



Contents lists available at ScienceDirect

EURO Journal on Computational Optimization

journal homepage: www.elsevier.com/locate/ejco



A reinforcement learning approach to the stochastic cutting stock problem



Anselmo R. Pitombeira-Neto^{a,b,c,*}, Arthur H.F. Murta^{b,c}

^a Department of Industrial Engineering, Federal University of Ceará, Campus do Pici, Fortaleza, Brazil

^b Graduate Program in Modeling and Quantitative Methods, Federal University of Ceará, Fortaleza, Brazil

^c OPL Lab – Operations Research in Production and Logistics, Federal University of Ceará, Fortaleza, Brazil

ARTICLE INFO

Keywords:

Cutting stock problem
Reinforcement learning
Approximate dynamic programming

ABSTRACT

We propose a formulation of the stochastic cutting stock problem as a discounted infinite-horizon Markov decision process. At each decision epoch, given current inventory of items, an agent chooses in which patterns to cut objects in stock in anticipation of the unknown demand. An optimal solution corresponds to a policy that associates each state with a decision and minimizes the expected total cost. Since exact algorithms scale exponentially with the state-space dimension, we develop a heuristic solution approach based on reinforcement learning. We propose an approximate policy iteration algorithm in which we apply a linear model to approximate the action-value function of a policy. Policy evaluation is performed by solving the projected Bellman equation from a sample of state transitions, decisions and costs obtained by simulation. Due to the large decision space, policy improvement is performed via the cross-entropy method. Computational experiments are carried out with the use of realistic data to illustrate the application of the algorithm. Heuristic policies obtained with polynomial and Fourier basis functions are compared with myopic and random

* Corresponding author at: Department of Industrial Engineering, Federal University of Ceará, Campus do Pici, Fortaleza, Brazil.

E-mail address: anselmo.pitombeira@ufc.br (A.R. Pitombeira-Neto).

URL: <http://www.opl.ufc.br/authors/anselmo/> (A.R. Pitombeira-Neto).

policies. Results indicate the possibility of obtaining policies capable of adequately controlling inventories with an average cost up to 80% lower than the cost obtained by a myopic policy.

© 2022 The Author(s). Published by Elsevier Ltd on behalf of Association of European Operational Research Societies (EURO). This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

1. Introduction

The cutting stock problem (CSP) is a well studied problem in combinatorial optimization with applications in many industries, such as: textile, furniture, paper, glass, construction and manufacturing. In its basic setting, the problem consists in cutting smaller items from larger objects in stock, given the demand from customers, with the objective of minimizing trim loss [1,2]. The problem has many variants which reflect specific features of real applications [3–6]. In this paper, we focus on the most common variant known as the one-dimensional CSP, which has many important applications in the manufacturing processes of steel bars, paper, precast concrete beams, tubes, among others.

The classic mathematical formulation of the CSP is the following: Let d_i be the demand for an item of type $i \in \{1, 2, \dots, m\}$. The items are obtained by cutting larger objects in stock according to predefined cutting patterns. A cutting pattern $j \in \{1, 2, \dots, n\}$ is defined by a vector $a_j = (a_{1j}, a_{2j}, \dots, a_{mj})$ in which a_{ij} is the number of items of type i produced when a stock object is cut according to pattern j . The problem is then to determine how many objects x_j to cut according to pattern j in order to meet the demand while minimizing total trim loss, and can be formulated as the following integer linear program:

$$\min \sum_{j=1}^n g_j x_j, \quad (1)$$

s.t.

$$\sum_{j=1}^n a_{ij} x_j \geq d_i, \quad i \in \{1, 2, \dots, m\},$$

$$x_j \in \mathbb{Z}_+, \quad j \in \{1, 2, \dots, n\}, \quad (2)$$

in which g_j is the trim loss cost associated with pattern j . The CSP can be solved to optimality by using branch-and-bound algorithms.

Although mathematically elegant, the classic formulation of the CSP given in Eqs. (1)–(2) has limited applicability in industry, since it ignores key characteristics

which occur in practice. First, it assumes that the demand is known at the time of decision, but in reality it is often unknown. Second, it is a static decision problem, while cutting stock operations in industry are typically dynamic, in which companies have to decide on a daily or weekly basis how many items to cut. Third, companies frequently keep items in inventory, so that, at a given decision time, there is an initial inventory of items that must be taken into account in the decision of cutting new items to fulfill the demand. This means that, in practice, cutting stock problems are mixed with lot-sizing problems with the added complication that production of individual items is coupled by the cutting patterns. We call this decision problem under a dynamic and uncertain scenario the *stochastic cutting stock problem* (stochastic CSP).

We approach the stochastic CSP from a stochastic optimal control standpoint. This is, we seek a decision policy which specifies at each decision epoch how many objects in stock to cut in each pattern, given current inventory levels of demanded items. We present an unprecedented formulation of the stochastic CSP as a discounted infinite-horizon Markov decision process. Ideally, we would like to find an optimal policy which minimizes expected discounted total cost. Although in theory we may obtain an optimal policy through exact methods, in practice it will likely be computationally infeasible to obtain a guaranteed optimal policy due to the large size of the state and decision spaces. To overcome this difficulty, we propose an approximate solution based on reinforcement learning techniques.

In particular, we develop a model-free off-policy approximate policy iteration algorithm. The action-value function is approximated by a linear combination of basis functions. Samples of state transitions, decisions and rewards (costs in our case) are collected by simulation, and policy evaluation is performed by approximately solving the projected Bellman equation. Policy improvement is carried out by approximately solving a nonconvex nonlinear integer programming problem via the cross-entropy method. We investigate the proposed approach by means of computational experiments with the use of realistic data.

The structure of this paper is the following: In Section 2, we comment on related work; in Section 3, we formulate the stochastic CSP as a discounted infinite-horizon Markov decision process; in Section 4, we detail the proposed approximate policy iteration solution approach; in Section 5, we discuss computational results; finally, in Section 6 we draw some concluding remarks.

2. Related work

Reinforcement learning (RL) is a computational approach to sequential decision problems based on the theory of Markov decisions processes and tries to overcome the well known *curse of dimensionality* through approximations of value functions or policies [7,8]. RL has recently shown impressive results in video games [9], board games [10] and robotics [11,12]. It is closely related to approximate dynamic programming (ADP), a

field with similar theory and methods but which evolved with different target applications [13,14].

There has been increasing interest in applying RL and ADP methods to stochastic combinatorial optimization problems which arise in industry. The applications have spanned diverse areas, such as: transportation [15,16], scheduling [17–19], energy storage management [20,21], ambulance dispatching [22], supply chain management [23,24] and blood bank management [25].

However, as far as we know, RL has not been applied to the stochastic CSP. Some work has addressed the CSP in a dynamic but deterministic environment, in which demand for items is known in advance and supplied in a make-to-order manner. Since the demand may often not be supplied in a single time period, cutting operations have to be scheduled so as to optimize makespan, tardiness or other performance measure. In this case, it has been referred to in the literature as a *combined cutting stock and scheduling problem* [26–30]. Other papers address the CSP in a dynamic deterministic make-to-stock setting, in which point forecasts of future demands are available for a finite horizon and decisions on how many objects to cut are taken in advance. (See e.g. [31–35].)

Just a few works approach a stochastic formulation of the CSP. The earliest we found is due to Sculli [36], who consider the sizes of the objects to be cut as random variables with known probability distributions. The problem is then to determine the position of the first cut of a fixed set of knives so as to minimize expected trim loss. Krichagina et al. [37] approach a CSP in the paper industry in which production and cutting decisions are made in a make-to-stock fashion. In a first stage, the average frequencies with which the cutting patterns are used are determined so as to meet average demand. Then, in a second stage, a scheduling policy decides if it shuts down a paper cutting machine to avoid building up too much inventory.

Some works apply the stochastic programming framework, which is fundamentally different from the stochastic control approach. Alem et al. [38] formulate a static CSP in which demands for items are unknown as a two-stage stochastic program with recourse. In the first stage, before seeing the demand, decisions are made on the items to be cut. In the second stage, after seeing the demand, costs are incurred for holding inventory or a penalty is incurred if demand is not met. The objective is to minimize the total expected cost in both stages. Beraldi et al. [39] propose a similar two-stage stochastic program with the difference that only the patterns to be used are decided upon before seeing the demand, while the actual amounts of objects cut according to each pattern are decided only after seeing the demand. Zanarini [40] also proposes a two-stage stochastic programming model in which demands are uncertain. In the first stage, decisions are made on the sizes of the objects to be used in the cutting process, while in the second stage, after demands are known, the classic deterministic CSP is solved. It is worth noting that these three works consider only static one-time decisions.

In Section 3, we propose an unprecedented formulation of the stochastic CSP as a discounted infinite-horizon Markov decision process. The main novelty of this formula-

tion is that, in contrast to previous works, it takes into account both the dynamic and stochastic nature of the problem in practice.

3. A Markov decision process formulation

Consider a cutting stock problem in a dynamic stochastic environment in which there are m different items that can be demanded and cut from larger objects in stock, and there are n different cutting patterns, where a_{ij} is the number of items of type $i \in \{1, 2, \dots, m\}$ obtained by cutting an object in the pattern $j \in \{1, 2, \dots, n\}$. At each decision epoch t , there is an initial inventory s_{it} for each type of item and a decision x_{jt} must be made on how many objects to cut in each pattern before the demand is known, constrained by the availability of objects.

After the decision is taken, the amount of items produced is added to the initial inventory making up the available inventory to meet the demand. The demanded quantities deplete the available inventory, which results in the final inventory which will also be the initial inventory at subsequent time $t + 1$. Demand which is not fulfilled is lost. There are costs related to trim loss, holding inventory to the next time period and lost sales. We would like to make decisions on which item quantities to cut at each decision epoch so as to minimize expected total cost over time. We name this problem the *stochastic cutting stock problem*, since the demand and costs are not fully known at the decision time.

We formulate the stochastic CSP as a discounted infinite-horizon Markov decision process (MDP). Let $s_t \in \mathcal{S}$ be a vector corresponding to the state of the system at time t and \mathcal{S} a finite set of states. We define $s_t = (s_{1t}, s_{2t}, \dots, s_{mt})$ where $s_{it} \in \{0, \dots, s_{\max}\}$ as the initial inventory of items of type i at time t . Let $x_t = (x_{1t}, x_{2t}, \dots, x_{nt})$ be a vector corresponding to the decision (we also use the term *action* interchangeably) at time t when the state is $s_t = s$, in which $x_{jt} \in \{0, 1, \dots, x_{\max}\}$ corresponds to the number of objects in stock that are cut in pattern j . Decisions are constrained to assume values in a set \mathcal{X}_s of feasible decisions given the state s and we also define the set of all possible decisions as $\mathcal{X} = \bigcup_{s \in \mathcal{S}} \mathcal{X}_s$. We denote as $d_{t+1} = (d_{1,t+1}, d_{2,t+1}, \dots, d_{m,t+1})$ the demand vector for items that occurs in the interval $(t, t + 1]$. We assume that this demand becomes known only *after* the decision to cut the objects in stock, so that d_{t+1} is a random discrete vector with a conditional probability function $\mathbb{P}(d_{t+1}|s_t, x_t)$.

The state transition function is defined by $s_{t+1} = f(s_t, x_t, d_{t+1})$, in which s_{t+1} is the initial inventory at time $t + 1$, whose components $s_{i,t+1}$ are given by

$$s_{i,t+1} = \left[s_{it} + \sum_{j=1}^n a_{ij} x_{jt} - d_{i,t+1} \right]^+, \quad i \in \{1, \dots, m\}, \quad (3)$$

in which $[x]^+ = \max(0, x)$. Notice in (3) that

$$\sum_{j=1}^n a_{ij}x_{jt}$$

corresponds to the quantity of item i added to the inventory by cutting x_{jt} stock objects according to patterns $j \in \{1, 2, \dots, n\}$ at time t . Notice also that s_{t+1} is a random vector, since it is a function of the random vector d_{t+1} .

The cost function (also called *reward function* in reinforcement learning) is the sum of costs associated with trim loss, holding inventory and lost sales, given by

$$\begin{aligned} c(s_t, x_t, d_{t+1}) = & \sum_{j=1}^n g_j x_{jt} + \sum_{i=1}^m h_i^+ \left[s_{it} + \sum_{j=1}^n a_{ij} x_{jt} - d_{i,t+1} \right]^+ \\ & + \sum_{i=1}^m h_i^- \left[d_{i,t+1} - \left(s_{it} + \sum_{j=1}^n a_{ij} x_{jt} \right) \right]^+, \end{aligned} \quad (4)$$

in which g_j is a cost term associated with trim loss by using pattern j , h_i^+ is the inventory holding cost of item i for a time period, and h_i^- is the lost sales cost for item i . Notice that the third term in the right-hand side of (4) corresponds to the total cost of lost sales, which is incurred when

$$d_{i,t+1} - \left(s_{it} + \sum_{j=1}^n a_{ij} x_{jt} \right) > 0,$$

i.e., when the demand for item i is greater than the available inventory. A noteworthy characteristic of the cost function (4) is that it is random at the decision time t , since the decision maker (also called the agent in reinforcement learning) does not know the demand d_{t+1} which will be realized during time interval $[t, t+1)$.

Let π be a stationary Markovian policy defined by $\pi := \{x(s)\}$, in which $x : \mathcal{S} \rightarrow \mathcal{X}$ is a function that associates with each state $s \in \mathcal{S}$ a decision $x = x(s)$. (Notice that we can drop the subscript t by convenience since we henceforth assume stationarity of the system and the policy.) A feasible policy associates with each state $s \in \mathcal{S}$ a decision x in the set \mathcal{X}_s of feasible decisions given by

$$\mathcal{X}_s = \left\{ \begin{array}{l} s_i + \sum_{j=1}^n a_{ij} x_j \leq s_{\max}, \quad i = 1, 2, \dots, m, \\ \sum_{j=1}^n x_j \leq x_{\max}, \\ x_j \in \mathbb{Z}_+, \quad j = 1, 2, \dots, n \end{array} \right\},$$

in which s_{\max} is the maximum inventory and x_{\max} is the maximum number of objects in stock which may be cut in a given time period.

Given an initial state $s_0 = s$, we define as $v_\pi : \mathcal{S} \rightarrow \mathbb{R}$ a value function associated with a policy π :

$$v_\pi(s) := \lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{t=0}^T \gamma^t c(s_t, x(s_t), d_{t+1}) \middle| s = s_0 \right], \quad \forall s \in \mathcal{S}, \quad (5)$$

in which $0 < \gamma < 1$ is a *discount factor* and the expected value is computed with regards to the joint probability function of the states (s_1, s_2, \dots, s_T) induced by policy π . The *Markov decision problem* corresponds to finding an optimal policy π^* which minimizes the value function (5) for all states $s \in \mathcal{S}$:

$$\pi^* \in \arg \min_{\pi \in \Pi_{\text{md}}} v_\pi(s), \quad \forall s \in \mathcal{S}, \quad (6)$$

in which Π_{md} denotes the class of deterministic Markovian policies.

Notice that the problem (6) is often too hard to solve directly, since the computation of the expected value is intractable. To mitigate this difficulty, we can decompose (5) as

$$v_\pi(s_0) = \mathbb{E}[c(s_0, x(s_0), d_1)] + \mathbb{E} \left[\lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{t=1}^T \gamma^t c(s_t, x(s_t), d_{t+1}) \middle| s_1 \right] \middle| s_0 \right], \quad \forall s_0 \in \mathcal{S}.$$

By further noticing that the term within the brackets of the second expectation operator corresponds the value $v_\pi(s_1)$ of starting at state s_1 at time $t = 1$ and following policy π discounted by γ , we obtain a form of the celebrated Bellman equation:

$$v_\pi(s) = \mathbb{E}[c(s, x(s), d) + \gamma v_\pi(s') | s], \quad \forall s \in \mathcal{S}, \quad (7)$$

in which we have dropped time subscripts since we assumed stationarity and the next state $s' = f(s, x, d)$ is given by Eq. (3). The expected value in (7) is computed with regards to the conditional probability function $\mathbb{P}(d|s, x)$ of the demand d . We notice that this includes the independence case, in which $\mathbb{P}(d|s, x) = \mathbb{P}(d)$, as a particular case. The possibility that the demand may be conditional on both the current state s and the decision x in our formulation is an advantage over other formulations used for stochastic problems (such as some stochastic programming formulations) which assume that the demand d to be observed is independent of the current state and decision.

In theory, an optimal policy may be obtained by exact methods such as value iteration, policy iteration or linear programming [7]. Value iteration starts from arbitrary values $v^{(0)}(s)$ assigned to each possible state in \mathcal{S} . Then, at an iteration of the method, the value of each one of the states is updated by assigning to it the right-hand side of Bellman equation at the greedy action:

$$v^{(k+1)}(s) = \min_{x \in \mathcal{X}_s} \mathbb{E}[c(s, x, d) + \gamma v^{(k)}(s') | s], \quad \forall s \in \mathcal{S},$$

for $k = 0, 1, \dots$. Under some conditions [7], this iteration converges to a fixed point v^* as $k \rightarrow \infty$, which is the optimal value function. An optimal policy π^* can then be computed by simply choosing the action at each state s which minimizes the right-hand side of Bellman equation relative to the optimal value function:

$$x^*(s) \in \arg \min_{x \in \mathcal{X}_s} \mathbb{E}[c(s, x, d) + \gamma v^*(s') | s], \quad \forall s \in \mathcal{S}.$$

We are particular interested in an alternative exact method called *policy iteration* [41], which swaps between two steps: policy evaluation and policy improvement. Given an arbitrary policy $\pi^{(k)}$ at iteration k of the method, in the evaluation step we obtain its value function $v_{\pi^{(k)}}$ by solving Bellman equation (7). Then, in the improvement step, a new policy $\pi^{(k+1)} := \{x^{(k+1)}(s)\}$ is obtained by acting greedily with respect to the value function $v_{\pi^{(k)}}$:

$$x^{(k+1)}(s) \in \arg \min_{x \in \mathcal{X}_s} \mathbb{E}[c(s, x, d) + \gamma v_{\pi^{(k)}}(s') | s], \quad \forall s \in \mathcal{S}.$$

It can be shown [7] that $\pi^{(k+1)}$ is better than $\pi^{(k)}$ in the sense that $v_{\pi^{(k+1)}}(s) \leq v_{\pi^{(k)}}(s)$, $\forall s \in \mathcal{S}$. (\geq if maximizing rewards.) Policy iteration begins with an initial policy $\pi^{(0)}$ and generates a sequence of policies $\pi^{(0)}, \pi^{(1)}, \pi^{(2)}, \dots$ which stops when $\pi^{(k+1)} = \pi^{(k)}$. Value iteration may be regarded as a particular case of policy iteration in which the policy evaluation and improvement steps are merged in a single step.

Although policy iteration (and value iteration) is computationally more efficient than alternative exact methods for solving MDPs, such as linear programming [8], it is impractical for problems with large states spaces, which is the case in the stochastic CSP. Its main limitation is that the value function must be computed for all states, and in vector-valued state variables the number of states increases exponentially with the dimension of the vector. (The well known *curse of dimensionality*.) Other limitation of policy iteration is that it relies on the exact computation of the expectation in the Bellman equation (7), which will not be possible in general.

Since exact policy iteration will be computationally infeasible in our case, in the next section we leverage methods developed in the field of reinforcement learning to obtain an approximate solution to the stochastic CSP.

4. An approximate policy iteration solution approach

In this section, we develop an approximate policy iteration method to solve the stochastic CSP. Exact policy iteration alternates between two steps: policy evaluation and policy improvement. However, in approximate policy iteration these two steps are intertwined and the boundaries of each step are not so crispy.

In order to cope with the large state space, we use a parameterized value-function approximation defined by a linear combination of basis functions. Given a greedy policy with respect to a current value function, we gather samples from state transitions,

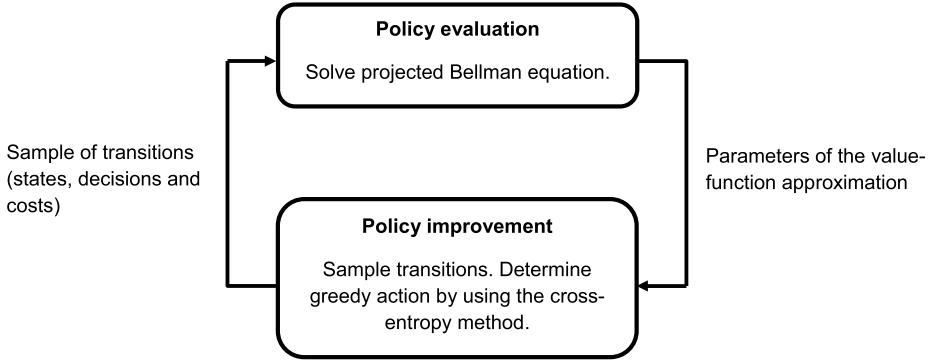


Fig. 1. General scheme of the proposed approximate policy iteration algorithm for the stochastic cutting stock problem.

decisions and costs by simulation. Computing the greedy action is nontrivial, since the decision space is also very large and corresponds to solving a nonconvex nonlinear integer programming problem. We then solve it heuristically by using the cross-entropy method. This sampling step corresponds to the policy improvement step.

Given the collected sample of transitions, the policy evaluation step is accomplished to solving the Bellman equation. As we use a value-function approximation, we solve the projected Bellman equation, whose solution is the fixed point of the Bellman operator projected into the space spanned by the basis functions. Fig. 1 illustrates the general sketch of the proposed algorithm. We detail its development in the next sections.

4.1. Policy evaluation via projected Bellman equation

Instead of working with *state-value functions* $v_\pi(s)$, we will work with *action-value functions* (also known as *q-functions*), defined as

$$q_\pi(s, x) := \lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{t=0}^T \gamma^t c(s_t, x(s_t), d_{t+1}) | s = s_0, x = x(s_0) \right], \quad \forall s \in \mathcal{S}, \forall x \in \mathcal{X},$$

which corresponds to the value of starting at state s , choosing initial action x and then following policy $\pi := \{x(s)\}$ for $t \in \{1, 2, \dots\}$. Action-value functions satisfy a form of Bellman equation, given by

$$q_\pi(s, x) = \mathbb{E}[c(s, x, d) + \gamma q_\pi(s', x(s')) | s, x], \quad \forall s \in \mathcal{S}, \forall x \in \mathcal{X}, \quad (8)$$

in which the expected value is computed with regards to the conditional probability function $\mathbb{P}(d|s, x)$ and $s' = f(s, x, d)$ is the transition function. The use of action-value functions simplifies the computation of a greedy policy, since in this case a greedy action may be determined without knowing the MDP model or computing an expectation. Given a state $s \in \mathcal{S}$, a greedy action $x(s)$ may be determined by

$$x(s) \in \arg \min_{x \in \mathcal{X}_s} q_\pi(s, x).$$

In order to cope with the large state space of the stochastic CSP, we use a linear model to approximate the action-value function of a policy π :

$$q_\pi(s, x; \theta) = \sum_{k=1}^K \phi_k(s, x) \theta_k, \quad (9)$$

in which $\phi_k(s, x)$ are K basis functions, which represent *features* associated with state s and action x , and $\theta = (\theta_1, \theta_2, \dots, \theta_K)$ is a parameter vector. This linear model is in fact a generalization of multiple linear regression—it is *linear in the parameters*, but the basis functions are allowed to be nonlinear. In this way, we can obtain nonlinear approximations by using this model. For example, let s be the current inventory of a single item and x the amount we decide to produce of it. We can then specify a quadratic approximation of the action-value function of a policy as

$$q_\pi(s, x; \theta) = \theta_1 + \theta_2 s^2 + \theta_3 x^2,$$

whose basis functions are $\phi_1(x, s) = 1$, $\phi_2(x, s) = s^2$ and $\phi_3(x, s) = x^2$. Many different basis functions have been used in the literature; we point to [43] for more details on basis functions.

Our objective is then to *train* the model (9)—that is, to estimate its parameters—from samples of state transitions, actions and costs obtained by simulation. Notice though that the exact action-value function does not necessarily lie in the space spanned by the basis functions. Proposed criteria to set the parameters include Bellman error minimization and projected Bellman error minimization [42]. We adopt the projected Bellman error criterion to choose suitable values for the parameters.

Consider the *Bellman operator* $T_\pi : \mathcal{Q} \rightarrow \mathcal{Q}$ associated with a policy $\pi := \{x(s)\}$, in which \mathcal{Q} is the space of action-value functions, defined by

$$\begin{aligned} (T_\pi q)(s, x) &:= \mathbb{E}[c(s, x, d) + \gamma q(s', x(s')) | s, x], \quad \forall s \in \mathcal{S}, \forall x \in \mathcal{X}, \\ &:= \sum_{d \in \mathcal{D}} \mathbb{P}(d | s, x) (c(s, x, d) + \gamma q(s', x(s'))), \end{aligned} \quad (10)$$

in which \mathcal{D} is the support of the probability function $\mathbb{P}(d | s, x)$. Notice that q is an arbitrary action-value function in the space \mathcal{Q} . It will be convenient to write the Bellman operator in vector notation:

$$T_\pi(q) := \bar{c} + \gamma \bar{q},$$

in which $\bar{c} \in \mathbb{R}^{|\mathcal{S} \times \mathcal{X}|}$ is a column-vector, whose components are given by

$$\bar{c}(s, x) = \sum_{d \in \mathcal{D}} \mathbb{P}(d|s, x) c(s, x, d),$$

and $\bar{q} \in \mathbb{R}^{|\mathcal{S} \times \mathcal{X}|}$ is a column-vector, whose components are given by

$$\bar{q}(s, x) = \sum_{d \in \mathcal{D}} \mathbb{P}(d|s, x) q(s', x(s')),$$

and $s' = f(s, x, d)$ is the transition function given by (3).

Now, given an arbitrary policy π , we want to obtain its action-value function q_π . We first notice that, from Bellman equation (8), the action-value function q_π is the fixed point of the Bellman operator (10). However, we cannot solve directly for the fixed point, since we are using an approximation for the action-value function given by (9) and there is no guarantee that q_π lies in the space spanned by the basis functions $\phi_k, k \in \{1, 2, \dots, K\}$. Instead, we solve the *projected Bellman equation* [14]:

$$\Phi \theta = \Pi_\xi T_\pi(\Phi \theta), \quad (11)$$

in which Φ is a matrix with $|\mathcal{S}| \times \sum_{s \in \mathcal{S}} |\mathcal{X}_s|$ rows and K columns whose rows are the feature column-vectors $\phi(s, x)^T = (\phi_1(s, x), \phi_2(s, x), \dots, \phi_K(s, x))$ for each pair (s, x) with $s \in \mathcal{S}$ and $x \in \mathcal{X}_s$:

$$\Phi = \begin{bmatrix} \text{---} \phi(s_1, x_1)^T \text{---} \\ \text{---} \phi(s_1, x_2)^T \text{---} \\ \vdots \\ \text{---} \phi(s_{|\mathcal{S}|}, x_{|\mathcal{X}_s|})^T \text{---} \end{bmatrix},$$

and Π_ξ is a projection matrix to the space spanned by the basis functions in which the projection is defined with relation to a weighted norm $\|\cdot\|_\xi$ and ξ is the weight vector. In the particular case of a Euclidean norm, the projection matrix corresponds to the least squares solution, given by [43]:

$$\Pi_\xi = \Phi(\Phi^T \Xi \Phi)^{-1} \Phi^T \Xi, \quad (12)$$

in which $\Xi = \text{diag}(\xi)$. Substituting (12) in (11), we have

$$\begin{aligned} \Phi \theta &= \Phi(\Phi^T \Xi \Phi)^{-1} \Phi^T \Xi (\bar{c} + \gamma \bar{\Phi} \theta), \\ \Phi^T \Xi \Phi \theta &= \Phi^T \Xi (\bar{c} + \gamma \bar{\Phi} \theta), \\ \Phi^T \Xi (\Phi - \gamma \bar{\Phi}) \theta &= \Phi^T \Xi \bar{c}, \end{aligned} \quad (13)$$

in which $\bar{\Phi}$ is a matrix with $|\mathcal{S}| \times \sum_{s \in \mathcal{S}} |\mathcal{X}_s|$ rows and K columns whose rows $\bar{\phi}(s, x)$ are given by

$$\bar{\phi}(s, x) = \sum_{d \in \mathcal{D}} \mathbb{P}(d|s, x) \phi(s', x(s')), \quad s' = f(s, x, d),$$

so that

$$\bar{\Phi} = \begin{bmatrix} \text{---} \bar{\phi}(s_1, x_1)^T \text{---} \\ \text{---} \bar{\phi}(s_1, x_2)^T \text{---} \\ \vdots \\ \text{---} \bar{\phi}(s_{|S|}, x_{|\mathcal{X}_s|})^T \text{---} \end{bmatrix}.$$

Notice that (13) corresponds to solving a linear system of equations $A\theta = b$ in which $A = \Phi^T \Xi (\Phi - \gamma \bar{\Phi})$ is a $K \times K$ matrix and $b = \Phi^T \Xi \bar{c}$ is a $K \times 1$ column-vector. Notice also that the linear system (13) depends on all states $s \in \mathcal{S}$ and actions $x \in \mathcal{X}$. However, we can solve it approximately if we have a sample of N transitions $\langle s_t, x_t, c_{t+1}, s_{t+1} \rangle$, $t = 0, \dots, N-1$. We then form the approximate matrix \hat{A} and vector \hat{b} :

$$\begin{aligned} \hat{A} &= \frac{1}{N} \sum_{t=0}^{N-1} \phi(s_t, x_t) (\phi(s_t, x_t) - \gamma \phi(s_{t+1}, x(s_{t+1})))^T, \\ \hat{b} &= \frac{1}{N} \sum_{t=0}^{N-1} \phi(s_t, x_t) c_{t+1}, \end{aligned}$$

in which $\phi(s, x)$ is the feature column-vector and transitions are sampled according to a probability distribution corresponding to the weights ξ . Therefore, we can obtain the parameter vector by solving the approximate linear system

$$\hat{A}\theta = \hat{b}. \quad (14)$$

Let $\hat{\theta}$ be a solution to the linear system (14). (We may use the pseudo-inverse if \hat{A} is singular.) Then the action-value function of policy π is approximated by

$$q_\pi(s, x; \hat{\theta}) = \sum_{k=1}^K \phi_k(s, x) \hat{\theta}_k.$$

In summary, evaluating a policy π reduces to solving the linear system (14) from a sample of transitions taken while using policy π . In the next section 4.2, we describe how to obtain an improved policy.

4.2. Policy improvement via the cross-entropy method

Given a parameter vector θ , an improved policy π' relative to a current policy π may be obtained by acting greedily with respect to the approximate action-value function $q_\pi(s, x; \theta)$:

$$x(s) \in \arg \min_{x \in \mathcal{X}_s} \sum_{k=1}^K \phi_k(s, x) \theta_k, \quad \forall s \in \mathcal{S}. \quad (15)$$

Notice that, in the case of the stochastic CSP, (15) will be in general a nonconvex integer nonlinear programming problem. These kind of mathematical programs are among the hardest to solve exactly. Currently available exact solvers are not sufficiently flexible or computationally efficient for our application, since we will have to solve problem (15) thousands of times during the simulation of a sample of transitions to evaluate a single policy. Moreover, solving (15) only approximately may not be detrimental to performance, since we are already working with an approximation of the action-value function.

We use the cross-entropy method to obtain a heuristic solution to (15). The cross-entropy method employs a parametric probability distribution over the space of solutions, which is used to generate a sample of candidate solutions. The objective values of the candidate solutions are computed and a fraction of the best solutions (the elite group) is selected. The elite group is then used to estimate new parameters to the probability distribution. This process is repeated until some stopping criterion is achieved. A detailed explanation of the method is given by [44].

Given a current state s , we generate candidate feasible decisions in the following way: We first sample the total number x_{total} of objects which will be cut from a discrete uniform distribution:

$$x_{\text{total}} \sim \text{DiscUnif}(0, x_{\text{max}}),$$

in which x_{max} is the maximum number of available objects in stock. We then use a multinomial probability distribution to generate a candidate solution x :

$$x \sim \text{Multinomial}(x_{\text{max}}; p_1, p_2, \dots, p_n), \quad (16)$$

in which p_j is the probability of using a cutting pattern $j \in \{1, 2, \dots, n\}$. Feasibility of the sampled candidate is then checked against the set \mathcal{X}_s . Infeasible candidate solutions are rejected.

Let $x^{(1)}, x^{(2)}, \dots, x^{(N)}$ be a sample of candidate solutions taken from (16). We compute the action values $q_\pi(s, x^{(1)}; \theta), q_\pi(s, x^{(2)}; \theta), \dots, q_\pi(s, x^{(N)}; \theta)$ from (9). We then sort the action values in increasing order. Let δ be the $\lceil \rho \times N \rceil$ -th order statistic, with $0 < \rho < 1$, such that the $\lceil \rho \times N \rceil$ best candidate solutions have action values $q_\pi(s, x; \theta) \leq \delta$. (For example $\rho = 0.10$ select the 10% best candidate solutions, also called elite solutions.) We then determine updated parameters $p' = (p'_1, p'_2, \dots, p'_n)$ of the multinomial distribution by maximizing the cross-entropy function, which in this case will be mathematically equivalent to determining the maximum likelihood estimate given only the elite solutions:

$$\max \quad \frac{1}{N} \sum_{i=1}^N I\{q_\pi(s, x^{(i)}; \theta) \leq \delta\} \ell(x^{(i)}; p), \quad (17)$$

s.t.

$$\begin{aligned} \sum_{j=1}^n p_j &= 1, \\ p_j &\geq 0 \quad j \in \{1, 2, \dots, n\}, \end{aligned}$$

in which $I\{\cdot\}$ is an indicator function and $\ell(x^{(i)}; p)$ is the multinomial log-likelihood, given by

$$\ell(x; p) = \sum_{j=1}^n x_j \ln p_j.$$

It can be shown, through the method of Lagrange multipliers, that $p' = (p'_1, p'_2, \dots, p'_n)$ which maximizes (17) is given by

$$p'_j = \frac{\sum_{i=1}^N I\{q_\pi(s, x^{(i)}; \theta) \leq \delta\} x_j^{(i)}}{\sum_{j=1}^n \sum_{i=1}^N I\{q_\pi(s, x^{(i)}; \theta) \leq \delta\} x_j^{(i)}}, \quad j \in \{1, 2, \dots, n\}. \quad (18)$$

Although a bit daunting, Eq. (18) means that p'_j is given by the sample frequency with which cutting pattern j was used by the elite candidate solutions. Initial probabilities are set at $p_j = 1/n, j \in \{1, 2, \dots, n\}$. Algorithm 1 summarizes the steps to obtain a heuristic solution to (15) with the use of the cross-entropy method.

4.3. Approximate policy iteration algorithm

Algorithm 2 describes the proposed approximate policy iteration to obtain a heuristic policy to the stochastic CSP. The algorithm has two loops: in the outer loop, each iteration i corresponds to the evaluation step of a greedy policy $\pi^{(i)}$ relative to the approximate action-value function given by the linear model (9); in the inner loop, a sample of transition is taken by simulating states, decisions and costs. The sample is used to compute the matrix \hat{A} and the vector \hat{b} in order to obtain a new parameter vector $\theta^{(i+1)}$. The outer loop is run for L_1 policy iterations, while the inner loop is run for L_2 iterations, which correspond to the size of the sample of transitions. It returns all computed parameter vectors $\{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(L_1)}\}$. This allows us to reevaluate all generated policies. The algorithm is *off-policy*, i.e., state transitions can be sampled arbitrarily without following the stationary probability distribution induced by the current policy being evaluated.

4.4. Features and basis functions

For the application of Algorithm 2, it is necessary to choose basis functions used in approximating the action-value function of a policy. There are many types of basis

Algorithm 1 Cross-entropy method for obtaining a heuristic greedy action.

```

1: input: Current state  $s$  and  $x_{\max}$ , basis functions  $\phi(\cdot)$ ,
2:   parameters  $\theta$ , algorithm parameters  $N_1, N_2, \rho$ ;
3: initial step: Set  $p_j^{(0)} = 1/n, j \in \{1, 2, \dots, n\}$ ,  $q_{\text{best}} \leftarrow +\infty$ ;
4: for  $k \leftarrow 0 \dots N_1 - 1$  do
5:   for  $i \leftarrow 1 \dots N_2$  do
6:     while True do
7:       Sample  $x_{\text{total}}^{(k,i)} \sim \text{DiscUnif}(0, x_{\max})$ ;
8:       Sample candidate  $x^{(k,i)} \sim \text{Multinomial}(x_{\text{total}}^{(k,i)}; p_1^{(k)}, p_2^{(k)}, \dots, p_n^{(k)})$ ;
9:       if  $x^{(k,i)} \in \mathcal{X}_s$  then
10:        break while;
11:       end if
12:     end while
13:     if  $q_\pi(s, x^{(k,i)}; \theta) < q_{\text{best}}$  then                                 $\triangleright$  Compute  $q_\pi(s, x^{(k,i)}; \theta)$  from equation (9)
14:        $x_{\text{best}} \leftarrow x^{(k,i)}$ ;                                           $\triangleright$  Candidate solution turns to current best solution
15:        $q_{\text{best}} \leftarrow q_\pi(s, x^{(k,i)}; \theta)$ ;
16:     end if
17:   end for
18:   Sort action values  $q_\pi(s, x^{(k,i)}; \theta), i \in \{1, 2, \dots, N_2\}$  in increasing order;
19:   Let  $\delta^{(k)}$  be the  $\lceil \rho \times N_2 \rceil$ -th order statistic in the sorted list of action values;
20:   Update probabilities

```

21: **end for**22: **return** x_{best} \triangleright Best found heuristic greedy action

functions, but in our application we focus on two particular ones: polynomial and Fourier basis functions, which we found to have convenient properties in our case. First, they are easy to specify and, second, they are computationally cheap to compute. This is critical, since we cannot spend too much time to obtain a greedy action when approximately solving (15). Polynomial basis functions are simply polynomial terms in the state s and decision x , while Fourier basis functions are sinusoidal functions such as sines and cosines [45].

In addition, instead of defining the features $\phi(s, x)$ as explicit functions of s and x , we take an intermediate step and aggregate state and decision in a *post-decision* state $s^x = f(s, x)$, in a way that each feature will be defined in terms of $s^x = (s_1^x, s_2^x, \dots, s_m^x)$. In the case of the stochastic CSP, we define as the post-decision state s_i^x related to item i the available inventory, given by the sum of the initial inventory s_i and the quantity of the item resulting from the decision x :

$$s_i^x = s_i + \sum_{j=1}^n a_{ij} x_j. \quad (19)$$

The main advantage of working with the notion of the post-decision state is that it naturally takes into account the relationship between the inventory of items and the cutting patterns, besides being a more compact form which is independent of the number of cutting patterns.

Algorithm 2 Approximate policy iteration for the stochastic CSP.

input: Basis functions $\phi(\cdot)$, $\theta^{(0)}$, γ , L_1 , L_2 ;
1: **for** $i \leftarrow 0 \dots L_1 - 1$ **do**
 Simulate the greedy policy using $\theta^{(i)}$: ▷ Policy improvement
2: $\hat{A}^{(i,0)} \leftarrow 0$, $\hat{b}^{(i,0)} \leftarrow 0$;
3: **for** $t \leftarrow 0 \dots L_2 - 1$ **do**
4: Sample initial inventory $s_t \in \mathcal{S}$;
5: Sample decision $x_t \in \mathcal{X}_{s_t}$;
6: Sample demand $d_{t+1} \sim p(d|s_t, x_t)$;
7: Compute cost $c_{t+1} = c(s_t, x_t, d_{t+1})$ from equation (4);
8: Compute transition $s_{t+1} = f(s_t, x_t, d_{t+1})$ from equation (3);
9: Compute decision $x_{t+1} = x(s_{t+1})$, by solving

$$x(s_{t+1}) \in \arg \min_{x \in \mathcal{X}_{s_{t+1}}} \sum_{k=1}^K \phi_k(s_{t+1}, x) \theta_k^{(i)},$$

 by means of the cross-entropy method (Algorithm 1);
10: Update matrices

$$\hat{A}^{(i,t+1)} \leftarrow \hat{A}^{(i,t)} + \phi(s_t, x_t)(\phi(s_t, x_t) - \gamma \phi(s_{t+1}, x_{t+1}))^T,$$

$$\hat{b}^{(i,t+1)} \leftarrow \hat{b}^{(i,t)} + \phi(s_t, x_t) c_{t+1};$$

11: **end for**
 Solve projected Bellman equation: ▷ Policy evaluation
12: Determine $\theta^{(i+1)}$ by solving $\hat{A}^{(i,L_2)} \theta = \hat{b}^{(i,L_2)}$;
 (Use the pseudo-inverse if $\hat{A}^{(i,L_2)}$ is singular.)
13:
14: **end for**
15: **return** $\{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(L_1)}\}$

In the case of a polynomial basis, the action-value function may be written as [8]:

$$q_{\pi}^{\text{poly}}(s, x; \theta) = \sum_{k=1}^K \prod_{i=1}^m (s_i^x)^{c_{ki}} \theta_k,$$

in which c_{ki} are integers which specify the order of the polynomial features.

Regarding the Fourier basis, it uses as basis functions the terms of the Fourier series, which is commonly used to approximate continuous functions. In our case, we will use a Fourier basis with cosine functions only, which according to [45] is sufficient to approximate any function defined in a non-negative domain, given by

$$q_{\pi}^{\text{Fourier}}(s, x; \theta) = \sum_{k=1}^K \cos(\pi c_k \cdot s^{\text{norm}}) \theta_k,$$

in which $c_k = (c_1, c_2, \dots, c_m)$ is a vector of integers which specify the frequencies of the cosine functions. Larger integer numbers generate cosine functions with higher frequencies. In addition, we use normalized post-decision states, such that $s^{\text{norm}} = (s_1^{\text{norm}}, s_2^{\text{norm}}, \dots, s_m^{\text{norm}})$ and $s_i^{\text{norm}} = s_i^x / s_{\max}$.

5. Results and discussion

In this section, we report the application of Algorithm 2 with the use of a realistic dataset. Our objective was to evaluate whether our approximate policy iteration algorithm can generate a decision policy for the stochastic CSP with acceptable performance, which we mean as a policy that is able of meeting the demand without maintaining excessive inventory levels. In other words, our *agent* has to learn which cutting patterns to use over time so that it maintains sufficient and not very high inventory levels. We compared the obtained heuristic policy with myopic and random policies.

5.1. A myopic policy

We benchmarked the trained policies with a myopic policy which is a reasonable heuristic to the stochastic CSP in practice. The myopic policy decides which items to cut at each time period according to the expected demand. Then, at each time period, the myopic policy solves a deterministic CSP model so that the sum of the current initial inventory at hand and the items cut is greater than the expected demand at minimal trim loss. Algorithm 3 describes the myopic policy, in which $\bar{d} = (\bar{d}_1, \bar{d}_2, \dots, \bar{d}_m)$ is the expected demand.

Algorithm 3 A myopic policy for the stochastic CSP.

input: Expected demand vector \bar{d} , cutting patterns a_j , trim loss costs g_j , $\forall j \in \{1, 2, \dots, n\}$;
1: Sample initial inventory $s_t \in \mathcal{S}$;
2: **for** $t \leftarrow 0 \dots$ **do**
3: Compute decision x_t by solving the following integer linear program:

$$\begin{aligned} \min \quad & \sum_{j=1}^n g_j x_{jt} \\ \text{s.t.} \quad & \sum_{j=1}^n a_{ij} x_{jt} + s_{it} \geq \bar{d}_i, \quad i \in \{1, 2, \dots, m\}, \\ & x_{jt} \in \mathbb{Z}_+; \end{aligned}$$

4: Sample demand $d_{t+1} \sim p(d|s_t, x_t)$;
5: Compute cost $c_{t+1} = c(s_t, x_t, d_{t+1})$ from equation (4);
6: Compute transition $s_{t+1} = f(s_t, x_t, d_{t+1})$ from equation (3);
7: **end for**

5.2. Dataset

In the experiments, we used data originating from a real enterprise which faces the problem of cutting steel bars for building construction. Seven different lengths of steel bars (Table 1) may be cut from stock bars with length 1500 cm. Demand for each bar length at each time period is stationary but random, so that the enterprise cannot predict

Table 1

Lengths of demanded item types. Length of stock objects is 1500 cm.

Item	1	2	3	4	5	6	7
Length (cm)	115	180	267	314	880	1180	1200

Table 2

Cutting patterns used in the numerical experiments. Each column shows a different cutting pattern with associated trim losses.

Item	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12	P13	P14	P15
1	10	13	3	3	2	2	1	1	0	0	0	0	0	0	0
2	0	0	1	1	2	0	1	1	0	0	0	0	0	0	1
3	0	0	0	0	0	0	0	0	1	1	1	2	2	3	0
4	1	0	0	3	0	0	0	0	0	0	1	0	3	2	4
5	0	0	1	0	1	0	0	0	0	0	1	1	0	0	0
6	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0
7	0	0	0	0	0	1	1	0	1	0	0	0	0	0	0
Trim loss (cm)	36	5	95	33	30	70	5	25	33	53	39	86	24	71	64

exactly how much of each bar length customers will order. In this particular case, a time period corresponds to one week. The enterprise wants to decide in anticipation which bar lengths to cut in order to meet demand without incurring high trim loss cost or maintaining high inventory.

The number of possible cutting patterns in this problem is large (all feasible solutions to an integer knapsack problem with length 1500 cm and 7 item types). From a practical perspective, enterprises often restrict the number of cutting patterns to a small manageable set of *good patterns* among efficient ones, i.e., those with small trim loss. Other selection criteria depend on particular factors of the manufacturing process, such as easiness of setup and ergonomics. We then handcrafted 15 cutting patterns given in Table 2, and as such they are not optimized to the enterprise operations. (Patterns used in practice were not disclosed by the enterprise.)

We assumed the demand follows a mixture distribution in which the vector $d = (d_1, d_2, \dots, d_m)$ of demands for items follow a multinomial distribution $d \sim \text{multinomial}(d_{\text{total}}, p)$ conditional on the total demand $d_{\text{total}} \sim \text{DiscUnif}(d_{\min}, d_{\max})$, which follows a discrete uniform probability distribution and $p = (p_1, p_2, \dots, p_m)$ is the vector of probabilities of each item being demanded. Thus, the marginal distribution of the vector d is

$$\begin{aligned}
 \mathbb{P}(d) &= \sum_{d_{\text{total}}=d_{\min}}^{d_{\max}} \mathbb{P}(d|d_{\text{total}}) \mathbb{P}(d_{\text{total}}) \\
 &= \sum_{d_{\text{total}}=d_{\min}}^{d_{\max}} \text{Multinomial}(d_{\text{total}}, p) \times \text{DiscUnif}(d_{\min}, d_{\max}). \quad (20)
 \end{aligned}$$

Samples from the probability distribution of the demand given in Eq. (20) were simulated by sampling the total demand $d_{\text{total}} \sim \text{DiscUnif}(d_{\min}, d_{\max})$ then sampling

Table 3

Probability distribution of the demand for items.

Item	1	2	3	4	5	6	7	d_{\min}	d_{\max}
Probability (p)	0.30	0.20	0.20	0.10	0.10	0.05	0.05	40	50

Table 4

Additional problem data.

Parameter	Value	Description
h_i^+	$0.01l_i$	Unit cost of holding inventory, where l_i is the length of item $i \in \{1, 2, \dots, 7\}$
h_1^-	$1.0l_i$	Penalty unit cost of not meeting the demand
g_j	$0.1g_j^+$	Trim loss cost, in which g_j^+ is the trim loss of pattern $j \in \{1, 2, \dots, 15\}$
s_{\max}	70	Maximum inventory for each item in a time period
x_{\max}	30	Number of stock objects available in each time period

Table 5

Parameters used in the training stage of the numerical experiments (Algorithm 2).

Parameter	Value	Description
γ	0.8	Discount factor
L_1	30	Number of policy iterations
L_2	50×10^3	Number of simulated state transitions
N_1	10	Number of iteration in the cross-entropy method
N_2	100	Sample size of candidate solutions in the cross-entropy method

$d \sim \text{Multinomial}(d_{\text{total}}, p)$. Table 3 shows the parameters of the probability function of the demand. Other parameters related to the problem are given in Table 4.

5.3. Experiment results

We implemented the algorithms in Python 3.7 with the aid of NumPy and Numba and the experiments were run in a Core i7 7700K machine with 8 GB RAM. The experiments were divided in a *training* stage and a *testing* (or application) stage. In the training stage, we applied Algorithm 2 to tune the parameters $\theta = (\theta_1, \theta_2, \dots, \theta_K)$ of our policy—these are the parameters of the basis-function model which approximates the value function of a policy (see Eq. (15)). This stage was carried out offline with the use of simulation. After that, in the testing stage, we fixed the parameter θ obtained in the training stage and applied the trained policies to a new sample path of demands.

Table 5 shows the parameter values used in the training experiments. Initial values for $\theta_k^{(0)}, k \in 1, 2, \dots, K$ were sampled from a Gaussian distribution with mean 0 and standard deviation 1. We set a sample size of 50×10^3 state transitions in the evaluation of a policy; this sample size was chosen based on trial and error during preliminary experiments. Algorithm 2 was run for 30 iterations, amounting to a total sample size of 1.5 million state transitions. The training stage took approximately 20 hours with our computational resources. (Obviously, this time can be reduced with more powerful computers.)

Due to simulation error and approximation of the action-value functions, the successive policies generated during training did not always exhibit monotonically better

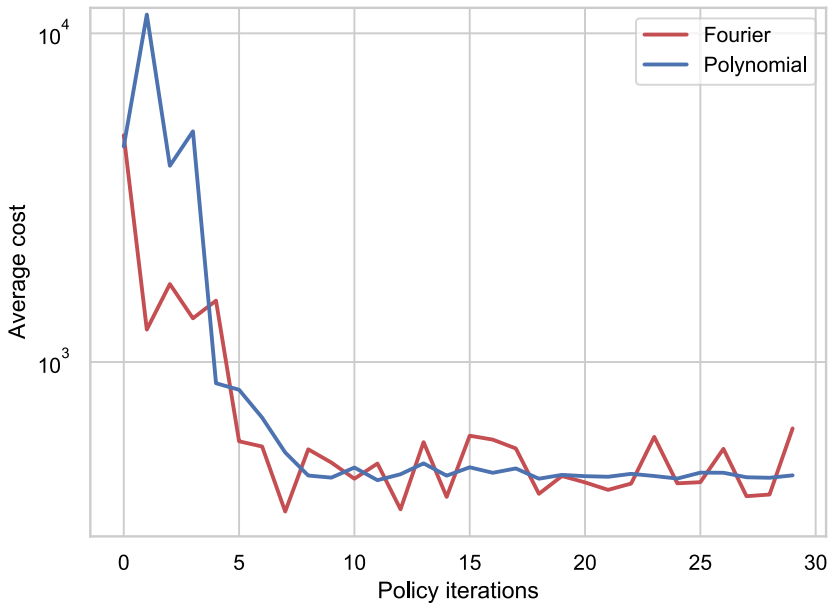


Fig. 2. Average cost (log scale) of heuristic policies trained by our algorithm, evaluated with 10 simulation replications after training. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

performance. Then, after training, we carried out a reevaluation of each of the 30 policies with 10 simulation replications in order to identify the best policy generated. Fig. 2 exhibits the performance of the successive policies after reevaluation, and it can be seen that policies improved in the first iterations and then started to oscillate. This policy oscillation phenomenon has been documented early in the literature from experimental studies in different applications of approximate policy iteration algorithms [14].

In the testing stage, trained policies were compared with the myopic policy described in Algorithm 3 and a random policy, in which decisions are sampled uniformly in the set \mathcal{X}_s at each time step. This testing stage is considerably faster in comparison with the training stage, taking a few seconds at each time step, since it involves only solving an optimization problem by using Algorithm 1. Thus, the computational times of both our proposed policies and the myopic policy are comparable in this testing stage.

We simulated the demand for 1000 time steps—i.e., a sample path of 1000 demand values—and applied each policy to this same sample path, so that variability due to sampling was controlled and variability in results was only due to the difference in policies. We further replicated the experiment 10 times and the average cost of the policies over the 10 replications are shown in Fig. 3.

Trained policies achieved better performance than a myopic policy, with the policy trained using the Fourier basis showing lower average cost than the policy trained using the polynomial basis. Final average cost for Fourier and polynomial basis policies is 363.4 and 441.3, respectively, while for the myopic policy is 2186.5 and for the random policy

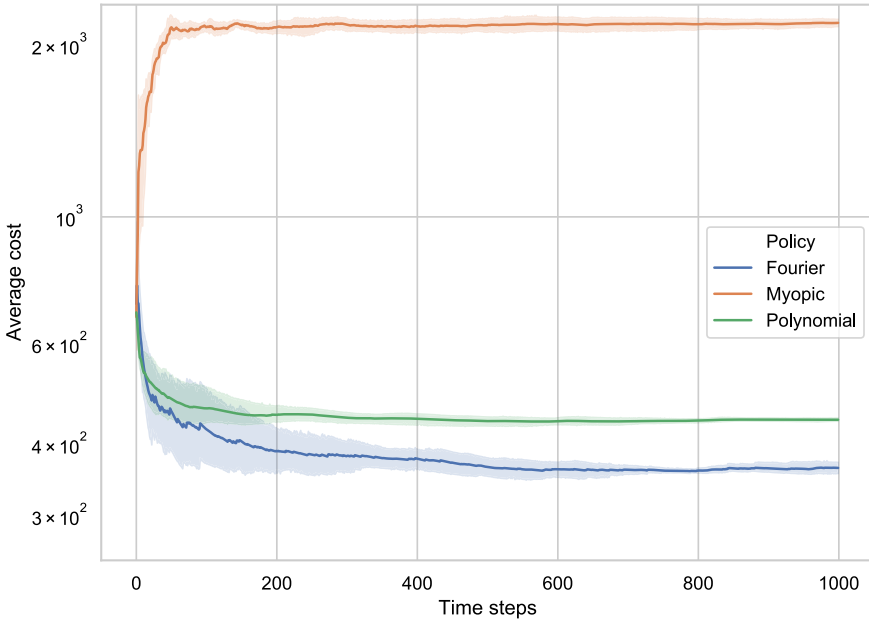


Fig. 3. Average cost of policies (log scale). Bands are 95% bootstrapped confidence intervals for a sample of 10 experiment replications. Final average cost for Fourier and polynomial basis policies is 363.4 and 441.3, respectively, while for the myopic policy is 2186.5. The cost of the random policy is 6955.8 (omitted in the graph).

is 6955.8. This represents an improvement in performance of approximately 80% over the myopic policy.

Fig. 4 shows the initial inventory—the state before taking a decision and observing the demand—at each time step for each policy applied in the testing stage. (Decisions x_t taken on how many cutting patterns to use at each time step are not shown in this figure.) Each line in the graph represents the inventory level $s_i, i = 1, 2, \dots, 7$ of each item type in Table 1. At time $t = 0$, the inventory for all items is equal to 20. Over the 300 time steps, each policy decides on how many objects in stock to cut at each available pattern in Table 2 given the current inventory. As a result, different amounts of each item type are cut and added to the inventory, which is then depleted by the realization of the demand.

We can see in Fig. 4 that both the myopic and random policies fail to adequately control inventory levels, maintaining very low levels for some items and very high levels for others. We suppose this is due to the myopic nature of these policies, which do not take into account the impact of decisions in the future states. (In fact, the random policy does not even look at the current inventory, while the myopic policy considers information of the current inventory and the expected demand.) In contrast, the trained policies using the Fourier and polynomial basis functions maintain sufficient levels to meet the demand while avoiding that the levels get high; however, as can be seen in

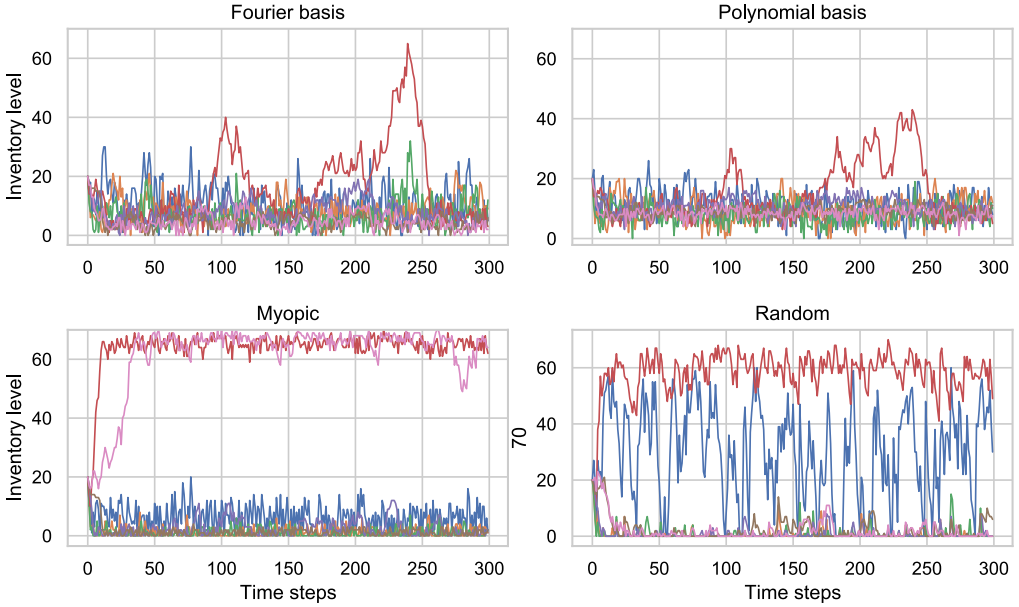


Fig. 4. Initial inventory (state before decision) of all items at each time step. The myopic policy fails to adequately control inventory levels, while the trained policies maintain low inventory levels.

Fig. 4, sometimes inventory levels for some items may digress to high levels before the policy manage to get them back to low levels.

Although both policies trained with a Fourier basis and a polynomial basis action-value function approximation exhibited lower average cost than a myopic policy, Fourier basis performed better. We can get some insight on this difference in performance when we look at the available inventory set by both policies over time. The available inventory is the post-decision state given by Eq. (19), i.e., it is the sum of the initial inventory at a time period and the items produced by cutting stock objects. Fig. 5 shows the available inventory for items 1 and 7. Notice that the policy trained with the Fourier basis maintains lower available inventory while still sufficient to fulfill the demand. Interestingly, the policy with the Fourier basis function approximation seems to have learned a *safety* inventory level of 5 units for item 7, which has very low demand.

Finally, we have also investigated the effect of using different values for the discount factor γ , which works as a hyperparameter of the model. It is not obvious at first which value to assign to the discount factor and its effect on the performance of the trained policy. We have run experiments with the same dataset and 6 different values for the discount factor with a Fourier basis function approximation, with 10 simulation replications for each discount factor value. The results are given in Fig. 6. It can be seen that the performance of the trained policy is very sensitive to different values of the discount factor. Lower levels (0.5 and 0.6) have considerably worse average cost than higher values, with lowest average cost at value 0.8. This indicates that there may be an optimal

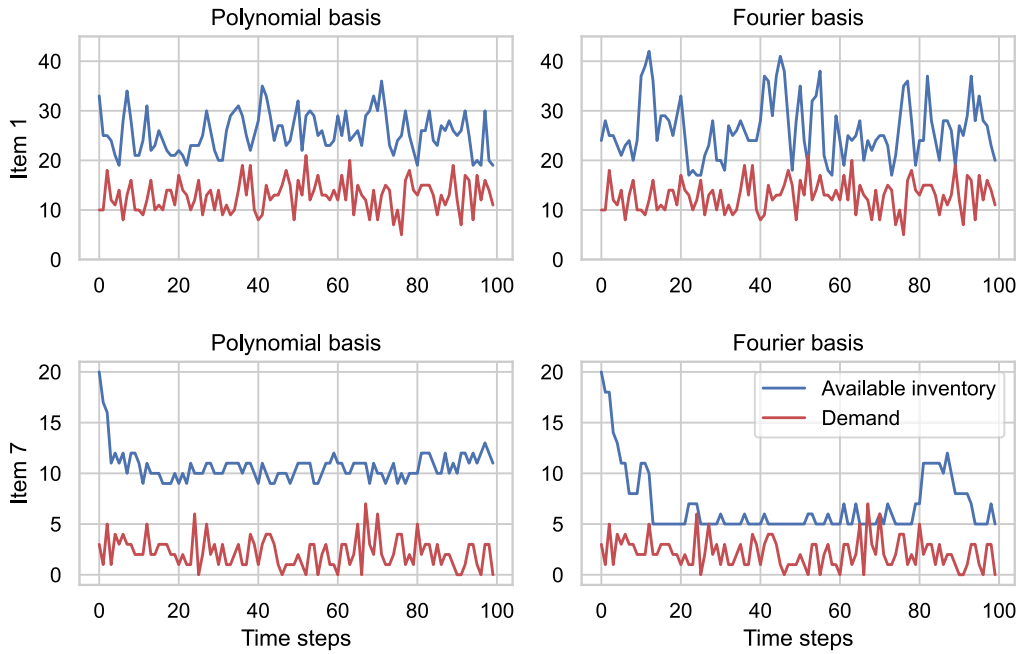


Fig. 5. Available inventory (sum of initial inventory and cut items) generated by the trained policies, and demand for items at each time step. It can be seen that the trained policies maintain sufficient available inventory in order to fulfill the demand.

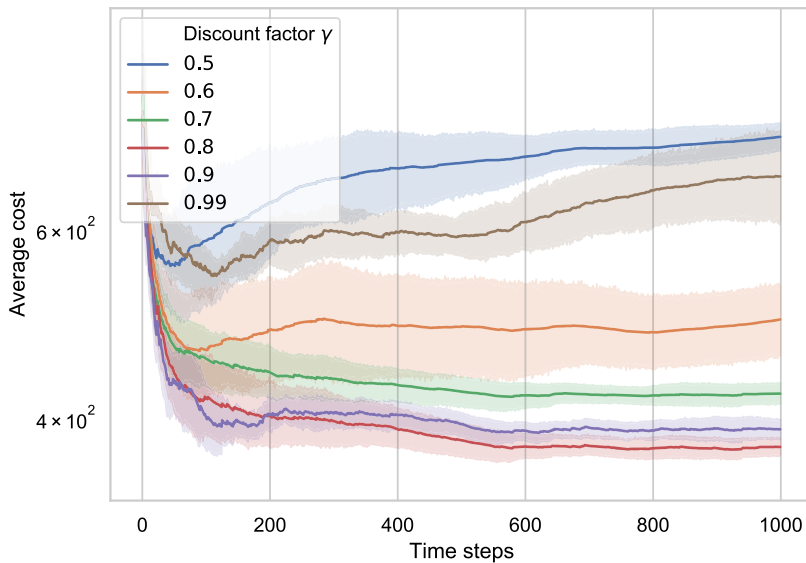


Fig. 6. Performance of trained policies with the Fourier basis for different values of the discount factor γ . Performance is critically affected by the choice of discount factor. Bands represent 95% confidence intervals.

discount factor which results in the best performance. In practical applications, the best discount factor may be searched by using an algorithm for hyperparameter optimization.

6. Conclusions

In this paper, we developed a solution approach to the stochastic cutting stock problem based on reinforcement learning. In the numerical experiments with realistic data, the obtained heuristic policies showed a considerable performance improvement of up to 80% over a myopic policy which repeatedly solves a deterministic formulation of the cutting stock problem. These promising results provide experimental evidence that a decision system for the cutting stock problem based on reinforcement learning may result in considerable cost reductions in industry, specially when the demand is uncertain.

Although we focused on the one-dimensional variant of the cutting stock problem, our pattern-based formulation may be adapted to other variants, including two- or three-dimensional ones; particular features of a variant may be embedded in the cutting patterns and the general mathematical structure of the problem remains approximately the same.

Finally, more studies are necessary to evaluate the quality and stability of policies under different conditions, such as nonstationary demand and varying number of items or cutting patterns. We also envisage as possible research directions the investigation of nonlinear approximators such as artificial neural networks and the application of hyperparameter optimization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

Funding: This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior – Brazil (CAPES) – Finance Code 001; and Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq; Grant No.: 422464/2016-3). We also thank NVIDIA Corporation for the GPU support.

References

- [1] P.C. Gilmore, R.E. Gomory, A linear programming approach to the cutting-stock problem, *Oper. Res.* 9 (1961) 849–859.
- [2] P.C. Gilmore, R.E. Gomory, A linear programming approach to the cutting stock problem—part ii, *Oper. Res.* 11 (1963) 863–888, <https://doi.org/10.1287/opre.11.6.863>.
- [3] G. Wäscher, H. Haußner, H. Schumann, An improved typology of cutting and packing problems, *Eur. J. Oper. Res.* 183 (2007) 1109–1130, <https://doi.org/10.1016/j.ejor.2005.12.047>.

- [4] S.C. Leung, D. Zhang, A fast layer-based heuristic for non-guillotine strip packing, *Expert Syst. Appl.* 38 (2011) 13032–13042, <https://doi.org/10.1016/j.eswa.2011.04.105>.
- [5] A.M. Del Valle, T.A. de Queiroz, F.K. Miyazawa, E.C. Xavier, Heuristics for two-dimensional knapsack and cutting stock problems with items of irregular shape, *Expert Syst. Appl.* 39 (2012) 12589–12598, <https://doi.org/10.1016/j.eswa.2012.05.025>.
- [6] A. Fernández, C. Gil, R. Baños, M.G. Montoya, A parallel multi-objective algorithm for two-dimensional bin packing with rotations and load balancing, *Expert Syst. Appl.* 40 (2013) 5169–5180, <https://doi.org/10.1016/j.eswa.2013.03.015>.
- [7] M.L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, John Wiley & Sons, Inc., USA, 1994.
- [8] R. Sutton, A. Barto, *Reinforcement Learning: An Introduction*, 2nd ed., MIT Press, 2018.
- [9] V. Mnih, K. Kavukcuoglu, D. Silver, A.A. Rusu, J. Veness, M.G. Bellemare, A. Graves, M. Riedmiller, A.K. Fidjeland, G. Ostrovski, et al., Human-level control through deep reinforcement learning, *Nature* 518 (2015) 529.
- [10] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al., A general reinforcement learning algorithm that masters chess, shogi, and go through self-play, *Science* 362 (2018) 1140–1144.
- [11] S. Gu, E. Holly, T. Lillicrap, S. Levine, Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates, in: 2017 IEEE International Conference on Robotics and Automation (ICRA), 2017, pp. 3389–3396.
- [12] Q. Liu, Z. Liu, B. Xiong, W. Xu, Y. Liu, Deep reinforcement learning-based safe interaction for industrial human–robot collaboration using intrinsic reward function, *Adv. Eng. Inform.* 49 (2021) 101360, <https://doi.org/10.1016/j.aei.2021.101360>.
- [13] W. Powell, *Approximate Dynamic Programming: Solving the Curses of Dimensionality*, 2nd ed., Wiley, 2011.
- [14] D. Bertsekas, J. Tsitsiklis, *Neuro-Dynamic Programming*, Athena Scientific, 1996.
- [15] M. Aslani, S. Seipel, M.S. Mesgari, M. Wiering, Traffic signal optimization through discrete and continuous reinforcement learning with robustness analysis in downtown Tehran, *Adv. Eng. Inform.* 38 (2018) 639–655, <https://doi.org/10.1016/j.aei.2018.08.002>.
- [16] L. Al-Kanj, J. Nascimento, W.B. Powell, Approximate dynamic programming for planning a ride-hailing system using autonomous fleets of electric vehicles, *Eur. J. Oper. Res.* (2020), <https://doi.org/10.1016/j.ejor.2020.01.033>.
- [17] M.A. Lopes Silva, S.R. de Souza, M.J. Freitas Souza, A.L.C. Bazzan, A reinforcement learning-based multi-agent framework applied for solving routing and scheduling problems, *Expert Syst. Appl.* 131 (2019) 148–171, <https://doi.org/10.1016/j.eswa.2019.04.056>.
- [18] C.-L. Liu, C.-C. Chang, C.-J. Tseng, Actor-critic deep reinforcement learning for solving job shop scheduling problems, *IEEE Access* 8 (2020) 71752–71762.
- [19] H. Wang, Q. Yan, S. Zhang, Integrated scheduling and flexible maintenance in deteriorating multi-state single machine system using a reinforcement learning approach, *Adv. Eng. Inform.* 49 (2021) 101339, <https://doi.org/10.1016/j.aei.2021.101339>.
- [20] D.R. Jiang, T.V. Pham, W.B. Powell, D.F. Salas, W.R. Scott, A comparison of approximate dynamic programming techniques on benchmark energy storage problems: does anything work?, in: 2014 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL), 2014, pp. 1–8.
- [21] F. Schneider, U.W. Thonemann, D. Klabjan, Optimization of battery charging and purchasing at electric vehicle battery swap stations, *Transp. Sci.* 52 (2018) 1211–1234, <https://doi.org/10.1287/trsc.2017.0781>.
- [22] M.S. Maxwell, S.G. Henderson, H. Topaloglu, Tuning approximate dynamic programming policies for ambulance redeployment via direct search, *Stoch. Syst.* 3 (2013) 322–361, <https://doi.org/10.1287/10-SSY020>.
- [23] A. Kara, I. Dogan, Reinforcement learning approaches for specifying ordering policies of perishable inventory systems, *Expert Syst. Appl.* 91 (2018) 150–158, <https://doi.org/10.1016/j.eswa.2017.08.046>.
- [24] R. Pourmoayed, L.R. Nielsen, An approximate dynamic programming approach for sequential pig marketing decisions at herd level, *Eur. J. Oper. Res.* 276 (2019) 1056–1070, <https://doi.org/10.1016/j.ejor.2019.01.050>.
- [25] U. Abdulwahab, M. Wahab, Approximate dynamic programming modeling for a typical blood platelet bank, *Comput. Ind. Eng.* 78 (2014) 259–270, <https://doi.org/10.1016/j.cie.2014.07.017>.

- [26] P. Trkman, M. Gradisar, One-dimensional cutting stock optimization in consecutive time periods, *Eur. J. Oper. Res.* 179 (2007) 291–301, <https://doi.org/10.1016/j.ejor.2006.03.027>.
- [27] H. Reinertsen, T. Vossen, The one-dimensional cutting stock problem with due dates, *Eur. J. Oper. Res.* 201 (2010) 701–711, <https://doi.org/10.1016/j.ejor.2009.03.042>.
- [28] C. Arbib, F. Marinelli, On cutting stock with due dates, *Omega* 46 (2014) 11–20, <https://doi.org/10.1016/j.omega.2014.01.004>.
- [29] B.A. Prata, A.R. Pitombeira-Neto, C.J.M. Sales, An integer linear programming model for the multiperiod production planning of precast concrete beams, *J. Constr. Eng. Manage.* 141 (2015) 04015029, [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000991](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000991).
- [30] A.R. Pitombeira-Neto, B.d.A. Prata, A matheuristic algorithm for the one-dimensional cutting stock and scheduling problem with heterogeneous orders, *TOP* 28 (2019) 178–192, <https://doi.org/10.1007/s11750-019-00531-3>.
- [31] S.L. Nonås, A. Thorstenson, A combined cutting-stock and lot-sizing problem, *Eur. J. Oper. Res.* 120 (2000) 327–342, [https://doi.org/10.1016/S0377-2217\(99\)00160-5](https://doi.org/10.1016/S0377-2217(99)00160-5).
- [32] K.C. Poldi, S.A. de Araujo, Mathematical models and a heuristic method for the multiperiod one-dimensional cutting stock problem, *Ann. Oper. Res.* 238 (2016) 497–520, <https://doi.org/10.1007/s10479-015-2103-2>.
- [33] G.M. Melega, S.A. de Araujo, R. Jans, Classification and literature review of integrated lot-sizing and cutting stock problems, *Eur. J. Oper. Res.* 271 (2018) 1–19, <https://doi.org/10.1016/j.ejor.2018.01.002>.
- [34] M.C.N. Gramani, P.M. França, The combined cutting stock and lot-sizing problem in industrial processes, *Eur. J. Oper. Res.* 174 (2006) 509–521, <https://doi.org/10.1016/j.ejor.2004.12.019>.
- [35] B. Durak, D.T. Aksu, Dynamic programming and mixed integer programming based algorithms for the online glass cutting problem with defects and production targets, *Int. J. Prod. Res.* 55 (2017) 7398–7411, <https://doi.org/10.1080/00207543.2017.1349951>.
- [36] D. Sculli, A stochastic cutting stock procedure: cutting rolls of insulating tape, *Manag. Sci.* 27 (1981) 946–952, <https://doi.org/10.1287/mnsc.27.8.946>.
- [37] E.V. Krichagina, R. Rubio, M.I. Taksar, L.M. Wein, A dynamic stochastic stock-cutting problem, *Oper. Res.* 46 (1998) 690–701, <https://doi.org/10.1287/opre.46.5.690>.
- [38] D.J. Alem, P.A. Munari, M.N. Arenales, P.A.V. Ferreira, On the cutting stock problem under stochastic demand, *Ann. Oper. Res.* 179 (2010), <https://doi.org/10.1007/s10479-008-0454-7>.
- [39] P. Beraldi, M. Bruni, D. Conforti, The stochastic trim-loss problem, *Eur. J. Oper. Res.* 197 (2009) 42–49, <https://doi.org/10.1016/j.ejor.2008.04.042>.
- [40] A. Zanarini, Optimal stock sizing in a cutting stock problem with stochastic demands, in: D. Salvagnin, M. Lombardi (Eds.), in: *Lecture Notes in Computer Science*, vol. 10335, Springer, Cham, 2017.
- [41] R.A. Howard, *Dynamic Programming and Markov Processes*, John Wiley, 1960.
- [42] M.G. Lagoudakis, R. Parr, Least-squares policy iteration, *J. Mach. Learn. Res.* 4 (2003) 1107–1149.
- [43] A. Geramifard, T.J. Walsh, S. Tellex, *A Tutorial on Linear Function Approximators for Dynamic Programming and Reinforcement Learning*, Now Publishers Inc., Hanover, MA, USA, 2013.
- [44] P.-T. De Boer, D.P. Kroese, S. Mannor, R.Y. Rubinstein, A tutorial on the cross-entropy method, *Ann. Oper. Res.* 134 (2005) 19–67.
- [45] G. Konidaris, S. Osentoski, P. Thomas, Value function approximation in reinforcement learning using the Fourier basis, in: *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*, AAAI Press, 2011, pp. 380–385.