

Solving Hard Bi-objective Knapsack Problems Using Deep Reinforcement Learning

Hadi Charkhgard^a, Hanieh Rastegar Moghaddam^a, Ali Eshragh^b, Sasan Mahmoudinazlou^a

^a*Department of Industrial and Management Systems Engineering, University of South Florida, Tampa, FL, USA*

^b*Carey Business School, Johns Hopkins University, Washington, D.C., USA*

Abstract

We study a class of bi-objective integer programs known as bi-objective knapsack problems (BOKPs). Our research focuses on the development of innovative exact and approximate solution methods for BOKPs by synergizing algorithmic concepts from two distinct domains: multi-objective integer programming and (deep) reinforcement learning. While novel reinforcement learning techniques have been applied successfully to single-objective integer programming in recent years, a corresponding body of work is yet to be explored in the field of multi-objective integer programming. This study is an effort to bridge this existing gap in the literature. Through a computational study, we demonstrate that although it is feasible to develop exact reinforcement learning-based methods for solving BOKPs, they come with significant computational costs. Consequently, we recommend an alternative research direction: approximating the entire nondominated frontier using deep reinforcement learning-based methods. We introduce two such methods, which extend classical methods from the multi-objective integer programming literature, and illustrate their ability to rapidly produce high-quality approximations.

Keywords: Multi-objective optimization; Bi-objective Knapsack Problem; Deep Reinforcement Learning; Q -learning

1. Introduction

Multi-objective optimization is an active area of research that primarily focuses on achieving a comprehensive understanding of the trade-offs between conflicting objective functions of an optimization problem. Its main aim is to identify Pareto-optimal solutions, which are solutions where it is impossible to improve one objective without adversely affecting at least one other objective. Over the past two decades, significant attention has been directed toward the development of numerous exact and heuristic solution methods designed to address a specific class of multi-objective optimization problems known as multi-objective integer linear programs [4, 5, 8, 9, 11, 14, 18, 20, 27, 28]. In this class of optimization problems, all constraints and objective functions are linear, and all decision variables are constrained to integer values exclusively.

The heightened interest in this area is, in part, a result of the substantial advancements in single-objective optimization solvers and the fact that a vast number of practical problems can be formulated as multi-objective integer programs. While substantial progress has been made in algorithmic advancements within multi-objective integer programming, there remain several underexplored research areas, primarily driven by the rapid evolution of machine learning. This trend is exemplified by the recent surge in efforts to apply machine learning techniques to enhance

1 single-objective integer linear programming methods [3]. In contrast, such endeavors have been
2 relatively limited in the field of multi-objective integer linear programming. This discrepancy is
3 not surprising, considering that similar initiatives typically reach the multi-objective optimization
4 domain with some delay due to its inherently greater complexity when compared to its single-
5 objective optimization counterparts. In fact, even the existing, albeit limited, efforts have primarily
6 been in the opposite direction, i.e., utilizing multi-objective optimization techniques to improve
7 machine learning methods [1, 23].

8 The gap becomes even more pronounced within a prominent machine learning paradigm known
9 as (deep) Reinforcement Learning (RL), which naturally aligns with sequential single-objective
10 decision-making problems involving uncertainty. Deep RL emerged approximately a decade ago,
11 demonstrating its efficacy in conquering Atari games [15]. Several years after its emergence, re-
12 searchers began to explore the prospect of employing RL to solve specific classes of single-objective
13 integer programming problems without uncertainty. Such optimization problems need to be trans-
14 formed into a sequential decision-making problem before being solved. One of the initial endeavors
15 in this direction was undertaken by [16], who sought to solve vehicle routing problems. Since then,
16 numerous studies have explored this research avenue [12, 21]. However, comparable endeavors in
17 the domain of multi-objective integer programming have been scarce, and this study aims to bridge
18 this gap.

19 In the context of single-objective integer linear programming, the absence of uncertainty presents
20 a considerable challenge when designing deep RL methods: Achieving superior performance com-
21 pared to state-of-the-art ‘custom-built’ heuristics in this domain can be hard, i.e., evaluations
22 typically involve comparisons with other RL-based methods or general-purpose solution methods.
23 Consequently, existing studies have predominantly focused on evaluating the effectiveness of these
24 methods in solving single-objective integer linear programming problems rather than their compet-
25 itive edge [10, 13]. The hope is that as this research domain matures, deep RL methods can either
26 stand on their own as competitive methods or contribute to the development of hybrid methods
27 that exhibit competitiveness. A similar trajectory is anticipated in the domain of multi-objective
28 integer linear programming.

29 With this consideration, in this study, we investigate how fundamental ideas from both multi-
30 objective optimization and reinforcement learning can be effectively combined to solve a specific
31 class of multi-objective integer programs known as Bi-objective Knapsack Problems (BOKPs).
32 Our aim is to present our proposed methods in a clear and understandable manner, making them
33 accessible even to readers who may not be familiar with multi-objective integer programming and/or
34 RL. This will also help readers envision potential directions for future research. To achieve this goal,
35 we begin by demonstrating the development of an exact RL-based method by blending a customized
36 adaptation of the ε -constraint method from multi-objective optimization with Q -learning from the
37 RL field. Unfortunately, as we show through a small computational experiment, exact RL-based
38 methods tend to be computationally slow and require individualized training for each problem
39 instance. Nevertheless, this detailed presentation of the exact RL-based approach serves as a
40 foundation for introducing the essential principles required for our proposed approximate RL-based
41 techniques.

42 Specifically, we introduce two approximate RL-based methods that enable the effective genera-
43 tion of a single deep neural network for approximating nondominated frontiers instead of training
44 many networks to produce points from different regions in the criterion space [25]. The first method
45 is an extension of the proposed exact RL-based approach which combines a customized variation of

1 the ε -constraint method with the *deep Q*-learning. The second method combines deep *Q*-learning
2 with the weighted sum method. A notable advantage of the proposed approximate RL-based
3 methods is their ability to be trained for a class of instances without the need for instance-specific
4 training. As mentioned earlier, it is not practical to benchmark RL-based methods against state-
5 of-the-art solutions, as these methods are still in their early stages of development. However, we
6 can illustrate their potential. To demonstrate this potential, we focus on a hard class of BOKPs,
7 which are known to pose significant challenges for existing solution methods even when there is
8 only a single objective [26, 22]. These hard BOKPs serve as an ideal testing ground for evaluating
9 the effectiveness of our proposed approximate RL-based methods. In a computational analysis, we
10 first demonstrate that our methods provide high-quality approximations for small-sized instances,
11 where the exact nondominated frontier can be computed. For larger instances, we compare the
12 performance of our methods against a state-of-the-art heuristic [30], and show that it can rapidly
13 generate better approximations.

14 The rest of this paper is organized as follows: In Section 2, we provide some preliminary
15 concepts. Section 3 elaborates on the proposed exact RL-based method. Section 4 provides an
16 in-depth explanation of the proposed approximate RL-based methods. Section 5 encompasses a
17 computational study. Finally, in Section 6, we offer concluding remarks.

18 2. Preliminaries

19 A biobjective optimization problem can be stated as follows,

$$20 \quad \begin{aligned} & \max \{z_1(\mathbf{x}), z_2(\mathbf{x})\} \\ & \text{s.t. } \mathbf{x} \in \mathcal{X} \end{aligned} \tag{1}$$

21 where \mathcal{X} represents the *feasible set in the decision space*, and $z_k(\mathbf{x})$ for each $k = 1, 2$ represents
22 an objective function. We denote all the vectors using bold fonts in this study. The image \mathcal{Y} of
23 \mathcal{X} under vector-valued function $\mathbf{z} = \{z_1, z_2\}$ represents the *feasible set in the objective/criterion*
24 *space*, i.e., $\mathcal{Y} := z(\mathcal{X}) := \{\mathbf{y} \in \mathbb{R}^2 : \mathbf{y} = \mathbf{z}(\mathbf{x}) \text{ for some } \mathbf{x} \in \mathcal{X}\}$.

25 **Definition 1.** A feasible solution $\mathbf{x}' \in \mathcal{X}$ is called *efficient or Pareto optimal*, if there is no other
26 $\mathbf{x} \in \mathcal{X}$ such that $z_k(\mathbf{x}) \geq z_k(\mathbf{x}')$ for $k = 1, 2$ and $\mathbf{z}(\mathbf{x}) \neq \mathbf{z}(\mathbf{x}')$. If \mathbf{x}' is efficient, then $\mathbf{z}(\mathbf{x}')$ is
27 called a *nondominated point*. The set of all efficient solutions $\mathbf{x}' \in \mathcal{X}$ is denoted by \mathcal{X}_E . The set
28 of all nondominated points $\mathbf{z}(\mathbf{x}') \in \mathcal{Y}$ for some $\mathbf{x}' \in \mathcal{X}_E$ is denoted by \mathcal{Y}_N and referred to as the
29 *nondominated frontier*.

30 The focus of this study is on a class of bi-objective optimization problems known as Bi-Objective
31 Knapsack Problem (BOKPs) where

$$32 \quad \mathcal{X} := \{\mathbf{x} \in \{0, 1\}^n : \sum_{i=1}^n w_i x_i \leq W\},$$

33 and

$$34 \quad z_k(\mathbf{x}) := \sum_{i=1}^n c_i^k x_i \quad \forall k \in \{1, 2\}.$$

1 It is assumed that all coefficients/parameters are non-negative integers, i.e., $\mathbf{w} \in \mathbb{Z}_+^n$, $W \in \mathbb{Z}_+$,
2 and $\mathbf{c}^k \in \mathbb{Z}_+^n$ for each $k = 1, 2$. A desirable characteristic of any Binary Knapsack Optimization
3 Problem (BOKP) lies in its finite nondominated frontier, leading to the availability of numerous
4 algorithms for solving any BOKP. In the remaining of this section, we introduce two classical
5 algorithms for solving BOKPs, which we later adapt to develop reinforcement learning-based al-
6 gorithms aimed at solving BOKPs. The first approach is an exact method, while the second is an
7 approximation method.

8 A well-known classical exact method to generate the nondominated frontier of a BOKP is the
9 augmented ε -constraint method [7, 17]. The algorithm solves the following *scalarization* problem,
10 i.e., single-objective integer linear program, at each iteration,

$$\begin{aligned} \max_{\mathbf{x} \in \mathcal{X}} \lambda z_1(\mathbf{x}) + z_2(\mathbf{x}) &= \sum_{i=1}^n (\lambda c_i^1 + c_i^2) x_i \\ \text{s.t. } \sum_{i=1}^n c_i^2 x_i &\geq C_2^* + \varepsilon, \end{aligned} \tag{2}$$

12 where C_2^* represents a parameter that acts as a lower bound on the second objective function's
13 value, while ε takes on a sufficiently small value, and λ adopts a sufficiently large value. It is
14 important to note that the value of C_2^* undergoes changes after each iteration, whereas the other
15 two parameters are fixed from the beginning. Specifically, considering our assumption that all
16 parameters are integers, it is safe to assign $\varepsilon = 1$ and $\lambda = 1 + \sum_{i=1}^n c_i^2$. The latter choice implies
17 that the second objective function should only be optimized if there is no further way to enhance
18 the first objective function's value, even by a single unit. Essentially, the primary objective of the
19 method during each iteration is to optimize the first objective function.

Algorithm 1: The augmented ε -constraint method

```

1 Inputs: An instance of BOKP
2  $C_2^* \leftarrow -1$ 
3  $\mathcal{Y}_N \leftarrow \emptyset$ 
4 SearchDone  $\leftarrow$  False
5 while SearchDone is False do
6    $\mathbf{x}^* \leftarrow \text{Eps}(C_2^*)$ 
7   if  $\mathbf{x}^* = \text{Null}$  then
8     SearchDone  $\leftarrow$  True
9   else
10     $\mathbf{z}^* \leftarrow (\sum_{i=1}^n c_i^1 x_i^*, \sum_{i=1}^n c_i^2 x_i^*)$ 
11    Add  $\mathbf{z}^*$  to  $\mathcal{Y}_N$ 
12     $C_2^* \leftarrow z_2^*$ 
13 return  $\mathcal{Y}_N$ 

```

20 Algorithm 1 details the procedure of the augmented ε -constraint method. In this algorithm, we
21 use the operation $\text{Eps}(C_2^*)$ to show that Problem (2) is solved for a given C_2^* . The output of the
22 operation $\text{Eps}(C_2^*)$ is an optimal solution \mathbf{x}^* for Problem (2) in each iteration. If Problem (2) is
23 *infeasible*, the operation $\text{Eps}(C_2^*)$ sets \mathbf{x}^* to *null*. Importantly, during any iteration, if $\mathbf{x}^* \neq \text{null}$,
24 then \mathbf{x}^* is known to be an efficient solution. Thus, in such cases, $\mathbf{z}^* := (\sum_{i=1}^n c_i^1 x_i^*, \sum_{i=1}^n c_i^2 x_i^*)$ is a

1 nondominated point, which is subsequently appended to the list of previously identified nondomi-
 2 nated points. Initially, C_2^* can be set to -1 (or $-\infty$). Following the discovery of a nondominated
 3 point \mathbf{z}^* within an iteration, the ensuing task involves locating the subsequent closest nondomi-
 4 nated point to it in the next iteration. This succeeding point’s second objective function value is
 5 strictly greater than z_2^* . Consequently, at the end of each iteration, C_2^* can be updated to z_2^* . The
 6 algorithm terminates as soon as $\text{Eps}(C_2^*)$ returns *null*, denoting the successful identification of all
 7 nondominated points.

Algorithm 2: The weighted sum method

```

1 Inputs: An instance of BOKP
2  $\hat{\mathcal{Y}}_N \leftarrow \emptyset$ 
3 SearchDone  $\leftarrow$  False
4 while SearchDone is False do
5   if there is no other value for  $\theta$  that needs to be explored then
6     SearchDone  $\leftarrow$  True
7   else
8     Choose a new value for  $\theta$ 
9      $\mathbf{x}^* \leftarrow \text{WSM}(\theta)$ 
10     $\mathbf{z}^* \leftarrow (\sum_{i=1}^n c_i^1 x_i^*, \sum_{i=1}^n c_i^2 x_i^*)$ 
11    Add  $\mathbf{z}^*$  to  $\hat{\mathcal{Y}}_N$ 
12 return  $\hat{\mathcal{Y}}_N$ 

```

8 The second classical method is known as *the weighted sum method* in the literature of multi-
 9 objective optimization. This method is known to be only an approximate method for non-convex
 10 bi-objective optimization problems such as BOKPs. Specifically, this method does not guarantee
 11 the generation of all nondominated points of BOKPs; instead, it produces a subset of nondominated
 12 points denoted by $\hat{\mathcal{Y}}_N$. The weighted sum method is outlined in Algorithm 2. At each iteration, it
 13 solves the following scalarization problem:

$$14 \quad \max_{\mathbf{x} \in \mathcal{X}} \sum_{i=1}^n \theta c_i^1 x_i + (1 - \theta) c_i^2 x_i = \sum_{i=1}^n (\theta c_i^1 + (1 - \theta) c_i^2) x_i \quad (3)$$

15 where $\theta \in (0, 1)$. We utilize the operation $\text{WSM}(\theta)$ to show that Problem (3) is solved in Algo-
 16 rithm 2 for a given θ . It is worth mentioning that selecting the value of θ can be systematically
 17 accomplished, requiring only a finite number of choices [2, 19]. However, it is not trivial how such
 18 systematic schemes can be utilized in the context of RL. So, instead, we adopt a straightforward
 19 yet effective scheme in this study for developing RL-based solution methods, see Section 4.

20 3. An Exact RL-based Method

21 In this section, we propose an exact RL-based method for solving BOKPs. To achieve this, we
 22 combine and utilize ideas from two distinct fields: Q -learning from the RL literature and a modified
 23 version of the ε -constraint method from the multi-objective optimization literature that we call it
 24 *pseudo augmented ε -constraint method*. It is important to note that we deliberately abstain from
 25 employing the weighted sum method in this section. The rationale is that developing an exact

1 RL-based method by building upon an approximate technique such as the weighted sum method
 2 poses non-trivial challenges. Throughout this section, our proposed method is referred to as *pseudo*
 3 *augmented ε -constraint Q-learning*. The proposed method is specifically designed to “learn” to
 4 generate an optimal solution of the following problem,

$$\begin{aligned} \max_{\mathbf{x} \in \mathcal{X}} \lambda z_1(\mathbf{x}) + z_2(\mathbf{x}) &= \sum_{i=1}^n (\lambda c_i^1 + c_i^2) x_i \\ \text{s.t. } \sum_{i=1}^n c_i^1 x_i &\leq C_1^*. \end{aligned} \tag{4}$$

6 for different values of C_1^* where C_1^* is a parameter that acts as an upper bound for the value of the
 7 first objective function. As an aside, in this study, we sometimes use the operation `PseudoEps(C_1^*)`
 8 to show that Problem (4) is solved for a given C_1^* . The learning process will be conducted using
 9 Q-learning, a method that is guaranteed to yield an optimal solution for Problem (4) with sufficient
 10 time allocation [29]. Therefore, instead of using standard exact optimization techniques to solve
 11 Problem (4), we choose to adopt a methodology centered on learning. It is important to highlight
 12 that Problem (4) stands apart from Problem (2), which constitutes the core of Algorithm 1. This
 13 distinction leads us to call it the pseudo ε -constraint optimization problem. The use of the term
 14 “pseudo” arises from the fact that its optimal solutions are not necessarily efficient (or Pareto op-
 15 timal) for a given C_1^* value. Nonetheless, in accordance with Proposition (1), when considering the
 16 complete spectrum of C_1^* values, the entire nondominated frontier can be derived (after eliminating
 17 dominated points).

18 **Proposition 1.** *For every efficient solution \mathbf{x}' , there exists a corresponding $C_1^* \in \mathbb{Z}$ such that \mathbf{x}'*
 19 *is an optimal solution for Problem (4).*

PROOF. We prove by contradiction. Suppose that there exists an efficient solution \mathbf{x}' but it is not
 an optimal solution for Problem (4) for any value of $C_1^* \in \mathbb{Z}$. We know that $z_1(\mathbf{x}') \in \mathbb{Z}$ because by
 assumptions all parameters of any BOKP are integers. So, we can set $C_1^* = z_1(\mathbf{x}')$, and consider
 any optimal solution of Problem (4), denoted by \mathbf{x}^* . We also note that by construction, because
 $\lambda = 1 + \sum_{i=1}^n c_i^2$, when solving Problem (4), $z_2(\mathbf{x})$ will be optimized only if there is no way to
 improve the value of $z_1(\mathbf{x})$ further even by one unit. With this in mind, since C_1^* is a bound on the
 value of $z_1(\mathbf{x})$ in Problem (4), we know that $z_1(\mathbf{x}^*) \leq C_1^*$. However, it is possible to achieve this
 bound because \mathbf{x}' is a feasible solution with $z_1(\mathbf{x}') = C_1^*$. So, we must have that $z_1(\mathbf{x}^*) = C_1^*$ too,
 i.e., $z_1(\mathbf{x}') = z_1(\mathbf{x}^*)$. Now, we know that \mathbf{x}' is not an optimal solution of Problem (4) for any C_1^* .
 So, we must have that

$$\lambda z_1(\mathbf{x}') + z_2(\mathbf{x}') = \sum_{i=1}^n (\lambda c_i^1 + c_i^2) x'_i < \sum_{i=1}^n (\lambda c_i^1 + c_i^2) x_i^* = \lambda z_1(\mathbf{x}^*) + z_2(\mathbf{x}^*).$$

20 Therefore, since $z_1(\mathbf{x}') = z_1(\mathbf{x}^*)$, we must have that $z_2(\mathbf{x}') < z_2(\mathbf{x}^*)$. So, \mathbf{x}' cannot be an efficient
 21 solution because it is dominated by \mathbf{x}^* which is a contradiction. \square

22 In summary, our proposed method systematically traverses the complete spectrum of C_1^* values,
 23 utilizing Q-learning to generate optimal solutions. After generating optimal solutions for all C_1^*
 24 values, inferior solutions are filtered out, leaving the entire nondominated frontier for reporting. It

1 is noteworthy that our decision to employ Problem (4) within our proposed technique stems from
 2 a key rationale. In the context of BOKPs, Problem (4) is always *feasible*. This is because it is
 3 always possible to create a feasible solution by setting all decision variables equal to zero. This
 4 intrinsic feasibility makes the development an RL approach for Problem (4) more convenient as
 5 such a method naturally focuses on maximizing total rewards rather than directly addressing issues
 6 related to infeasibility.

7 3.1. A Markov Decision Process Formulation

In order to develop an RL approach for learning to generate an optimal solution of Problem (4),
 it is important to explain how Problem (4) can be formulated as a Markov Decision Process (MDP)
 for any arbitrary value of C_1^* . Note that since Problem (4) is deterministic, the MDP formulation
 will also be deterministic. In other words, the proposed MDP is basically a dynamic program. Our
 proposed MDP formulation follows an episodic structure, spanning n steps. During each time step
 t , a decision is made regarding the value of the decision variable x_t in Problem (4). Consequently,
 only two potential actions are considered: setting x_t to either zero or one. At each time step t , the
 agent receives a reward denoted by R , contingent upon the action taken. Specifically, when the
 action is $x_t = 0$, the reward is $R = 0$, and when the action is $x_t = 1$, the reward is $R = \lambda c_t^1 + c_t^2$.
 During time step t , the state of the system is represented as $S = (t, \mathcal{W}_t, \mathcal{C}_t)$, wherein \mathcal{W}_t signifies
 the remaining portion of W at time t , computed as $\mathcal{W}_t = W - \sum_{i=1}^{t-1} w_i x_i$. Similarly, \mathcal{C}_t indicates
 the remaining part of C_1^* at time step t , calculated as $\mathcal{C}_t = C_1^* - \sum_{i=1}^{t-1} c_i^1 x_i$. Note that the collective
 set of all states, denoted by \mathbb{S} , can be defined as

$$\mathbb{S} := \{(t, \mathcal{W}, \mathcal{C}) \in \mathbb{Z} \times \mathbb{Z} \times \mathbb{Z} : 1 \leq t \leq n \text{ and } \max(0, W - \sum_{i=1}^{t-1} w_i) \leq \mathcal{W} \leq W \text{ and } 0 \leq \mathcal{C} \leq \sum_{i=1}^n c_i^1\}.$$

8 It is also important to consider that in certain states, taking action 1 might not be feasible. This
 9 restriction applies when either $w_t > \mathcal{W}_t$ or $c_t^1 > \mathcal{C}_t$. Under these circumstances, at state S , only
 10 action 0 can be selected. Given this consideration, for every state $S \in \mathbb{S}$, the Bellman optimality
 11 equation can be expressed as follows:

$$12 \quad v^*(S) = \begin{cases} v^*(S^0) & \text{Taking action 1 at state } S \text{ is infeasible.} \\ \max\{R(S) + v^*(S^1), v^*(S^0)\} & \text{Otherwise.} \end{cases} \quad (5)$$

13 In this context, $v^*(S)$ represents the optimal return (or total rewards) when initiating from
 14 state S , and $R(S)$ corresponds to the reward attained at state S upon selecting action 1. It is
 15 important to recognize that, by design, a maximum of two other states can be observed from state
 16 S : namely, S^0 and S^1 . Here, S^0 signifies the state observed following the choice of action 0 at state
 17 S , while S^1 denotes the state observed after selecting action 1 at state S . In cases where S serves
 18 as the terminal state (i.e., $t = n$), no further states can be derived thereafter. In such instances, we
 19 simply establish $v^*(S^0) = v^*(S^1) = 0$ to denote the absence of new states following the terminal
 20 state S .

21 3.2. Training: The Proposed RL-based Method

22 Within the field of RL, several methodologies are available to solve the proposed MDP. A
 23 portion of these methods involves running a sufficient number of simulation episodes to learn how
 24 to solve an MDP formulation, i.e., solving its Bellman optimality equation. In this study, we adopt

1 a specific approach known in RL literature as Q -learning [29]. Therefore, it is crucial to clarify the
 2 fundamental concept behind Q -learning. In our proposed MDP formulation, an alternate method to
 3 express $v^*(S^0)$ and $R(S) + v^*(S^1)$ is by employing the notation $Q^*(S, 0)$ and $Q^*(S, 1)$, respectively.
 4 In other words, we define

$$5 \quad Q^*(S, 0) = v^*(S^0), \quad (6)$$

$$6 \quad Q^*(S, 1) = R(S) + v^*(S^1). \quad (7)$$

7 In simple terms, $Q^*(S, A)$ represents the optimal return (or total rewards) that can be achieved
 8 when initiating from state S and taking action $A \in \{0, 1\}$ in that state. Observe that one can use
 9 the provided Bellman optimality Equation (5) to replace $v^*(S^0)$ and $v^*(S^1)$ with their equivalent Q -
 10 functions in Equation (6)-(7). This replacement leads to the formulation of the Bellman optimality
 11 equations based on Q -functions, which can be expressed as follows:

$$12 \quad Q^*(S, 0) = \begin{cases} Q^*(S^0, 0) & \text{Taking action 1 at state } S^0 \text{ is infeasible,} \\ \max\{Q^*(S^0, 1), Q^*(S^0, 0)\} & \text{Otherwise,} \end{cases} \quad (8)$$

$$13 \quad Q^*(S, 1) = \begin{cases} R(S) + Q^*(S^1, 0) & \text{Taking action 1 at state } S^1 \text{ is infeasible,} \\ R(S) + \max\{Q^*(S^1, 1), Q^*(S^1, 0)\} & \text{Otherwise.} \end{cases} \quad (9)$$

14 The underlying idea behind Q -learning is to learn the values of $Q^*(S, 1)$ and $Q^*(S, 0)$ for each
 15 state $S \in \mathbb{S}$ based on Equations (8)-(9) through running a series of simulation episodes. The
 16 learning process begins with arbitrary initializations of Q -values and progresses by refining them
 17 across simulation episodes until they converge. In this context, convergence refers to the point where
 18 the left-hand-side values in Equations (8)-(9) approach (ideally, become equal to) the corresponding
 19 right-hand-side values for each state $S \in \mathbb{S}$. Once the Q -values are learned, determining an optimal
 20 solution for Problem (4) becomes straightforward. Specifically, when observing state $S \in \mathbb{S}$, if
 21 $Q^*(S, 1) \geq Q^*(S, 0)$, then action 1 is optimal; otherwise, action 0 is the optimal choice. With
 22 this foundation established, we next introduce our proposed method, i.e., the pseudo augmented
 23 ε -constraint Q -learning, for solving any BOKP.

24 Our proposed method is illustrated in Algorithm 3. For a given BOKP instance, our approach
 25 initiates by learning the optimal Q -values for the maximum value of C_1^* , which corresponds to
 26 $\sum_{i=1}^n c_i^1$ (Line 3). Once this is learned, the algorithm decreases the value of C_1^* by 1 (as indicated
 27 in Line 21), and the process is repeated to learn optimal Q -values for the new C_1^* value. The process
 28 will be repeated as long as the value of C_1^* remains non-negative. To learn optimal Q -values for a
 29 specific C_1^* value, the algorithm commences with initializing the Q -values (Line 2). We note that
 30 in terms of implementation, it is more efficient to conduct the initialization step on the fly, i.e.,
 31 we initialize a state if the state needs to be created. This strategy is employed because not all
 32 states within \mathbb{S} may be feasible for a BOKP instance, hence there is no requirement to create all
 33 states upfront. We also note that, in theory, we can initialize the Q -values randomly, so in our
 34 implementation, we set them to random small values within the range $[0.2, 0.3]$.

35 To learn the optimal Q -values for a specific C_1^* value, we run E number of episodes (Line 5).
 36 Note that E is a user-defined parameter, and it should be sufficiently large to ensure the convergence
 37 of the Q -values. Each episode encompasses exactly n time steps, mirroring the number of variables
 38 in the BOKP instance. At each time step t , the available actions narrow down to two choices:

1 assigning x_t to zero or one. At the beginning of each episode (Line 6), the initial state is defined as
2 $(1, \mathcal{W}, \mathcal{C})$, with \mathcal{W} set to W and \mathcal{C} to C_1^* . During time step t , an action denoted by A is selected to
3 represent the value of the decision variable x_t . Naturally, if $w_t > \mathcal{W}$ or $c_t^1 > \mathcal{C}$, action A must be
4 set to zero. Otherwise, its value can be randomly determined using the current Q -values following
5 the well-established epsilon-greedy strategy (Line 8). This strategy implies that with a probability
6 of $1 - \epsilon$, the action with the highest Q -value is chosen, while with a probability of ϵ , a random
7 action is selected. It is important to note that ϵ is a small positive value, distinct from the one
8 introduced for the ϵ -constraint method in Section 2. In our implementation, we dynamically adjust
9 the value of ϵ to strike a better balance between exploration and exploitation, where, specifically,
10 we set $\epsilon = 1 - 0.99e^{\frac{\text{episode} - E}{E}}$.

Algorithm 3: The pseudo augmented ϵ -constraint Q -learning

```

1 Inputs: An instance of BOKP, a step-size for learning denoted by  $\alpha$ , and the number of episodes
   denoted by  $E$ 
2 Initialize  $Q^*(S, A)$  randomly for all  $S \in \mathbb{S}$  and  $A \in \{0, 1\}$ 
3  $C_1^* \leftarrow \sum_{i=1}^n c_i^1$ 
4 while  $C_1^* \geq 0$  do
5     foreach  $\text{episode} \in \{1, \dots, E\}$  do
6          $\mathcal{W} \leftarrow W$  and  $\mathcal{C} \leftarrow C_1^*$  and  $S \leftarrow (1, \mathcal{W}, \mathcal{C})$ 
7         foreach  $\text{step } t \in \{1, \dots, n\}$  do
8             Choose an action at state  $S$  and denote it by  $A$ . Specifically, if  $w_t > \mathcal{W}$  or  $c_t^1 > \mathcal{C}$  then
               choose  $A = 0$ . Otherwise, use the epsilon-greedy approach
9              $R \leftarrow (\lambda c_t^1 + c_t^2)A$  where  $R$  is the reward observed
10            if  $t = n$  then
11                 $G \leftarrow 0$ 
12                 $Q^*(S, A) \leftarrow Q^*(S, A) + \alpha[(R + G) - Q^*(S, A)]$ 
13            else
14                 $\mathcal{W} \leftarrow \mathcal{W} - w_t A$  and  $\mathcal{C} \leftarrow \mathcal{C} - c_t^1 A$  and  $S' \leftarrow (t + 1, \mathcal{W}, \mathcal{C})$ 
15                if  $w_{t+1} > \mathcal{W}$  or  $c_{t+1}^1 > \mathcal{C}$  then
16                     $G \leftarrow Q^*(S', 0)$ 
17                else
18                     $G \leftarrow \max \{Q^*(S', 0), Q^*(S', 1)\}$ 
19                 $Q^*(S, A) \leftarrow Q^*(S, A) + \alpha[(R + G) - Q^*(S, A)]$ 
20                 $S \leftarrow S'$ 
21             $C_1^* \leftarrow C_1^* - 1$ 

```

11 It is evident that by selecting action A in state S , two outcomes follow: a reward (denoted by
12 R) is obtained, and a new state (denoted by S') is reached. The observed reward is a direct result
13 of the contribution of the value of x_t to the objective function of the scalarization Problem (4). In
14 simple terms, R can be calculated as $(\lambda c_t^1 + c_t^2)A$, Line 9. It is important to note that the new
15 state S' can be observed only if $t < n$, meaning we are not yet at the terminal state. When $t < n$,
16 the new state is straightforwardly represented as $(t + 1, \mathcal{W} - w_t A, \mathcal{C} - c_t^1 A)$, as indicated in Line 14.
17 Subsequently, the Q -values are updated through an exponential moving average procedure (Lines
18 12 and 19):

19
$$\text{NewEstimate} = \text{OldEstimate} + \alpha(\text{Target} - \text{OldEstimate})$$

1 where $\alpha \in [0, 1]$ is user-defined parameter representing the step-size for learning (default: $\alpha := \epsilon$).
 2 ‘Target’ can be defined as $R + G$ where R is the reward obtained after taking action A at state S ,
 3 and G is an estimate of the maximum total rewards that can be achieved initiating from the new
 4 state, i.e., S' . We note that one can view the ‘OldEstimate’ as the estimation of the left-hand-side
 5 value and ‘Target’, i.e., $R + G$, as the estimation of the right-hand-side value in Equations (8)-(9).
 6 So, we observe that $G = 0$ if $t = n$, i.e., we are at the terminal state (Line 11). Otherwise, G
 7 is equal to $\max \{Q^*(S', 0), Q^*(S', 1)\}$, see Line 18. However, we know that taking action 1 is not
 8 feasible at the state S' if $w_{t+1} > \mathcal{W} - w_t A$ or $c_{t+1}^1 > \mathcal{C} - c_t^1 A$. So, in that case, G can be set to
 9 $Q^*(S', 0)$, see Line 16, instead of $\max \{Q^*(S', 0), Q^*(S', 1)\}$. Finally, after updating the Q -values,
 10 the new state S can be set to S' and the next time step can start, i.e., Line 20.

11 3.3. Testing: Greedy Decoding

12 Upon completing the learning process for a given BOKP instance, we can employ a greedy
 13 decoding approach to generate its entire nondominated frontier. To begin, an initially empty list
 14 of nondominated points, denoted by \mathcal{Y}_N , is created. This list will ultimately encompass all non-
 15 dominated points upon the termination of the greedy decoding approach. To populate this list, we
 16 run $1 + \sum_{i=1}^n c_i^1$ number of episodes to generate $1 + \sum_{i=1}^n c_i^1$ number of solutions. Each episode can
 17 yield a potentially efficient (or Pareto-optimal) solution, which translates to a prospective nondom-
 18 inated point in the criterion space. To start episode $e \in \{0, \dots, \sum_{i=1}^n c_i^1\}$, the initial state is set to
 19 $S = (t, \mathcal{W}, \mathcal{C})$, with $t = 1$, $\mathcal{W} = W$, and $\mathcal{C} = \sum_{i=1}^n c_i^1 - e$. Each episode encompasses precisely n
 20 steps. During time step $t \in \{1, \dots, n\}$, the value of x_t needs to be determined. If $w_t > \mathcal{W}$ or $c_t^1 > \mathcal{C}$,
 21 then x_t is set to 0. Otherwise, if $Q^*(S, 1) > Q^*(S, 0)$, x_t is set to 1; otherwise, it is set to 0. Subse-
 22 quently, the state is updated by setting its parameters to $t + 1$, $\mathcal{W} - w_t x_t$, and $\mathcal{C} - c_t^1 x_t$, respectively.
 23 The next time step is then initiated. Once episode $e \in \{0, \dots, \sum_{i=1}^n c_i^1\}$ terminates, a potential
 24 efficient solution, denoted by $\tilde{\mathbf{x}}$, is available. Consequently, its objective values, $(z_1(\tilde{\mathbf{x}}), z_2(\tilde{\mathbf{x}}))$,
 25 must be computed. Given that $(z_1(\tilde{\mathbf{x}}), z_2(\tilde{\mathbf{x}}))$ might represent a nondominated point, it is added
 26 to \mathcal{Y}_N if no existing point in \mathcal{Y}_N dominates it. Moreover, when adding $(z_1(\tilde{\mathbf{x}}), z_2(\tilde{\mathbf{x}}))$ to \mathcal{Y}_N , if any
 27 point in \mathcal{Y}_N is dominated by $(z_1(\tilde{\mathbf{x}}), z_2(\tilde{\mathbf{x}}))$, that dominated point should be removed from \mathcal{Y}_N .

28 3.4. Numerical Results

29 In general, a main drawback of exact RL-based methods lies in their computational inefficiency.
 30 However, in the next section, we explain how this challenge can be effectively addressed. Meanwhile,
 31 for the purpose of illustrating the capability of the pseudo augmented ϵ -constraint Q -learning to
 32 derive the complete nondominated frontier of BOKPs, we perform a computational study involving
 33 40 small instances, where the size of the state space \mathbb{S} remains manageable. The set of instances
 34 is divided into four classes based on their number of decision variables. Each class is labeled by
 35 Cn where $n \in \{5, 10, 15, 20\}$ is the number of decision variables. To generate each instance, we
 36 randomly select c_i^1 , c_i^2 , and w_i from the discrete uniform distribution in the interval $[1, 20]$. Also,
 37 set $W = \lfloor \frac{\sum_{i=1}^n w_i}{2} \rfloor$. Table 1 summarizes the performance of the proposed exact RL-based method.
 38 Specifically, the table provides details about the minimum, average, and maximum numbers of
 39 points along with their corresponding runtimes (in minutes) for each class of instances. These results
 40 illustrate that the proposed algorithm successfully generates the entire nondominated frontier for all
 41 instances. It is worth noting that in certain cases, a significant number of episodes were necessary
 42 for the convergence of the Q -values. The runtimes serve as an indicator of the episode count required
 43 for Q -value convergence. As expected, the runtime correlates directly with both the instance size

1 and the quantity of nondominated points. Notably, as instance size increases, the expansion of \mathbb{S}
 2 has a direct impact on the runtime.

Table 1: Performance of the proposed exact RL-based method

Class	Info	Min	Average	Max
C5	#Points	1	2	3
	Runtime (minutes)	5	7	10
C10	#Points	1	4	7
	Runtime (minutes)	36	72	125
C15	#Points	3	6	8
	Runtime (minutes)	767	1,989	3,409
C20	#Points	3	9	14
	Runtime (minutes)	6,543	14,485	22,400

3 4. Approximate Deep RL-based Methods

4 The proposed exact RL-based method in Section 3 has two main drawbacks: (1) It requires
 5 individual training for each instance, and (2) it is computationally inefficient due to the so-called
 6 “curse of dimensionality”, where the number of states grows exponentially. As a result, it is not
 7 directly applicable to solving hard BOKP instances. To address these challenges, we propose two
 8 deep Q -learning-based techniques in this section. In both of these approaches, a deep learning
 9 model will be trained through reinforcement learning that its job is to estimate the optimal Q -
 10 values for various states across different instances. The main advantage of the proposed deep
 11 RL-based methods is that they do not focus on generating the optimal Q -values exactly rather
 12 they aim to visit a small fraction of states in order to train a deep learning model for predicting the
 13 optimal Q -values. Our first proposed approach is a customized version of the method proposed in
 14 Section 3. However, in the second method, we replace the pseudo augmented ε -constraint method
 15 with the weighted sum method. Building upon the empirical insights derived from our analysis of
 16 these two methodologies, as detailed in Section 5, we offer insights into prospective pathways for
 17 future research in the domain of RL-driven algorithmic advancements tailored for multi-objective
 18 optimization in Section 6.

19 4.1. Data Normalization

20 The normalization of data plays an important role in shaping the effectiveness of our proposed
 21 deep Q -learning-based methods. Following a series of trials and experiments, we advocate adopting
 22 a mechanism that numerically worked well during the course of our research. In this mechanism,
 23 each BOKP instance is assumed to be generated/sorted in such a way that $\frac{w_1}{c_1} \geq \frac{w_2}{c_2} \geq \dots \geq \frac{w_n}{c_n}$.
 24 Note that in the literature of single-objective knapsack problem, the set of so-called *utility-weight*
 25 *ratios*, i.e., $\{\frac{c_1^1}{w_1}, \dots, \frac{c_n^1}{w_n}\}$, denoted by ρ , plays an important role in developing different variations
 26 of the well-known Dantzig greedy algorithm [6]. Such greedy algorithms can find solutions that are
 27 theoretically guaranteed to be as good as $\frac{1}{2}$ -approximation. So, we are basically trying to exploit
 28 this property in our proposed mechanism.

29 Instances need to be normalized further once being transformed into their suitable scalariza-
 30 tion form. Our proposed Q -learning-based methods are designed to solve different scalarization

1 problems. Specifically, one is designed to generate an optimal solution of the weighted sum prob-
 2 lem, i.e., (3), for any value of θ and the other is designed to generate an optimal solution of the
 3 pseudo augmented ε -constraint problem, i.e., (4), for any value of C_1^* . Regardless of the specific
 4 scalarization problem at hand, it is crucial to normalize the corresponding parameters. For the nor-
 5 malization of objective function parameters, division by their maximum value suffices. Likewise,
 6 the normalization of each constraint involves dividing all constraint parameters by the maximum
 7 value found on the left-hand side of the constraint equation. For the remainder of this section, it
 8 is assumed that this normalization procedure has been executed.

9 *4.2. The Proposed Deep Learning Model*

10 Before providing the details of our proposed deep Q -learning-based methods, we first explain
 11 their underlying deep learning model and some details on its training process. The complexity of
 12 a deep learning model architecture depends on its specific applications. The architecture deemed
 13 appropriate for addressing BOKPs, following a sequence of trial and experimentation, is illustrated
 14 in Figure 1. In the proposed architecture, we pass the state features denoted by \mathcal{S} through two
 15 hidden-layers of fully connected layers (i.e., Feed Forward) to generate the outputs. The outputs
 16 are basically the estimation of optimal Q -values for state \mathcal{S} , i.e., $Q^*(\mathcal{S}, 0)$ and $Q^*(\mathcal{S}, 1)$. We also use
 17 ReLU activation function for the first hidden layer and the Leaky ReLU for the second hidden layer.
 18 We note that we are intentionally using the notation \mathcal{S} rather than S because unlike Section 3, we are
 19 predicting the optimal Q -values through a deep learning model in this section. So, more features are
 20 needed to describe a state in order to ensure that the deep learning model can generate reasonable
 21 estimations for the optimal Q -values even for the states that it has not observed during the training
 22 process. Next, we explain how \mathcal{S} is defined for each of our proposed deep Q -learning-based method.

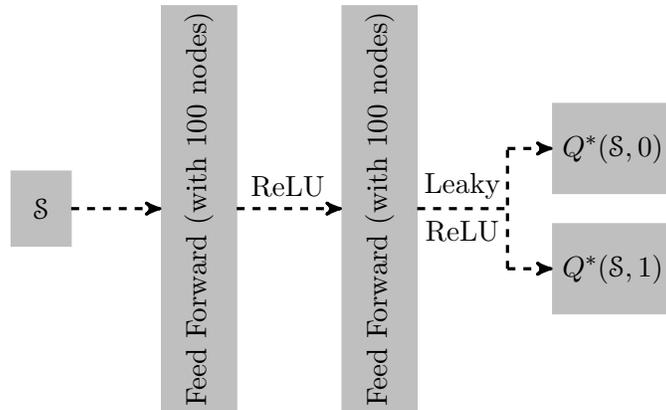


Figure 1: The proposed architecture for the deep learning model

23 Note that by construction at time t , our proposed methods determine the value of x_t . Keeping
 24 this in consideration, for the method tailored to solve the weighted sum scalarization problem, i.e.,
 25 (3), we define \mathcal{S} in the following manner at time t :

$$\mathcal{S} = (n^{[t]}, \mathcal{W}^{[t]}, c_t, w_t, \bar{\rho}^{[t]}, \bar{c}^{[t]}, \bar{w}^{[t]})$$

where $n^{[t]}$ is the number of variables that their values have not yet been determined, $\mathcal{W}^{[t]}$ is the remaining knapsack budget, c_t is the coefficient of x_t in the objective function of the scalarization

problem under consideration (which is Problem (3) here), w_t is the coefficient of x_t in the knapsack constraint. Finally, a notation with “bar” has a specific meaning in this study which is designed to ensure that the cardinality of \mathcal{S} is constant at all times. With this in mind, $\boldsymbol{\rho}^{[t]}$, $\mathbf{c}^{[t]}$, and $\mathbf{w}^{[t]}$ are the vectors of utility-weight ratios, the parameters of the objective function of the scalarization problem under consideration, and the knapsack weights for the variables that their values have not yet been determined at time t , respectively. Note that the cardinality of each of these vectors depends on time t and that makes the cardinality of \mathcal{S} non-constant. So, instead of using them directly, we use their corresponding vectors of their *percentile* values. This implies that $\bar{\boldsymbol{\rho}}^{[t]}$ refers to

$$\text{Percentile}(\boldsymbol{\rho}^{[t]}),$$

1 where the function always returns a vector with 100 components. Component $j \in \{1, \dots, 100\}$ is
 2 the j -th percentile among values in $\boldsymbol{\rho}^{[t]}$. Vectors $\bar{\mathbf{c}}^{[t]}$, and $\bar{\mathbf{w}}^{[t]}$ can be defined similarly, and each
 3 one has 100 components at any time too.

4 The state \mathcal{S} can be established in a similar manner for the method designed to solve the pseudo
 5 augmented ε -constraint scalarization problem, i.e., (4). However, it necessitates the inclusion of
 6 three additional components due to the presence of an extra constraint in problem (4). Therefore,
 7 at time t , \mathcal{S} can be described as follows:

$$\mathcal{S} = (n^{[t]}, \mathcal{W}^{[t]}, c_t, w_t, \bar{\boldsymbol{\rho}}^{[t]}, \bar{\mathbf{c}}^{[t]}, \bar{\mathbf{w}}^{[t]}, C_1^{*[t]}, c_t^1, \bar{\mathbf{c}}^{1,[t]})$$

8 where $C_1^{*[t]}$ is the remaining value from C_1^* , c_t^1 is the coefficient of x_t in $z_1(\mathbf{x})$, and $\mathbf{c}^{1,[t]}$ is the vector
 9 of parameters of $z_1(\mathbf{x})$ for the variables that their values have not yet been determined at time t .
 10 Note that $\bar{\mathbf{c}}^{1,[t]}$ is the vector of the percentile values of $\mathbf{c}^{1,[t]}$. It is also worth mentioning c_t and $\bar{\mathbf{c}}^{[t]}$
 11 are linked to the specific scalarization problem at hand. As a result, they differ from the definitions
 12 provided earlier for the state in the context of the weighted sum scalarization problem.

13 To train our proposed deep learning model, we require a loss function. In this study, we employ
 14 the mean squared error as our selected loss function. As previously explained in Section 3.2,
 15 within the context of Q -learning, the Q -values are updated through an exponential moving average
 16 procedure:

$$\text{NewEstimate} = \text{OldEstimate} + \alpha(\text{Target} - \text{OldEstimate}).$$

17
 18 Observe that “Target - OldEstimate” can be viewed as the *error*. Thus, the **Loss** function in
 19 our study is defined as the mean squared of these errors. Specifically, for each observation, the
 20 error is computed as $(R + G) - Q^*(\mathcal{S}, A)$. Therefore, our objective is to minimize the mean of
 21 $[(R + G) - Q^*(\mathcal{S}, A)]^2$ across batches of observations generated during simulation episodes in our
 22 proposed methods.

23 4.3. Training: The Proposed Deep RL-based Methods

24 In this section, we provide a detailed description of our proposed Deep RL-based methods.
 25 We have named them “deep augmented pseudo ε -constraint Q -learning” and “deep weighted sum
 26 Q -learning”. While these methods share similarities during the training process, they diverge
 27 significantly during testing which we explain them in Section 4.4. The primary distinction between
 28 these methods in the training phase lies in a single key step, along with the way states are defined, as
 29 discussed earlier in Section 4.2. The proposed methods in training mode are outlined in Algorithm 4,
 30 which shares key steps with Algorithm 3. Specifically, the mechanism for computing G remains

Algorithm 4: The Proposed Deep Q -learning-based Methods

```
1 Inputs: A set of BOKP instances, and the number of episodes denoted by  $E$ 
2 Make the size of all instances equal to  $n'$  where  $n'$  is the largest number of decision variables that
  exist in the BOKP instances
3 For each BOKP, generate several random instances for its corresponding scalarization problem, i.e.,
  either (3) or (4), by randomly selecting values for either  $C_1^*$  or  $\theta$ , and denote the index set of all
  such instances by  $\{1, \dots, M\}$ 
4 Generate the DeepLearningModel and initialize its parameters randomly
5 Generate a copy of the DeepLearningModel and denote it by TargetModel
6 Generate an empty experience replay list and denote it by ReplayList
7 foreach  $episode \in \{1, \dots, E\}$  do
8   Partition the set  $\{1, \dots, M\}$  into batches of fixed size, and denote the index set of batches by
      $\{1, \dots, B\}$ 
9   foreach  $batch \in \{1, \dots, B\}$  do
10    Generate the initial state  $S$  of each instance in the batch
11    foreach  $step t \in \{1, \dots, n'\}$  do
12      Call DeepLearningModel( $S$ ) to estimate  $Q^*(S, 0)$  and  $Q^*(S, 1)$  for each instance in the
        batch
13      Choose an action  $A$  (to determine the value  $x_t$ ) following the epsilon-greedy approach
        for each instance in the batch
14      Take the selected action  $A$  and observe the reward  $R$ , and generate the new state  $S'$  for
        each instance in the batch
15      Record the experience observed, i.e.,  $(S, A, Q^*(S, A), R, S')$ , in the ReplayList for each
        instance of the batch
16      Take a random sample set with an arbitrary size from the ReplayList
17      Call TargetModel( $S'$ ) to estimate  $Q^*(S', 0)$  and  $Q^*(S', 1)$  for each experience in the
        sample set
18      Compute  $G$  based on  $\max\{Q^*(S', 0), Q^*(S', 1)\}$  and feasibility conditions for each
        experience in the sample set
19      Compute the squared estimation error  $[(R + G) - Q^*(S, A)]^2$  of each experience in the
        sample set, and denote the mean of squared errors by the LOSS
20      Apply the backpropagation to the DeepLearningModel based on the LOSS
21      Apply one step of the stochastic gradient decent method to the DeepLearningModel
22      Set  $S \leftarrow S'$  for each instance in the batch
23      Once a while set TargetModel  $\leftarrow$  DeepLearningModel
```

1 identical to the one presented in Algorithm 3. Therefore, we do not reiterate this mechanism in
2 this section.

3 Keeping this in view, our proposed methods take a set of BOKP instances and the number of
4 training episodes, denoted by E , as inputs (Line 1). Subsequently, in Line 2, the methods make the
5 size of all BOKP instances the same, i.e, they all will have n' decision variables after transformation
6 where n' is the largest number of decision variables among all provided instances (default: $n' = 200$).
7 It is important to note that even if the initial BOKP instances possess different sizes, they can be
8 transformed into BOKP instances with the fixed size. For instance, if we intend to convert a BOKP
9 instance with 25 variables into an instance with 200 variables, we can achieve this by setting $w_j = 0$
10 and assigning suitably large negative values to c_j^1 and c_j^2 for each $j \in \{26, \dots, 200\}$. This ensures
11 that the RL agent receives a sufficiently large negative reward for setting the value of a decision
12 variable to one if the variable does not exist. As explained in Section 4.1, we assume that the
13 instances are already normalized in this section. Consequently, we adopt the default value of -2
14 as the “sufficiently large negative value”. This choice is substantiated by the fact that the largest
15 objective coefficient of the scalarization problem under consideration, following normalization, is
16 one, making -2 a significant negative reward.

17 Next, in Line 3, the methods start to generate training instances (based on the provided BOKP
18 instances) for their corresponding scalarization problem. This step constitutes the primary dis-
19 tinction between the deep weighted sum Q -learning and the deep augmented pseudo ε -constraint
20 Q -learning during the training process. This divergence arises because they are based on two dis-
21 tinct scalarization problems. The former employs Problem (3), whereas the latter is grounded in
22 Problem (4). Generating training instances can be achieved by generating diverse random values
23 for either θ or C_1^* for each BOKP instance. After generating all training instances, their index set
24 can be denoted by $\{1, \dots, M\}$, and they will be used for training a deep learning model. Therefore,
25 in Line 4, the methods create a random deep learning model, denoted by `DeepLearningModel`,
26 based on the proposed architecture outlined in Figure 1. The subsequent objective is to train this
27 `DeepLearningModel`. However, in the field of deep Q -learning, the state-of-the-art approach in-
28 volves introducing a delayed mechanism (see for example [15]) for updating the parameters of the
29 deep learning model, particularly when estimating Targets (i.e., $R + G$). This mechanism enhances
30 the smooth training and convergence of the deep learning model. With this in mind, the methods
31 establish a copy of `DeepLearningModel`, denoted by `TargetModel`, in Line 5. `TargetModel` will be
32 utilized for estimating Target values, and its parameters will undergo updates with delays in Line
33 23 (default: every n' training time steps).

34 Subsequently, in Line 6, the methods create an empty list referred to as `ReplayList`, commonly
35 known as the list of *experience replay* observations in the deep Q -learning literature [15]. This list
36 possesses a substantial capacity (default: 1,000,000) and continuously maintains only the most
37 recent observations when it is full. When it becomes necessary to update the `DeepLearningModel`,
38 a random sample set of arbitrary size will be drawn from this list for both forward and backward
39 propagations. With `ReplayList` in place, the main training loop starts at Line 7. Within each
40 episode, the training instance set $\{1, \dots, M\}$ is partitioned into batches of fixed size (Line 8, default:
41 50). For each batch, the training process is executed (Line 9). For every instance in the batch,
42 the initial state is extracted and set to \mathcal{S} (Line 10). Subsequently, the methods iterate through
43 n' training time steps, corresponding to the number of decision variables in each training instance
44 (Line 11).

45 At each training time step $t \in \{1, \dots, n'\}$, the methods call `DeepLearningModel(\mathcal{S})` to estimate

1 $Q^*(S, 0)$ and $Q^*(S, 1)$ for each instance in the batch (Line 12). Subsequently, the methods employ
 2 the epsilon-greedy approach to select an action $A \in \{0, 1\}$ to determine the value of x_t for each
 3 instance in the batch. Following the action selection, the methods take the chosen action A to
 4 observe a reward R and generate the next state S' for each instance in the batch (Line 14). It is
 5 worth noting that the way we compute the reward in Algorithm 4 differs from Algorithm 3. This
 6 difference arises from the necessity of ensuring that the deep learning model learns that certain
 7 actions are infeasible. Additionally, it needs to learn that sometimes it is more advantageous to
 8 wait and not set the value of a decision variable to one, even when there are sufficient resources
 9 available. With this in mind, we define the reward obtainable at time t by taking action A at state
 10 S as follows:

$$11 \quad R = \begin{cases} c_t & \text{The proposed method chose } A = 1 \text{ and taking action 1 at time } t \text{ is feasible,} \\ -\beta_1|c_t| & \text{The proposed method chose } A = 1 \text{ but taking action 1 at time } t \text{ is not feasible,} \\ -\beta_2c_t & \text{The proposed method chose } A = 0 \text{ and } c_t \geq 0, \\ 0 & \text{The proposed method chose } A = 0 \text{ and } c_t < 0, \end{cases} \quad (10)$$

12 where β_1 and β_2 are user-defined non-negative penalty values (default $\beta_1 = \beta_2 = 10$).

13 For each instance in the batch, $(S, A, Q^*(S, A), R, S')$ constitutes an experience observed at time
 14 step t . Specifically, it signifies that at time t , we initially occupied state S , proceeded to select action
 15 A , and estimated its corresponding optimal Q -value. We also received an immediate reward of R
 16 and ultimately transitioned to a new state S' . Consequently, the experience $(S, A, Q^*(S, A), R, S')$
 17 is appended to the `ReplayList` in Line 15. Following this, a random sample set of arbitrary
 18 size (default: 50) is drawn from `ReplayList` in Line 16. This sample set serves as input for
 19 computing the loss function value used to train the `DeepLearningModel`. To calculate the loss
 20 function value, it is crucial to determine G for each experience within the sample set. Therefore, in
 21 Line 17, the methods call `TargetModel(S')` to estimate $Q^*(S', 0)$ and $Q^*(S', 1)$ for each experience
 22 in the sample set. Subsequently, in Line 18, G is computed based on $\max\{Q^*(S', 0), Q^*(S', 1)\}$
 23 and the feasibility conditions for each experience in the sample set, mirroring the process in Lines
 24 10-18 of Algorithm 3. Once G is computed, the loss function, which is the mean squared error
 25 denoted by `LOSS`, is calculated. To achieve this, in Line 19, the methods compute the squared error
 26 estimation, i.e., $[(R + G) - Q^*(S, A)]^2$, for each experience in the sample set, and subsequently
 27 compute the average to obtain `LOSS`. Afterwards, in Line 20, backpropagation is applied to the
 28 `DeepLearningModel` based on `LOSS`, and in Line 21, one step of the gradient descent method is
 29 executed to update the parameters of the `DeepLearningModel`. Before starting another training
 30 step, the methods set S to S' for each instance in the batch, as indicated in Line 22.

31 4.4. Testing: Beam Search

32 Once the training process is complete, the proposed methods can be used to generate an ap-
 33 proximate nondominated frontier of any BOKP. The deep weighted sum Q -learning and the deep
 34 augmented pseudo ε -constraint Q -learning employ distinct search mechanisms to approximate the
 35 entire nondominated frontier of a given BOKP. However, when aiming to approximate an optimal
 36 solution for their respective scalarization problems at a specific value of θ or C_1^* , they share a com-
 37 mon procedure. In particular, both methods utilize a technique known as the *beam search* to solve

1 the scalarization problem. The beam search extends the concept of greedy search or decoding by
 2 maintaining a record of the best $L \in \mathbb{Z}_+$ candidates during the process of approximating an optimal
 3 solution for a scalarization problem. Here, L represents a user-defined parameter that determines
 4 the size or length of the beam search. Keeping this in mind, when solving a scalarization problem,
 5 the methods execute a single episode with n time steps.

6 In the context of the greedy decoding approach, when we reach time step $t \in \{1, \dots, n\}$, we
 7 can extract the system’s state, denoted by \mathcal{S} , and then employ the trained deep learning model to
 8 estimate both $Q^*(\mathcal{S}, 0)$ and $Q^*(\mathcal{S}, 1)$. By selecting $A = \arg \max\{Q^*(\mathcal{S}, 0), Q^*(\mathcal{S}, 1)\}$ and identifying
 9 its feasibility (i.e., not violating the constraints of the scalarization problem under consideration),
 10 we can determine the value of x_t at time t . In contrast, the beam search approach maintains a
 11 record of the best L partial solutions at each time step t . To decide which partial solutions should
 12 be maintained at any given time t , a utility function is required for evaluating the value of each
 13 partial solution. A natural choice for this utility function involves summing the rewards obtained
 14 so far by the known part of each partial solution and its corresponding Q -values (which estimate
 15 the maximum total future rewards). Consequently, partial solutions with higher utility values need
 16 to be maintained.

17 Upon completing the beam search to solve a scalarization problem, we have (at most) L complete
 18 solutions in memory. It is evident that the solution with the best objective value for the scalarization
 19 problem under consideration can be viewed as an approximation of its optimal solution. However, it
 20 is important to note that since we have computed L complete solutions, we can also evaluate these
 21 solutions with respect to the original objective functions, namely, $z_1(\mathbf{x})$ and $z_2(\mathbf{x})$. These solutions
 22 might be nondominated in relation to each other. Consequently, we can maintain these solutions
 23 and utilize them to approximate the entire nondominated frontier. So, a desirable property of the
 24 beam search method is its ability to discover several (locally) Pareto-optimal solutions within the
 25 same neighborhood when solving a scalarization problem. This feature is particularly valuable in
 26 the context of deep weighted sum Q -learning. As mentioned in Section 2, the weighted-sum method
 27 theoretically falls short in generating the complete nondominated frontier for non-convex problems
 28 such as BOKPs. Specifically, it fails to find some points which are referred to as *unsupported*
 29 *nondominated points* in the literature of multi-objective optimization [7]. However, thanks to the
 30 beam search, it becomes possible to identify some of these unsupported nondominated points when
 31 solving a scalarization problem.

32 With the beam search defined, we can now provide the specifics of the proposed methods in
 33 the testing mode. The details of the deep augmented pseudo ε -constraint Q -learning in the testing
 34 mode are shown in Algorithm 5. The method starts by accepting a BOKP instance and the beam
 35 search length, denoted by L , as inputs (Line 1). Subsequently, it creates an empty list, denoted by
 36 $\tilde{\mathcal{Y}}_N$, to maintain the set of approximate nondominated points (Line 2). The method’s initial goal
 37 is to approximate the endpoints of the nondominated frontier. To estimate the first endpoint, it
 38 proceeds to solve **PseudoEps**($+\infty$) using the beam search (Line 3). The beam search generates a
 39 maximum of L solutions, which are then added to $\tilde{\mathcal{Y}}_N$ (Line 4). To determine the other endpoint,
 40 the method temporarily switches the objective functions of the BOKP by relabeling them in reverse
 41 order (Line 5). Following this relabeling, it uses the beam search again to solve **PseudoEps**($+\infty$),
 42 see Line 6. This process generates up to L solutions that are then added to $\tilde{\mathcal{Y}}_N$ (Line 7). After
 43 estimating the second endpoint of the nondominated frontier, the method restores the original
 44 labels of the objective functions of the BOKP (Line 8).

45 It is worth noting that estimating both endpoints serves a crucial purpose in determining the

Algorithm 5: The Deep Augmented Pseudo ε -constraint Q -learning – Testing Mode

```
1 Inputs: A BOKP instance, and the beam search length denoted by  $L$ 
2 Create an empty list denoted by  $\tilde{\mathcal{Y}}_N$  to maintain the approximate nondominated points
3 Solve  $\text{PseudoEps}(+\infty)$  using the beam search to compute at most  $L$  solutions
4 Add the solutions found to  $\tilde{\mathcal{Y}}_N$ 
5 Relabel the objective functions of the BOKP in reverse manner, i.e., use  $z_1(\mathbf{x})$  for the second
  objective function and  $z_2(\mathbf{x})$  for the first objective function
6 Solve  $\text{PseudoEps}(+\infty)$  using the beam search to compute at most  $L$  solutions
7 Add the solutions found to  $\tilde{\mathcal{Y}}_N$ 
8 Restore the original labels of the BOKP's objective functions
9  $C_1^* \leftarrow +\infty$ 
10  $\text{SEARCHDONE} \leftarrow \text{False}$ 
11 while  $\text{SEARCHDONE}$  is False do
12    $\text{UPPERBOUND} \leftarrow$  the maximum value of  $z_1(\mathbf{x})$  among all solutions in  $\tilde{\mathcal{Y}}_N$  with  $z_1(\mathbf{x}) \leq C_1^*$ 
13    $C_1^* \leftarrow \text{UPPERBOUND} - 1$ 
14    $\text{LOWERBOUND} \leftarrow$  the minimum value of  $z_1(\mathbf{x})$  among all solutions in  $\tilde{\mathcal{Y}}_N$ 
15   if  $C_1^* < \text{LOWERBOUND}$  then
16      $\text{SEARCHDONE} \leftarrow \text{True}$ 
17   else
18     Solve  $\text{PseudoEps}(C_1^*)$  using the beam search to compute at most  $L$  solutions
19     Add all the solutions found in this iteration to  $\tilde{\mathcal{Y}}_N$ 
20 return  $\text{Pareto}(\tilde{\mathcal{Y}}_N)$ 
```

1 termination condition for the proposed method. This is because the endpoints of the nondominated
2 frontier define the minimum and maximum values of $z_1(\mathbf{x})$ among all nondominated solutions.
3 Consequently, we can constrain the values of C_1^* within this minimum and maximum range. With
4 this in mind, in Lines 9-11, we initialize C_1^* to $+\infty$ and initiate a search loop (Lines 10-11). In each
5 iteration of the loop, the method aims to compute a new UPPERBOUND for C_1^* . This is achieved
6 by searching through all solutions in $\tilde{\mathcal{Y}}_N$ with $z_1(\mathbf{x}) \leq C_1^*$ and identifying the maximum value of
7 $z_1(\mathbf{x})$ among those solutions (Line 12). After computing UPPERBOUND , C_1^* is updated in Line 13.
8 Specifically, it is set to $\text{UPPERBOUND} - 1$ in the hope of discovering solutions distinct from those
9 already present in $\tilde{\mathcal{Y}}_N$. Subsequently, the algorithm computes a LOWERBOUND for C_1^* . This is
10 achieved by searching through all solutions in $\tilde{\mathcal{Y}}_N$ to identify the minimum value of $z_1(\mathbf{x})$ (Line 14).
11 The search terminates when $C_1^* < \text{LOWERBOUND}$ is met (Lines 15-16). If not, the beam search is
12 utilized to solve $\text{PseudoEps}(C_1^*)$ and the solutions found are added to $\tilde{\mathcal{Y}}_N$ (Lines 17-19). Finally, a
13 filtering process is conducted over $\tilde{\mathcal{Y}}_N$ to keep and report only the solutions that are nondominated
14 in relation to each other. This process is denoted by $\text{Pareto}(\tilde{\mathcal{Y}}_N)$, see Line 20.

15 The testing process for the deep weighted sum Q -learning is more straightforward and is outlined
16 in Algorithm 6. The method starts by accepting a BOKP instance and the beam search length,
17 denoted by L , as inputs (Line 1). It then creates an empty list, denoted by $\tilde{\mathcal{Y}}_N$, to manage the
18 set of approximate nondominated points (Line 2). Subsequently, the method initiates a loop that
19 systematically increases the value of θ from 0 to 1 (Line 3). The increment size for θ is a user-
20 defined parameter, which is set to 0.01 in this study. For each value of θ , the method calls the
21 beam search to solve $\text{WSM}(\theta)$, and all the solutions found are added to $\tilde{\mathcal{Y}}_N$ (Lines 4-5). Finally, in

Algorithm 6: The Deep Weighted Sum Q -learning – Testing Mode

```
1 Inputs: A BOKP instance, and the beam search length denoted by  $L$ 
2 Create an empty list denoted by  $\tilde{\mathcal{Y}}_N$  to maintain the approximate nondominated points
3 foreach  $\theta \in \{0, 0.01, 0.02, \dots, 1\}$  do
4   | Solve WSM( $\theta$ ) using the beam search to compute at most  $L$  solutions
5   | Add all the solutions found in this iteration to  $\tilde{\mathcal{Y}}_N$ 
6 return Pareto( $\tilde{\mathcal{Y}}_N$ )
```

1 Line 6, the method filters and reports the results using **Pareto**($\tilde{\mathcal{Y}}_N$).

2 5. Computational Study

3 In this section, we conduct a computational study to show the efficacy of the proposed deep RL-
4 based methods. All computational experiments are executed on a Dell PowerEdge R360 equipped
5 with dual Intel Xeon E5-2650 2.2 GHz 12-Core Processors (30MB), 128GB of RAM, and the RedHat
6 Enterprise Linux 7.0 operating system. We employ Python and the PyTorch package to implement
7 all the proposed methods in this study. Our instances and Python source codes are accessible
8 through GitHub (<https://github.com/Multi-Objective-Optimization-Laboratory/DeepRL-BiobjectiveKP>). For the
9 training of each of our proposed methods, we employ approximately 10,000,000 training steps,
10 i.e., $E \times B \times n' \simeq 10,000,000$, as outlined in Algorithm 4. Attaining this can be accomplished
11 through the assignment of $E = 76$ and $M = 35,200$. It is important to note that in Algorithm 4,
12 M corresponds to the count of scalarization training instances. Each instance encompasses 25 to
13 200 variables in this study. To elaborate, 200 scalarization instances are generated for every value
14 of n within the range $\{25, 26, \dots, 200\}$. As a result, the total value of M can be calculated as
15 $176 \times 200 = 35,200$.

16 As mentioned in the Introduction, this study centers on hard instances of BOKPs. This choice
17 is rooted in the potential for these complex cases to illuminate the advantages of utilizing deep RL-
18 based methods for solving deterministic combinatorial optimization problems compared to existing
19 highly effective exact and heuristic techniques. To address this, we adopt a technique proposed by
20 Pisinger [22] to create a specific category of challenging BOKP instances referred to as “inverse
21 strongly correlated instances”. To elaborate, we randomly generate vectors \mathbf{c}^1 and \mathbf{c}^2 from a discrete
22 uniform distribution within the range $[1, R]$, where $R = 100$ in this study. The coefficients for the
23 knapsack constraint are determined as $w_i = c_i^1 + \frac{R}{10}$ for each $i = 1, \dots, n$. Finally, the value W
24 is drawn from a discrete uniform distribution within the interval $[0.25 \sum_{i=1}^p w_i, 0.5 \sum_{i=1}^p w_i]$. We
25 generate two sets of testing instances, each containing 6 classes. In each class, there are 10 instances
26 of the same size, labeled as C_n , where n represents the number of decision variables. For the first
27 set, we consider values of n from the set $\{25, 30, 35, 40, 45, 50\}$; for the second set, we use values of
28 n from $\{75, 100, 125, 150, 175, 200\}$.

29 The rationale behind having two sets of testing instances stems from the division of our com-
30 putational analysis into two distinct sections. The initial segment revolves around the utilization
31 of the first set of instances. Its primary focus lies in demonstrating the efficacy of the proposed
32 methods using instances where computing the exact nondominated frontier can be achieved within
33 a 24-hour time limit. For generating the exact nondominated frontier of each instance, Algorithm 1
34 is employed in conjunction with Gurobi 10.0.1 to solve the corresponding single-objective opti-

1 mization problems. The subsequent section pertains to the utilization of the second, larger set
 2 of instances. Its objective is to assess the effectiveness of our proposed methods in contrast to a
 3 state-of-the-art heuristic approach known as the Multi-Directional Local Search (MDLS) [30].

4 In light of the above, for a given instance, the true nondominated frontier \mathcal{Y}_N may be either
 5 known or unknown. In the latter scenario, we combine the points generated by our proposed
 6 methods and MDLS to create what we term the “best-known approximate nondominated frontier”,
 7 which we treat as the true \mathcal{Y}_N . To assess the efficacy of any approximation method (e.g., our
 8 proposed deep RL-based approaches or MDLS), the evaluation revolves around their respective
 9 reported approximate nondominated frontiers, denoted by $\tilde{\mathcal{Y}}_N$, using various metrics. Given that
 10 there is not a singular measure to accomplish this task, the common practice is to employ multiple
 11 metrics [24]. In the following sections, we explore several metrics adopted in this study. Interested
 12 readers may refer to [4, 20] for further details.

13 To introduce the metrics, we require some new notation. The Euclidean distance between two
 14 points \mathbf{y} and \mathbf{y}' is represented as $d(\mathbf{y}, \mathbf{y}')$. Let $k(\mathbf{y})$ denote the nearest point in the approximate
 15 frontier to a true nondominated point $\mathbf{y} \in \mathcal{Y}_N$. Lastly, for each $\mathbf{y} \in \tilde{\mathcal{Y}}_N$, $n(\mathbf{y})$ is defined as the
 16 count of true nondominated points $\mathbf{y}' \in \mathcal{Y}_N \setminus \tilde{\mathcal{Y}}_N$ where $k(\mathbf{y}') = \mathbf{y}$.

- **Hypervolume gap.** An established measure for evaluating and comparing approximate
 nondominated frontiers is the *hypervolume indicator* (or *S-metric*), which assesses the volume
 of the dominated region within the criterion space relative to a reference point [31]. In
 this study, to ensure non-negativity of the hypervolume, the reference point, denoted by
 $\mathbf{y}^R := (y_1^R, y_2^R)$, is determined by setting y_i^R to the minimum value of $z_i(\mathbf{x})$ among all solutions
 reported by different methods (used in this study) for each $i \in \{1, 2\}$. With this consideration,
 we define the hypervolume gap as follows:

$$\frac{100 \times (\text{Hypervolume of } \mathcal{Y}_N - \text{Hypervolume of } \tilde{\mathcal{Y}}_N)}{\text{Hypervolume of } \mathcal{Y}_N}.$$

17 As a result, an approximate nondominated frontier with a smaller hypervolume gap is deemed
 18 more favorable.

- **Cardinality.** We define the cardinality of an approximate nondominated frontier straightfor-
 wardly as the proportion of true nondominated points it encompasses, expressed as $\frac{100 \times |\tilde{\mathcal{Y}}_N \cap \mathcal{Y}_N|}{|\mathcal{Y}_N|}$.
 Accordingly, an approximate nondominated frontier with a larger cardinality is considered
 21 better.
- **Coverage.** A straightforward coverage metric, which gauges the average distance between
 23 true nondominated points not included in the approximate frontier and their nearest coun-
 24 terparts within the approximate frontier, can be described as follows:

$$26 \quad f^a := \frac{\sum_{\mathbf{y} \in \mathcal{Y}_N \setminus \tilde{\mathcal{Y}}_N} d(k(\mathbf{y}), \mathbf{y})}{|\mathcal{Y}_N \setminus \tilde{\mathcal{Y}}_N|},$$

27 Note that f^a can be interpreted as a gauge of how dispersed the points are within the approx-
 28 imate nondominated frontier. Smaller values of f^a imply that the nondominated points in the
 29 approximate frontier are spread across different areas of the criterion space. Consequently,
 30 an approximate nondominated frontier with a lower f^a is more preferable.

- **Uniformity.** An indicator of uniformity should reflect how effectively the points within an approximate frontier are distributed. Points clustered closely together do not contribute positively to the quality of an approximate frontier. Thus, we formulate the uniformity indicator as follows:

$$\mu := \frac{\sum_{\mathbf{y} \in \tilde{\mathcal{Y}}_N} n(\mathbf{y})}{|\tilde{\mathcal{Y}}_N|}.$$

Therefore, an approximate nondominated frontier with a lower value of μ is considered more favorable.

It is worth mentioning that while the training phase of the proposed deep RL-based methods demands a considerable amount of time, their efficiency during the testing phase is noteworthy. To be specific, it takes less than 2 seconds for them to solve a BOKP instance with $n = 25$ variables when the beam search size is set to one, i.e., $L = 1$. It is important to recognize that deep RL methods are inherently suited for parallelization, which can substantially enhance their speed during testing. Nevertheless, assuming a lack of parallelization, the runtime of the proposed deep RL-based methods demonstrate an almost linear correlation with the values of L and n . For instance, if we consider a BOKP instance with $n = 200$ and $L = 20$, the projected runtime can be approximated as 320 seconds using the formula $2 \times \frac{n}{25} \times L$. Given these considerations, we abstain from providing specific solution times for the proposed methods within this study.

5.1. Comparisons with an Exact Method

The goal of this section is to assess the potential effectiveness of the proposed deep RL-based methods in solving BOKPs. To achieve this goal, we evaluate the performance of the proposed methods using the first testing set. This set consists of instances where the true nondominated points are already identified. We conduct these evaluations under four different beam search sizes, i.e., L , selected from the set $\{5, 10, 15, 20\}$. It is important to note, as explained in Section 4.4, that the size of the beam search significantly influences the generation of nondominated points, including unsupported ones. Thus, a larger value of L is anticipated to yield a better approximation.

We employ the labels $W5$, $W10$, $W15$, and $W20$ to represent the deep weighted sum Q -learning method under varying beam search sizes. Similarly, the labels $E5$, $E10$, $E15$, and $E20$ denote the deep augmented pseudo ε -constraint Q -learning approach across different beam search sizes. The performance of the proposed methods across diverse metrics is depicted in Figure 2. It is evident that the performance of both methods exhibits enhancement across all metrics as the beam search size increases. Nevertheless, the deep weighted sum Q -learning method consistently outperforms the deep augmented pseudo ε -constraint Q -learning method, regardless of the specific beam search size employed.

Table 2 better highlights the superior effectiveness of the proposed deep weighted sum Q -learning method (for $L = 20$). The values presented in this table are averages computed over 10 instances. Notably, the deep weighted sum Q -learning has identified approximately 3.61 times more approximate nondominated points ($\tilde{\mathcal{Y}}_N$) on average. This method has successfully generated around 50.5% of the true nondominated points (\mathcal{Y}_N) on average, which is about 10 times better than the alternate method, i.e., the deep augmented pseudo ε -constraint Q -learning. As a consequence, the average hypervolume gap of the deep weighted sum Q -learning method is about 2.9%, providing a much more accurate approximation. In contrast, the alternate method has an average hypervolume gap of 62.7%, indicating a significantly poorer performance. When considering metrics such as

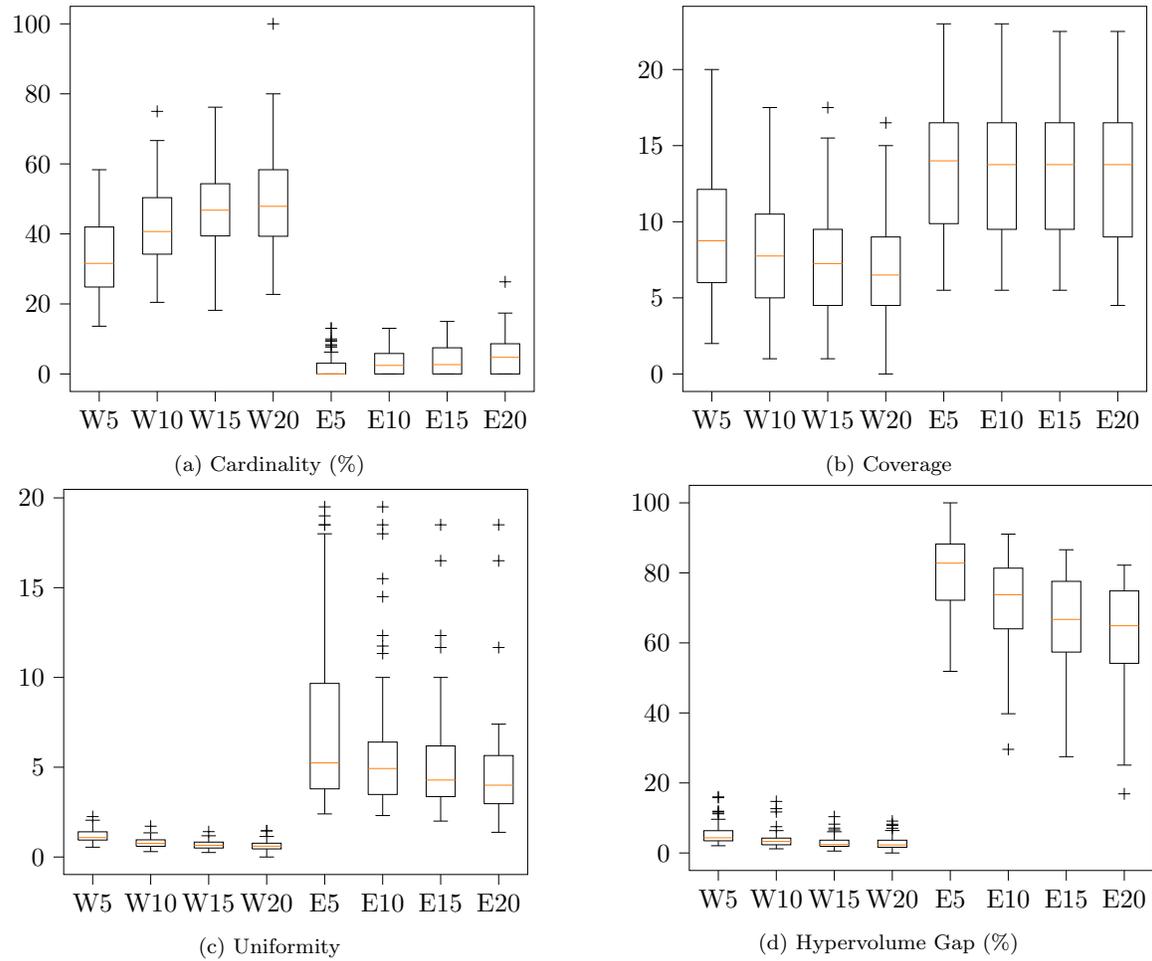


Figure 2: Performance comparisons between the proposed methods for different beam search sizes

Table 2: Performance comparisons between the proposed methods for $L = 20$

Class	\mathcal{Y}_N	Deep Weighted Sum Q -Learning					Deep ε -Constraint Q -Learning				
		$\tilde{\mathcal{Y}}_N$	Car.	Cov.	Uni.	Hyp.	$\tilde{\mathcal{Y}}_N$	Car.	Cov.	Uni.	Hyp.
C25	19.4	15.0	57.9%	3.9	0.6	3.9%	7.5	9.3%	8.3	2.7	43.1%
C30	24.4	19.1	54.8%	5.2	0.6	3.6%	7.3	6.7%	11.0	3.4	54.4%
C35	25.9	20.5	47.7%	6.3	0.7	3.4%	7.0	5.5%	11.7	3.6	63.8%
C40	30.8	24.7	51.7%	7.2	0.6	2.2%	7.4	6.6%	14.0	4.0	65.7%
C45	35.5	30.1	45.4%	9.6	0.7	2.2%	6.6	3.2%	16.7	5.5	71.6%
C50	35.5	31.8	45.7%	9.2	0.6	2.1%	4.7	1.5%	17.0	9.0	77.4%
Avg	28.6	23.5	50.5%	6.9	0.6	2.9%	6.8	5.5%	13.1	4.7	62.7%

1 uniformity and coverage, the deep weighted sum Q -learning method surpasses the alternate method
 2 by approximately 2-fold and 8-fold, respectively.

3 The poor performance of the proposed deep augmented pseudo ε -constraint Q -learning method
 4 can be attributed, in part, to the nature of its corresponding scalarization problem, expressed in
 5 Equation (4). This scalarization problem places a distinct emphasis on the first objective function,
 6 evident through the influence of the constant λ . This skewed focus prompts the deep learning
 7 model to exhibit a strong bias towards the first objective function. Consequently, the significance
 8 of the second objective value diminishes, rendering it akin to an approximation error. As a re-
 9 sult, the deep augmented pseudo ε -constraint Q -learning method has the potential to yield points
 10 with significantly inaccurate second objective values. Consequently, the generated approximate
 11 nondominated frontier through this approach tends to be of low quality.

12 Naturally, one potential solution to address this issue is to train the model with a smaller
 13 value of λ . During the course of our research, we indeed explored this avenue and encountered
 14 two challenges. Firstly, determining an appropriate value for λ in such cases is nontrivial. More
 15 crucially, while the choice of smaller λ values leads the learning process to favor the second objective
 16 function, there is a risk of disregarding numerous nondominated points as this procedure is closely
 17 linked with determining the value of C_1^* . Hence, the effectiveness of the deep weighted sum Q -
 18 learning stems from its ability to adapt to various weights instead of relying on a fixed weight.

19 *5.2. Comparisons with a Heuristic Method*

20 Given that the deep weighted sum Q -learning has proven effective among our proposed meth-
 21 ods, we employ it in this section. The goal of this section is to determine whether the proposed
 22 method, when $L = 20$, can compete with a state-of-the-art technique called MDLS [30] in solving
 23 challenging BOKPs. To accomplish this goal, we assess the performance of the proposed methods
 24 using the second testing set. This set comprises instances in which the true nondominated points
 25 are unknown. Notably, even finding a single nondominated point in this set can be exceedingly dif-
 26 ficult when employing an exact method such as Algorithm 1. As previously mentioned, to address
 27 this challenge, we combine the sets of points identified by both methods and consider the best ones
 28 as the true set of nondominated points, i.e., \mathcal{Y}_N .

Table 3: Performance comparisons with MDLS for $L = 20$

Class	\mathcal{Y}_N	Deep Weighted Sum Q -Learning					MDLS				
		$\tilde{\mathcal{Y}}_N$	Car.	Cov.	Uni.	Hyp.	$\tilde{\mathcal{Y}}_N$	Car.	Cov.	Uni.	Hyp.
C75	45.1	45.0	98.0%	0.1	0.0	0.5%	24.2	3.2%	21.4	1.8	26.9%
C100	55.1	53.7	97.3%	0.3	0.0	0.0%	28.7	3.4%	26.1	1.9	26.6%
C125	60.5	59.7	98.7%	0.1	0.0	0.0%	31.3	1.5%	29.3	1.9	27.1%
C150	71.2	70.4	98.8%	0.1	0.0	0.0%	33.7	1.4%	34.6	2.1	27.2%
C175	80.8	80.0	99.0%	0.2	0.0	0.0%	35.3	1.1%	39.5	2.3	26.9%
C200	57.9	44.7	75.1%	6.8	0.3	5.2%	38.2	25.7%	20.9	1.1	11.9%
Avg	61.8	58.9	94.5%	1.3	0.1	0.9%	31.9	6.1%	28.6	1.9	24.4%

29 Table 3 presents a comprehensive summary of the performance comparisons between the pro-
 30 posed method and MDLS, effectively showcasing the potential of the proposed approach. To im-
 31 prove the quality of solutions produced by MDLS, we adjusted certain parameters, such as the
 32 number of iterations and the quantity of initial locally nondominated points, to values larger than

1 their default settings. Consequently, MDLS takes approximately 800 seconds to solve each of these
 2 instances. The values provided in Table 3 are averages derived from 10 instances. Notably, the deep
 3 weighted sum Q -learning has identified approximately 1.9 times more approximate nondominated
 4 points, i.e., $\tilde{\mathcal{Y}}_N$ on average. This method has managed to generate approximately 94.5% of the
 5 “true” nondominated points, i.e., \mathcal{Y}_N on average, marking a 3.3-fold improvement compared to
 6 MDLS. Consequently, the average hypervolume gap of the deep weighted sum Q -learning hovers
 7 around 0.9%, signifying a notably accurate approximation. In contrast, the alternate method ex-
 8 hibits an average hypervolume gap of 24.4%, illustrating a relatively lesser effectiveness. Regarding
 9 metrics such as uniformity and coverage, the deep weighted sum Q -learning method excels over the
 10 alternate method by approximately 22-fold and 19-fold, respectively.

11 6. Final Remarks

12 In this study, we introduced new exact and approximate solution methods for BOKPs by merg-
 13 ing concepts from the (deep) Q -learning in RL with two classical scalarization methods from multi-
 14 objective integer programming: the ε -constraint method and the weighted sum method. Our
 15 findings demonstrated that combining a scalarization method capable of generating the exact non-
 16 dominated frontier in multi-objective integer programming, e.g., the ε -constraint method, with an
 17 exact reinforcement learning method such as Q -learning allows us to construct an exact RL-based
 18 approach for BOKPs. However, it is essential to note that this exact method comes with compu-
 19 tational inefficiencies. On the other hand, we also explored the possibility of creating an efficient
 20 approximate RL-based approach for BOKPs by uniting scalarization methods with an approximate
 21 RL method, i.e., the deep Q -learning.

22 Our numerical results brought to light three key observations: Firstly, to develop a high-quality
 23 approximate RL-based approach, it is preferable to choose a method from multi-objective optimiza-
 24 tion that dynamically adjusts the objective function weights during the search, such as the weighted
 25 sum method. Secondly, we found that the performance of our proposed methods experiences signif-
 26 icant enhancements when we employ beam search instead of greedy decoding. This strategic shift
 27 enables the identification of numerous (locally) nondominated points within the vicinity of certain
 28 points. Lastly, our research revealed that RL-based methods exhibit superior performance when
 29 dealing with hard instances of BOKPs.

30 These insights have paved the way for several promising avenues for future research. One such
 31 direction involves extending the proposed deep weighted Q -learning approach to handle hard classes
 32 of bi-objective integer programs, particularly those where solving their single-objective optimiza-
 33 tions is computationally intensive, e.g., vehicle routing problem. Achieving this is anticipated to
 34 require more advanced deep learning models at the heart of our proposed method, including but
 35 not limited to pointer networks and transformer models. The primary contribution of this research
 36 direction lies in the development of a suitable deep learning model. Another intriguing research
 37 path involves extending our approach to scenarios with more than two objective functions and
 38 exploring various scalarization techniques from the literature, customizing them to accommodate
 39 dynamic weights for the objective functions. This expansion can possibly enable the development
 40 of more effective RL-based methods. Lastly, it is worth exploring alternative deep RL methods,
 41 such as deep REINFORCE or deep Actor-critic approaches, as potential substitutes for the deep
 42 Q -learning. This exploration can help amplify the capabilities of the deep RL-based methods when
 43 solving hard multi-objective integer programs.

1 References

- 2 [1] Aickelin, U., Khorshidi, H. A., Qu, R., Charkhgard, H., 2023. Guest editorial special issue on
3 multiobjective evolutionary optimization in machine learning. *IEEE Transactions on Evolutionary
4 Computation* 27 (4), 746–748.
- 5 [2] Aneja, Y. P., Nair, K. P. K., 1979. Bicriteria transportation problem. *Management Science* 27,
6 73–78.
- 7 [3] Bengio, Y., Lodi, A., Prouvost, A., 2021. Machine learning for combinatorial optimization: A
8 methodological tour d’horizon. *European Journal of Operational Research* 290 (2), 405–421.
- 9 [4] Boland, N., Charkhgard, H., Savelsbergh, M., 2016. The L-shape search method for triobjective
10 integer programming. *Mathematical Programming Computation* 8 (2), 217–251.
- 11 [5] Dächert, K., Klamroth, K., Lacour, R., Vanderpooten, D., 2017. Efficient computation of
12 the search region in multi-objective optimization. *European Journal of Operational Research*
13 260 (3), 841 – 855.
- 14 [6] Dantzig, G. B., 1957. Discrete-variable extremum problems. *Operations Research* 5 (2), 266–
15 288.
- 16 [7] Ehrgott, M., 2005. *Multicriteria optimization*, 2nd Edition. Springer, New York.
- 17 [8] Gadegaard, S. L., Nielsen, L. R., Ehrgott, M., 2019. Bi-objective branch-and-cut algorithms
18 based on lp relaxation and bound sets. *INFORMS Journal on Computing* 31 (4), 790–804.
- 19 [9] Helfrich, S., Perini, T., Halffmann, P., Boland, N., Ruzika, S., 2023. Analysis of the weighted
20 tchebycheff weight set decomposition for multiobjective discrete optimization problems. *Journal
21 of Global Optimization* 86, 417–440.
- 22 [10] Hu, Y., Yao, Y., Lee, W. S., 2020. A reinforcement learning approach for optimizing multiple
23 traveling salesman problems over graphs. *Knowledge-Based Systems* 204, 106244.
- 24 [11] Kirlik, G., Sayın, S., 2014. A new algorithm for generating all nondominated solutions of mul-
25 tiobjective discrete optimization problems. *European Journal of Operational Research* 232 (3),
26 479 – 488.
- 27 [12] Kool, W., van Hoof, H., Gromicho, J., Welling, M., 2022. Deep policy dynamic program-
28 ming for vehicle routing problems. In: *International conference on integration of constraint
29 programming, artificial intelligence, and operations research*. Springer, pp. 190–213.
- 30 [13] Li, K., Zhang, T., Wang, R., 2021. Deep reinforcement learning for multiobjective optimization.
31 *IEEE Transactions on Cybernetics* 51 (6), 3103–3114.
- 32 [14] Lokman, B., Köksalan, M., 2013. Finding all nondominated points of multi-objective integer
33 programs. *Journal of Global Optimization* 57 (2), 347–365.
- 34 [15] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller,
35 M., 2013. Playing atari with deep reinforcement learning. In: *Proceedings of the 26th Inter-
36 national Conference on Neural Information Processing Systems - Volume 1*. pp. 201–220.

- 1 [16] Nazari, M., Oroojlooy, A., Takáč, M., Snyder, L. V., 2018. Reinforcement learning for solving
2 the vehicle routing problem. In: Proceedings of the 32nd International Conference on Neural
3 Information Processing Systems. NIPS'18. Curran Associates Inc., Red Hook, NY, USA, p.
4 9861–9871.
- 5 [17] Özlen, M., Azizoglu, M., 2009. Multi-objective integer programming: A general approach for
6 generating all non-dominated solutions. *European Journal of Operational Research* 199, 25–35.
- 7 [18] Özlen, M., Burton, B. A., MacRae, C. A. G., 2013. Multi-objective integer pro-
8 gramming: An improved recursive algorithm. *Journal of Optimization Theory and*
9 *Applications*10.1007/s10957-013-0364-y.
- 10 [19] Özpeynirci, Ö., Köksalan, M., 2010. An exact algorithm for finding extreme supported non-
11 dominated points of multiobjective mixed integer programs. *Management Science* 56 (12),
12 2302–2315.
- 13 [20] Pal, A., Charkhgard, H., 2019. Fpbh: A feasibility pump based heuristic for multi-objective
14 mixed integer linear programming. *Computers & Operations Research* 112, 104760.
- 15 [21] Peng, B., Wang, J., Zhang, Z., 2020. A deep reinforcement learning algorithm using dynamic
16 attention model for vehicle routing problems. In: *Artificial Intelligence Algorithms and Ap-*
17 *plications: 11th International Symposium, ISICA 2019, Guangzhou, China, November 16–17,*
18 *2019, Revised Selected Papers 11.* Springer, pp. 636–650.
- 19 [22] Pisinger, D., 2005. Where are the hard knapsack problems? *Computers & Operations Research*
20 32 (9), 2271–2284.
- 21 [23] Reiners, M., Klamroth, K., Heldmann, F., Stiglmayr, M., 2022. Efficient and sparse neural
22 networks by pruning weights in a multiobjective learning approach. *Computers & Operations*
23 *Research* 141, 105676.
- 24 [24] Sayın, S., 2000. Measuring the quality of discrete representations of efficient sets in multiple
25 objective mathematical programming. *Mathematical Programming* 87 (3), 543–560.
- 26 [25] Shao, Y., Lin, J. C.-W., Srivastava, G., Guo, D., Zhang, H., Yi, H., Jolfaei, A., 2021. Multi-
27 objective neural evolutionary algorithm for combinatorial optimization problems. *IEEE trans-*
28 *actions on neural networks and learning systems.*
- 29 [26] Smith-Miles, K., Christiansen, J., noz, M. A. M., 2021. Revisiting where are the hard knapsack
30 problems? via instance space analysis. *Computers & Operations Research* 128, 105184.
- 31 [27] Soyly, B., Yıldız, G. B., 2016. An exact algorithm for biobjective mixed integer linear pro-
32 gramming problems. *Computers & Operations Research* 72, 204–213.
- 33 [28] Stidsen, T., Andersen, K. A., Dammann, B., 2014. A branch and bound algorithm for a class
34 of biobjective mixed integer programs. *Management Science* 60 (4), 1009–1032.
- 35 [29] Sutton, R. S., Barto, A. G., 2018. *Reinforcement Learning: An Introduction*, 2nd Edition. The
36 MIT Press.

- 1 [30] Tricoire, F., 2012. Multi-directional local search. *Computers & Operations Research* 39 (12),
2 3089 – 3101.
- 3 [31] Zitzler, E., Brockhoff, D., Thiele, L., 2007. The hypervolume indicator revisited: On the design
4 of Pareto-compliant indicators via weighted integration. In: Obayashi, S., Deb, K., Poloni, C.,
5 Hiroyasu, T., Murata, T. (Eds.), *Evolutionary Multi-Criterion Optimization*. Vol. 4403 of
6 *Lecture Notes in Computer Science*. pp. 862–876.