVaR constraint

In most service systems, staffing and scheduling determine both cost and service qualities [8]. This is especially true in health care systems.

A common approach is to treat staffing level as a minimum coverage constraint that needs to be strictly met. This approach has been applied widely with Dantzig's formulation in the literature, and in some of our previous work [10,20]. This approach is appealing yet less flexible under uncertainty. A more flexible approach is using a probability constraint to limit the expected probability that the number of patients exceeds the nurse-to-patient ratio as proposed in [30]. Another innovative approach is to adopt the concept of robust optimisation and choose the worst case to determine the smallest number of required staff [6]. Kim and Mehrotra [14] adopt a big M method by setting a sufficiently large penalty to track the nurse-to-patient level. Understaffing needs to be specially attended to, as it can severely deteriorate the quality of care [23,26]. In this paper, we introduce an additional constraint adapted from the financial domain to control the risk of understaffing.

Value at Risk (VaR) has been widely used in finance to estimate the exposure to risk by estimating the loss of a finance product. VaR represents the maximum loss associated with a specific confidence level. However, it does not explain the magnitude of loss when the VaR limit is exceeded. Conditional Value-at-Risk (CVarR) is firstly proposed by Rockafellar and Uryasev [22]. It is defined as the expected value of losses strictly exceeding VaR, shown in Fig. 1. CVaR is a coherent risk measure, and has superior mathematical properties compared with VaR. It can be applied either as an objective function or a constraint to control the risk of loss. In both cases it can be reduced to a set of linear functions, which are very easy to optimise in mathematical programming.

In general, we can simply add a CVaR constraint in a model to control the loss under a user specified threshold value $\mu$ by specifying that

$$CVaR \leq \mu$$

## 4. Solution approach

In this section, we describe how the two-stage stochastic programming models are solved. We first discuss how to linearise the CVaR constraint and solve the problem as an integer linear program. Then, we discuss how to construct scenarios for the problem. The resulted stochastic integer programming is finally solved by the IBM ILOG CPLEX solver.

The objective function of a general two-stage stochastic programming (SP) model is to minimise the first stage cost and the expected value of the second stage cost. Recourse variables can be continuous or integer. Our model is a two-stage SP model with integer recourse, which is very challenging to solve. The solution approaches to two-stage SP with integer recourse can be generally grouped into exact methods and heuristic methods. Ahmed et al. [2] adopt an L-shape method based on Bender's decomposition, which incorporates a Branch-

and-Bound procedure to achieve optimality. Kim and Mehrotra [14] developed a modified multi-cut approach in the L-shape algorithm with a prioritising branching strategy. Ahmed and Shapiro [1] reported a general sample average approximation algorithm for stochastic integer optimisation, which can provide an exact optimal solution with a large enough sample size. However, for a relatively small sample size, only a good approximation solution can be obtained. In our problem, we will first linearise the CVaR constraint, and then a relatively small number of scenarios will be generated. IBM CPLEX will be used to solve the resulting stochastic integer programming.

### 4.1. Linearisation of CVaR constraint

As we described in Section 3.2.3, CVaR is used to control the loss under a user-specific value $\mu$. The appealing property of CVaR is that it can be re-written into a set of easy-to-solve linear functions. With the notations defined in Table 2 we demonstrate how a CVaR constraint can be written as a set of linear constraints.

According to Rockafellar and Uryasev [22], the general $CVaR \leq \mu$ term can be replaced by the following linear constraints (19)–(21):

$$\xi + \frac{1}{(1-\sigma)} \sum_{\omega} p^{\omega} z^{\omega} \leq \mu \tag{19}$$

$$z^{\omega} \geq 0, \forall \omega \tag{20}$$

$$z^{\omega} \geq f(\mathbf{x}, \mathbf{y}^{\omega}) - \xi, \forall \omega \tag{21}$$

We can now express the constraints using the notations defined in Table 2 in the context of nurse scheduling as follows:

$$\min c_1 \sum_{i} \sum_{j} \sum_{k} sr_{ijk} + c_2 \sum_{i} \sum_{j} \sum_{k} so_{ijk} + c_3 \sum_{i} \sum_{j} dev1_{ij}$$

$$+ c_4 \sum_{i} \sum_{j} \sum_{k} dev2_{ijk} + \sum_{\omega} p^{\omega} \sum_{j} \sum_{k} \left( q^+ \alpha_{jk}^{\omega} + q^- \beta_{jk}^{\omega} \right)$$

s.t. (1)–(18)

$$\xi + \frac{1}{(1-\sigma)} \sum_{\omega} p^{\omega} z^{\omega} \leq \mu \tag{19}$$

$$z^{\omega} \geq 0, \ \forall \omega \tag{20}$$

$$z^{\omega} \geq \sum_{j} \sum_{k} \left( R_{jk}^{\omega} - \sum_{i} (sr_{ijk} + so_{ijk}) - \alpha_{jk}^{\omega} + \beta_{jk}^{\omega} \right) - \xi, \ \forall \omega \tag{22}$$

We denote it as the stochastic demand model with CVaR constraint or SDM-CVaR.

### 4.2. Scenario generation

We used two different approaches to generate scenarios:

(1) Based on historical data:

In Parisio and Jones [18], demand vectors were generated based on historical data, and then fed into a pool. From the pool, a small number of vectors were randomly selected to construct the demand scenarios. In our model, $R_{jk}^{\omega}$ represents the number of nurses required for a given shift $k$ on day $j$ under a given scenario $\omega$. Inspired by Parisio and Jones [18], we considered the historical number of patients occupying beds as a representation of the true distribution of patient demand and collected this weekly over a 12-month period (52 weekly data). Then these patient demand patterns were fed into a pool. From this pool, we randomly selected 4 weekly data values to construct one monthly demand scenario - each of the scenarios has equal probability. Then we used the nurse-to-patient ratio to convert the number of patients into the number of nurses required. We have taken this approach because it is practical and normally provides a good approximation of the true demand.
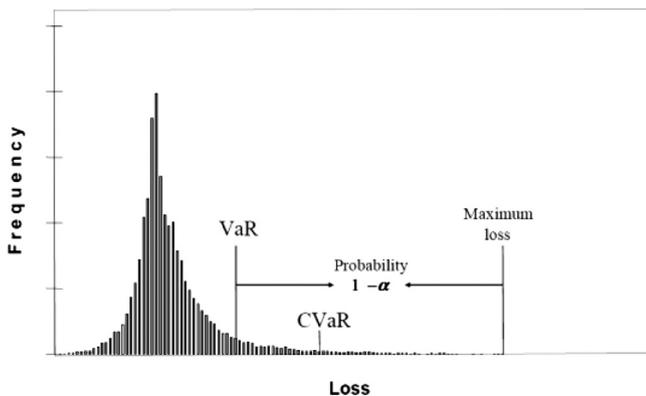


**Fig. 1.** VaR and CVaR [22].

**Table 2**

Notations for a general and Nurse Rostering Problem (NRP) models with CVaR.

| | |
|---|---|
| $x$ | Decision vector. |
| | In our NRP model, $x$ consists of the decision variables $sr_{ijk}$ and $so_{ijk}$. |
| $y^\omega$ | Random vector that influences the loss of decision $x$. |
| | In our NRP model, it is the patient demand uncertainty denoted by $R_{jk}^\omega$, i.e. the number of nurses required on day $j$ with shift of type $k$ under scenario $\omega$ |
| $f(x, y^\omega)$ | A loss function that is generated by $x$ and $y$. |
| | In our NRP model, $R_{jk}^\omega - \sum_i (sr_{ijk} + so_{ijk})$ is the nurse shortage function on day $j$ with shift type $k$ under scenario $\omega$. Thus $f(x,$ |
| | $y^\omega) = \sum_j \sum_k (R_{jk}^\omega - \sum_i (sr_{ijk} + so_{ijk}) - \alpha_{jk}^\omega + \beta_{jk}^\omega)$ |
| $p^\omega$ | The probability of scenario $\omega$. |
| | In our NRP model, it is the probability of patient demand scenario $\omega$. |
| $\xi$ | The VaR value in the optimal solution. |
| $z^\omega$ | Auxiliary variables in the linear programming formulation which represent the loss (i.e. $f(x, y^\omega)$) in excess of the VaR value (i.e. $\xi$). |
| $\sigma$ | A user specified percentile value, i.e. the confidence level, 95% in our case |
| $\mu$ | A user specified threshold value of loss |

**Table 3**

Shift types, durations and baseline demand. Each shift covers 9 h including one hour resting time, except for night shifts that contain no resting time. Demand is based on historical data.

| Shift type | Start time | End time | Demand | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| Early | 07:00 | 16:00 | 3 | 3 | 3 | 3 | 3 | 2 | 2 |
| Day | 08:00 | 17:00 | 3 | 3 | 3 | 3 | 3 | 2 | 2 |
| Late | 14:00 | 23:00 | 3 | 3 | 3 | 3 | 3 | 2 | 2 |
| Night | 23:00 | 07:00 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

(1) Based on Auto-Regressive Integrated Moving Average (ARIMA) forecasts:

In Zinouri's work [32], empirical forecast errors were used to generate demand scenarios for the model. More specifically as in Kim and Mehrotra [14], long-term forecasts were obtained using the ARIMA method. We applied this method to generate demand scenarios. We defined 52-week as a time window and rolled it forward one day to create a new time window. Thus, 364 ($52 \times 7$) time windows were created over a year period. For each time window, a forecast is generated using ARIMA. Then the forecast error vectors were fed into a pool. From this pool, we randomly selected a certain number of error vectors and added them to the mean point forecast to generate demand scenarios.

Kaut and Wallace [13] stated that in stochastic programming we only solve an approximation of the stochastic programming model with a finite number of scenarios. The quality of this approximation is directly linked to the quality of the scenarios. The quality of the scenarios is problem dependent. The number of scenarios is also important. We would like the number of constructed scenarios to be relatively modest so that the resulting model can be solved with reasonable computational effort.

These scenario generation methods have been applied previously in several nurse-scheduling problems and simply adapted to our problem here, given their suitability. While an in-depth investigation of scenario generation methods would be worthy on its own in the literature, it is out of the scope of this manuscript.

## 5. A case study

The development of the approach is based on several benchmark nurse rostering problems, e.g. GPOST and ORTEC, which are publicly available at http://www.schedulingbenchmarks.org. They are monocyclic problems and different from each other with respect to the parameters, such as the number of nurses, number of shift types and length of scheduling periods. These are notably simplified problems, only serving the purpose of developing modelling and solution approaches. These problems preserve the generic constraints, such as

coverage constraints and shift pattern constraints, in general nurse rostering problems. The model and the solution approach can thus be adapted and applied to problems with similar features and constraints. The main problem, i.e. integrated nurse-scheduling problem, in this section is based on a variant of the ORTEC problem [10] with the additional characteristics described in Section 3 and historical patient demand patterns. We use this problem as a case study to explore the benefits of SDM and SDM-CVaR models. The models are developed in C++ with concert technology in CPLEX on top of the CPLEX solver.

### 5.1. Problem instances and input data

We created the problem instances covering the period from January to December 2013. Shift types, start times, and end times are presented in Table 3. Fulltime nurses work 36 h regular time per week. The working regulations of the hospital in our case study state that a nurse may work at most one extra 9-hour overtime shift per week. However, on a single day, if a nurse has already taken a regular shift, she/he cannot take an overtime shift on the same day.

Daily patient census data from January to December over the study period were applied to create the problem instances and scenarios. A 1:4 nurse-to-patient ratio was applied to the patient demand during the daytime to obtain the baseline of nurses required shown in Table 3 as an example.

The two methods based respectively on historical data and ARIMA forecasts and described in Section 4.2 were tested to generate demand scenarios in our experiment. The model parameters were set as follows: the number of shift types was 4, the number of weeks in the period was 4, $n_1 = 24$, $n_2 = 3$, $n_3 = 16$, $n_4 = 4$, $c_1 = 10$, $c_2 = 15$, $c_3 = c_4 = 5$, $q^+ = 18$, $q^- = 2$.

### 5.2. Evaluation of SDM model and SDM-CVaR model solutions and computational time

Table 4 presents the performance of the SDM model based on the two scenario generation methods, i.e. based respectively on historical data and on ARIMA forecasts. For each of the methods, we tested 50 and 200 scenarios. The choice of using 50 and 200 scenarios in our empirical study is arbitrary. The first two rows report the numbers of regular and overtime shifts assigned at the first stage, and the third and fourth rows report the adjustment made (added and cancelled shifts) at the second stage to meet demands. It can be seen that scenarios generated with ARIMA forecasts required more adjustments, maybe due to larger fluctuations in demand. The CPU time needed to solve the models are similar as well as the optimality gap of the final solution. The fundamental reason of the similar CPU time is that the CVaR constraint has been transformed to linear constraints, which do not increase the complexity of the problem. The CPLEX parameter CPX_PARAM_EPGAP (gap to the optimum) was set to 0.01.

Given the concerns about nurses' job satisfaction, many hospitals

**Table 4**
SDM model solution evaluation.

| | Historical data based scenario generation | | Forecast data based scenario generation | |
|---|---|---|---|---|
| | 50 scenario solution | 200 scenario solution | 50 scenario solution | 200 scenario solution |
| No. of regular shift ($sr_{ijk}$) | 231 | 231 | 231 | 231 |
| No. of overtime shift ($so_{ijk}$) | 25 | 25 | 25 | 25 |
| No. of added shift ($\alpha_{jk}^\omega$) | 36 | 38 | 168 | 189 |
| No. of cancelled shift ($\beta_{jk}^\omega$) | 6 | 6 | 29 | 35 |
| Soft constraint violation (*dev1,dev2*) | 0 | 0 | 0 | 0 |
| CPU time | 61.64 s | 80.09 s | 66.10 s | 93 s |
| Optimality gap | 1.02% | 0.90% | 0.43% | 0.26% |

**Table 5**
SDM-CVaR model solution evaluation.

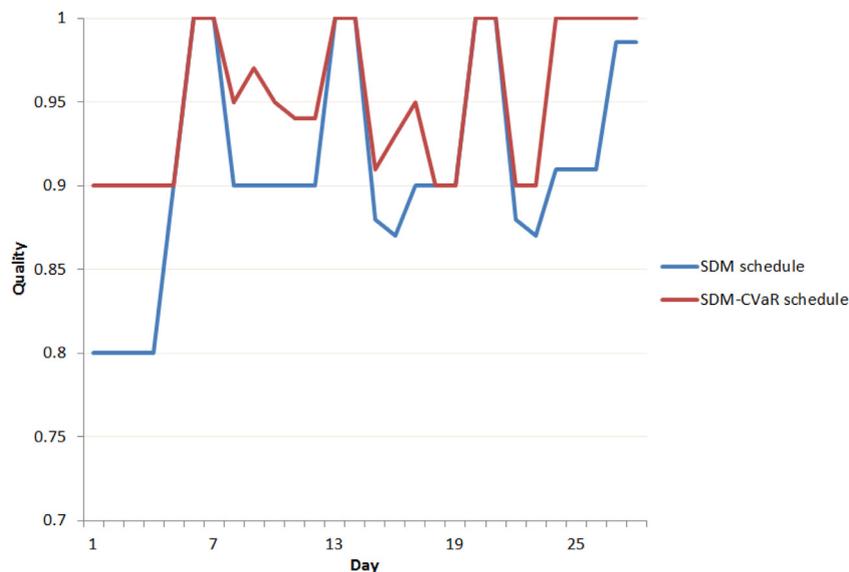| | Historical data based scenario generation | | Forecast data based scenario generation | |
|---|---|---|---|---|
| | 50 scenario solution | 200 scenario solution | 50 scenario solution | 200 scenario solution |
| No. of regular shift ($sr_{ijk}$) | 231 | 231 | 231 | 231 |
| No. of overtime shift ($so_{ijk}$) | 25 | 25 | 25 | 25 |
| No. of added shift ($\alpha_{jk}^\omega$) | 19 | 20 | 90 | 90 |
| No. of cancelled shift ($\beta_{jk}^\omega$) | 15 | 14 | 50 | 35 |
| Soft constraint violation (*dev1,dev2*) | 0 | 0 | 0 | 0 |
| CPU time | 103.82 s | 139.58 s | 125.30 s | 155.56 s |
| Optimality gap | 0.6% | 0.34% | 0.82% | 0.20% |



Fig. 2. Quality factors of SDM and SDM-CVaR schedule.

**Table 6**
SDM-CVaR model with different parameter $\mu$.

| | $\mu = 50, \sigma = 95\%$ | $\mu = 20, \sigma = 95\%$ |
|---|---|---|
| No. of regular shift ($sr_{ijk}$) | 231 | 231 |
| No. of overtime shift ($so_{ijk}$) | 25 | 25 |
| No. of added shift ($\alpha_{jk}^\omega$) | 19 | 95 |
| No. of cancelled shift ($\beta_{jk}^\omega$) | 15 | 10 |
| Soft constraint violation (dev1,dev2) | 0 | 0 |
| CPU time | 103.82 s | 104.56 s |
| Optimality gap | 0.6% | 0.93% |

are actively seeking ways to improve the situation. We now show how the SDM model can be used to improve the overall desirability of the schedule in terms of reduction of personnel constraint violations.

The violations of soft constraints such as unwanted shift patterns were measured using $c_3 \sum_i \sum_j dev1_{ij} + c_4 \sum_i \sum_j \sum_k dev2_{ijk\text{in}}$ in the

objective function defined in Section 3.2.2. We penalised unwanted shift patterns by *dev1, dev2*. We observed that the penalty cost was 0, as shown in the fifth row in Table 4. This demonstrates that the solutions satisfied all these constraints. This finding demonstrates that better work satisfaction for nurses can be achieved in the SDM model solution.

Table 5 presents the performance of the SDM-CVaR model also based on the two scenario generation methods. The user-specified parameters in the model were set to $\mu = 50, \sigma = 95\%$. The difference between the SDM and SDM-CVaR models mainly exists in the scheduled adjustment shifts, i.e. $\alpha_{jk}^\omega$, $\beta_{jk}^\omega$. We will investigate the SDM-CVaR model in more details in the next section.

### 5.3. Comparison of SDM and SDM-CVaR schedule

To observe a clearer comparison on the basic SDM and SDM-CVaR models, we compared the schedule quality of the two models for the same set of historical data based scenarios. The actual historical number

**Table 7**
Comparisons of the SDM and SDM-CVaR models on the 12 instances.

| | Baseline No. | SDM | | | | SDM-CVaR | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | No. of regular | No. of overtime | No. of adjustments | Soft constraint violations | No. of regular | No. of overtime | No. of adjustments | Soft constraint violations |
| Jan | 286 | 231 | 25 | 42 | 0 | 231 | 25 | 35 | 0 |
| Feb | 256 | 212 | 22 | 35 | 4 | 212 | 22 | 30 | 4 |
| Mar | 280 | 230 | 25 | 40 | 3 | 230 | 25 | 32 | 3 |
| Apr | 286 | 231 | 25 | 44 | 0 | 231 | 25 | 33 | 0 |
| May | 301 | 235 | 32 | 46 | 1 | 235 | 32 | 40 | 1 |
| Jun | 294 | 231 | 25 | 42 | 0 | 231 | 25 | 40 | 0 |
| Jul | 289 | 230 | 25 | 42 | 0 | 230 | 25 | 42 | 0 |
| Aug | 301 | 236 | 30 | 53 | 2 | 236 | 30 | 47 | 2 |
| Sep | 298 | 234 | 28 | 43 | 1 | 234 | 28 | 43 | 1 |
| Oct | 286 | 230 | 26 | 38 | 0 | 230 | 26 | 38 | 0 |
| Nov | 283 | 230 | 26 | 32 | 2 | 230 | 26 | 28 | 2 |
| Dec | 260 | 210 | 20 | 55 | 1 | 210 | 20 | 42 | 1 |
| Avg | 285 | 228.33 | 25.75 | 42.67 | 1.17 | 228.33 | 25.75 | 37.5 | 1.17 |
| s.d. | 13.77 | 8.01 | 3.06 | 6.29 | 1.28 | 8.01 | 3.06 | 5.60 | 1.28 |

of nurse shifts required was applied as the baseline for comparison, in contrast to the schedules obtained by solving the SDM model (SDM schedule) and the SDM-CVaR model (SDM-CVaR schedule).

To measure the quality of a schedule $s$, we applied a quality factor as defined in Parisio and Jones [18] as follows:

$$\theta = 1 - \frac{\sum_t \varepsilon_t}{\sum_t d_t}$$

where $\varepsilon_t = |s_t - d_t|$ is the deviation between the nurse shifts assigned at time $t$ by schedule $s$ and the actually required number of nurse shifts (demand) at time $t$. The quality factor, which is always between 0 and 1, is plotted in Fig. 2 for both SDM and SDM-CVaR schedules. This factor ranges from 0.8 to 1 with an average of 90% in the SDM schedule, and from 0.9 to 1 with an average of 95% in the SDM-CVaR schedule. Therefore, the average quality of the SDM-CVaR schedule is approximately 5% better than that of the SDM schedule.

*5.4. Evaluation of the CVaR constraint*

The coverage constraint $\sum_i (sr_{ijk} + so_{ijk}) + \alpha_{jk}^{\omega} - \beta_{jk}^{\omega} = R_{jk}^{\omega}, \forall \omega, j, k$ includes adjustment shifts from the nurse pool applied to meet demand fluctuations. By adding extra shifts $\alpha_{jk}^{\omega}$, the downside risk of the schedule, i.e. the nurse shift shortage, can be restricted by constraints (19)–(21). That is, we can control the number of shortage shifts by setting different user specified values for $\mu$ in Eq. (19) in the CVaR constraint. For instance, if we want to keep the shift shortage under 50 (i.e. set $\mu = 50$) with 95% confidence in the CVaR constraint, we need about 20 extra shifts as shown in the third row of Table 6. If we want to have a tighter control on nurse shortage, we can set the CVaR constraint to a smaller value e.g. $\mu = 20$, with 95% confidence. However, a large number of extra shifts is required (95 for $\mu = 20$ vs 20 for $\mu = 50$) to achieve this target.

Table 7 compares the results of SDM and SDM-CVaR based on 50 scenarios generated from historical data for the 12-instance (monthly) set. The second column presents the baseline requirement of shifts required in the first stage. The regular and overtime shifts are constructed based on this baseline requirement. The number of adjustments consists of the number of shifts added and cancelled. We also report the number of soft constraints violations. From the results we can see that there is no difference between the two models in terms of soft constraints violations. The average and standard deviations in Table 7 show that the difference lies in the number of adjustments. The incorporation of the CVaR constraint into the SDM-CVaR model leads to less adjustments.

## 6. Conclusions

Healthcare systems show non-stationary and uncertain demand. Therefore, decision making on nurse staffing and scheduling should be considered together in a stochastic setting. When patient demand fluctuates, overstaffing or understaffing may occur. In particular, understaffing needs to be paid more attention to as it can severely deteriorate the quality of care.

In this paper, two integrated nurse scheduling models with patient demand uncertainty have been proposed and analysed. The experimental results showed that the Stochastic Demand Model (SDM) with CVaR constraint is able to control the number of shortage shifts at a user-specified confidence level. Our research showed how a nurse schedule can be adjusted with respect to the level of risk the decision maker is willing to take.

Using our model could potentially lead to healthcare quality improvement and result in cost benefit. In this paper, we used historical patient data and forecast error vectors based on the ARIMA method to generate scenarios with rough estimate of patient demand for the SDM and SDM-CVaR models. With more accurate demand scenarios, we could potentially have a better scheduling of nurse shifts. Further work will investigate the generation of a larger number of scenarios using different methods, such as Sample Average Approximation method, based on a more detailed study of patient demand distributions.

## References

[1] Ahmed S, Shapiro A. The sample average approximation method for stochastic programs with integer recourse School of Industrial and Systems Engineering, Georia Institute of Technology; 2002. Technical report.

[2] Ahmed S, Tawarmalani M, Sahinidis NV. A finite branch-and-bound algorithm for two-stage stochastic integer programes. Math Program 2004;100:355–77.

[3] Appleby J, Galea A, Murray R. The NHS productivity challenge-Experience from the front line. 2014. https://www.kingsfund.org.uk/publications/nhs-productivity-challenge.

[4] Bagheri M, Devin AG, Izanloo A. An application of stochastic programming method for nurse scheduling problem in real word hospital. Comput Ind Eng 2016;96:192–200.

[5] Bard JF, Purnomo HW. Short-term nurse scheduling in response to daily fluctuations in supply and demand. Health Care Manag Sci 2005;8:315–24.

[6] Chen P, Lin Y, Peng N. A two-stage method to determine the allocation and scheduling of medical staff in uncertain environments. Comput Ind Eng 2016;99:174–88.

[7] Dantzig G. A comment on Edies traffic delay at toll booths. Oper Res 1954;2:339–41.

[8] Defraeye M, Nieuwenhuyse IV. Staffing and scheduling under nonstationary demand for service: a literature review. Omega 2016;58:4–25.

[9] Easton FF, Rossin DF, Borders WS. Analysis of alternative scheduling policies for hospital nurses. Prod Oper. Manag 1992;1(2):159–74.

[10] He F, Qu R. A constraint programming based column generation approach to nurse rostering problems. Comput Oper Res 2012;39(12):3331–43.

[11] Jonas Ingels J, Maenhout B. The impact of reserve duties on the robustness of a personnel shift roster: an empirical investigation. Comput Oper Res

2015;61:153–69.

[12] Ingolfsson A, Haque A, Umnikov A. Accounting for time-varying queueing effects in workforce scheduling. Eur J Oper Res 2002;139:585–97.

[13] Kaut M, Wallace SW. Evaluation of scenario-generation methods for stochastic programming. Pac J Optim 2007;3(2):257–71.

[14] Kim K, Mehrotra S. A two-stage stochastic integer programming approach to integrated staffing and scheduling with application to nurse management. Oper Res 2015;12:1431–51.

[15] Kishimoto S, Yamashita M. A successive LP approach with C-VaR type constraints for IMRT optimization. Oper Res Health Care 2018;17:55–64.

[16] Maenhout B, Vanhoucke M. An integrated nurse staffing and scheduling analysis for longer-term nursing staff allocation problems. Omega 2013;41(2):485–99.

[17] Najjarbashi A, Lim GJ. A variability reduction method for the operating room scheduling problem under uncertainty using CVaR. Oper Res Health Care 2019;20:25–32.

[18] Parisio A, Jones CN. A two-stage stochastic programming approach to employee scheduling in retail outlets with uncertian demand. Omega 2015;53:97–103.

[19] Punnakitikashem P, Rosenberber JM, Buckley-Behan DF. A stochastic programming approach for integrated nurse staffing and assignment. IIE Trans 2013;45:1059–76.

[20] Qu R, He F. A hybrid constraint programming approach for nurse rostering problems. In: Allen T, Ellis R, Petridis M, editors. applications and innovations in intelligent systems XVI. London: Springer; 2008. p. 211–24.

[21] Rafferty A, Maben J, West E, Robbinson D. What makes a good employer? Global Nurs Rev Inititiva (WHO) 2005;3:1–84.

[22] Rockafellar RT, Uryasev S. Optimization of conditional Value-at-Risk. J Risk 2000;2(3):21–41.

[23] Rogers AE, Hwang WT, Scott LD, Aiken LH, Dinges DF. The working hours of hospital staff nurses and patient safety. Health Affair 2004;23(4):202–12.

[24] Shapiro A, Philpott A. A Tutorial on Stochastic Programming. 2007. https://www2.isye.gatech.edu/people/faculty/Alex_Shapiro/TutorialSP.pdf.

[25] Tam B, Ehrgott M, Ryan D, Zakeri G. A comparison of stochastic programming and bi-objective optimisation approaches to robust airline crew scheduling. OR Spectr 2011;33:49–75.

[26] Ulrich CM, Wallen G, Grady C. The nursing shortage and the quality of care. N Engl J Med 2002;347:1118–9.

[27] Van den Bergh J, Beliën J, De Bruecker P, Demeulemeester E, De Boeck L. Personnel scheduling: a literature review. Eur J Oper Res 2013;226(3):367–85.

[28] Venkataraman R, Brusco MJ. An integrated analysis of nurse staffing and scheduling policies. Omega 1996;24(1):57–71.

[29] Wright PD, Bretthauer KM. Strategies for addressing the nursing shortage: coordinated decision making and workforce flexibility. Decis Sci 2010;41(2):373–401.

[30] Wright PD, Mahar S. Centralized nurse scheduling to simultaneously improve schedule cost and nurse satisfaction. Omega 2013;41(6):1042–52.

[31] Zhu X, Sherali HD. Two-stage workforce planning under demand fluctuation and uncertainty. J Oper Res Soc 2007;60:94–103.

[32] Zinouri N. Improving healthcare resource management through demand prediction and staff scheduling PhD thesis 2016