

# Deep Reinforcement Learning for Exact Combinatorial Optimization: Learning to Branch

Tianyu Zhang<sup>§</sup>

University of Alberta  
Edmonton, Canada

tianyu.zhang@ualberta.ca

Amin Banitalebi-Dehkordi

Huawei Technologies Canada Co., Ltd.  
Vancouver, Canada

amin.banitalebi@huawei.com

Yong Zhang

Huawei Technologies Canada Co., Ltd.  
Vancouver, Canada

yong.zhang3@huawei.com

**Abstract**—Branch-and-bound is a systematic enumerative method for combinatorial optimization, where the performance highly relies on the variable selection strategy. State-of-the-art handcrafted heuristic strategies suffer from relatively slow inference time for each selection, while the current machine learning methods require a significant amount of labeled data. We propose a new approach for solving the data labeling and inference latency issues in combinatorial optimization based on the use of the reinforcement learning (RL) paradigm. We use imitation learning to bootstrap an RL agent and then use Proximal Policy Optimization (PPO) to further explore global optimal actions. Then, a value network is used to run Monte-Carlo tree search (MCTS) to enhance the policy network. We evaluate the performance of our method on four different categories of combinatorial optimization problems and show that our approach performs strongly compared to the state-of-the-art machine learning and heuristics based methods.

## I. INTRODUCTION

Combinatorial optimization is a broad topic covering several areas of computer science, operations research, and artificial intelligence. The fundamental goal of combinatorial optimization is to find optimal configurations from a finite discrete set that satisfy all given conditions, which involves enormous discrete search spaces. Examples include internet routing [1], scheduling [2], protein structure prediction [3], combinatorial auctions [4]. Many real-life problems can also be formalized as combinatorial optimization problems, including the travelling salesman [5], the vertex colouring [6], and the vehicle routing problems [7], [8]. As combinatorial optimization includes various NP-hard problems, there is a significant demand for efficient combinatorial optimization algorithms.

Several exact combinatorial optimization algorithms have been proposed to provide theoretical guarantees on finding optimal solutions or determining the non-existence of a solution. The core idea is to prune the candidate solution set by confidently introducing new conditions. Branch-and-bound (B&B) [9] is an example of an exact method to solve the combinatorial problem, which recursively divides the candidate solution set into disjoint subsets and rules out subsets that cannot have any candidate solutions satisfying all conditions. It has shown a reliable performance in the domain of mixed-integer linear programs (MILPs) to which many combinatorial problems can be reduced [10]. Several commercial optimization solvers (e.g. CPLEX, Gurobi) use a B&B algorithm to solve

MILP instances. However, two decisions must be made at each iteration of B&B: *node selection* and *variable selection*, which determine the next solution set to be partitioned, and the variable to be used as the partition rule, respectively. Most state-of-the-art optimizers use heuristics hard-coded by domain experts to improve the performance [11]. However, such heuristics are hard to develop and require adjustment for different problems [12].

In recent years, an increasing number of studies have been focusing on training machine learning (ML) algorithms to solve MILP problems. The idea is that some procedural parts of the solvers may be replaced by ML models that are trained with historical data. However, most ML models are trained through supervised learning, which requires the mapping between training inputs and outputs. Since the optimal labels are typically not accessible, supervised learning is not capable for most MILP problems [13]. In contrast, reinforcement learning (RL) algorithms show a potential benefit to the B&B decision-making, thanks to the fact that the B&B decision-making process can be modelled as a Markov decision process (MDP) [12], [14]. This offers an opportunity to use statistical learning for decision-making.

In this work, we provide an RL-based approach to learn a variable selection strategy, which is the core of the B&B method. Our agent is trained to maximize the improvement of dual bound integral with respect to time in the B&B method. We adopt the design of Proximal Policy Optimization (PPO) [15], combining the idea of imitation learning to improve the sample efficiency and advance imitated policy. We imitate the Full Strong Branching (FSB) [16] variable selection rule to discourage the exploration of unpromising directions. We also introduce a Monte Carlo Tree Search (MCTS) like approach [17] to encourage exploration during the training phase and reinforce the action selection strategy.

We evaluate our RL agent with four kinds of widely adopted combinatorial optimization problems. The experiments show that our approach can outperform state-of-the-art methods under multiple metrics. In summary, our contribution is threefold:

- We implement and evaluate an RL-based agent training framework for B&B variable selection problem and achieve comparable performance with the state-of-the-art GCNN approach using supervised learning.
- To facilitate the decision quality, we propose a new MDP formulation that is more suitable for the B&B method.

<sup>§</sup>Work done during an internship at Huawei Technologies Canada Co., Ltd.

- We use imitation learning to accelerate the training process of our PPO agent and propose an MCTS policy optimization method to refine the learned policy.

## II. RELATED WORK

B&B [9] is one of the most general approaches for global optimization in nonconvex and combinatorial problems, which combines partial enumeration strategy with relaxation techniques. B&B maintains a provable upper and lower (primal and dual) bound on the optimal objective value and, hence, provides more reliable results than heuristic approaches. However, the B&B method can be slow depending on the selection of branching rules, which may grow the computational cost exponentially with the size of the problem [18].

Several attempts have been made to derive good branching strategies. Current branching strategies can be categorized into hand-designed approaches that make selections based on heuristic scoring functions; and statistical approaches that use machine learning models to approximate the scoring functions. Most modern MILP solvers use hand-designed branching strategies, including most infeasible branching [19], pseudocost branching (PC) [20], strong branching (SB) [21], [16], reliability branching (RB) [19], and more. Strong branching provides the local optimal solution with the highest computational cost by experimenting with all possible outcomes. Pseudocost branching keeps a history of the success of performed branchings to evaluate the quality of candidate variables, which provides a less accurate but computationally cheaper solution. Reliability branching integrates both strong and pseudocost branching to balance the selection quality and time.

Given the fact that strong branching decisions provide a minimum number of explored nodes among all other hand-designed branching strategies but have a high computational cost, several studies have come up with the idea of approximating and speeding up strong branching strategies using statistical approaches. In [22], a regressor is learned to predict estimated strong branching scores using offline data collected from similar instances. Similarly, a learning-to-rank algorithm that estimates the rank of variables can also provide reliable result [23], [24], which is more reliable than mimicking the score function. However, these statistical approaches suffer from extensive feature engineering.

One common approach to reducing the feature engineering effort is to use the graph convolutional neural network (GCNN). Reference [25] first proposed a GCNN model to solve combinatorial optimization problems, and reference [12] extended the structure to the context of B&B variable selection, which is the closest line of work to ours. In [12], authors show the GCNN can provide accurate estimation of strong branching with the shortest solving time in most of the considered instances.

However, most recent statistical approaches for variable selection in B&B use supervised learning techniques, which require a mapping between training inputs and expected labels. The quality of the model highly depends on the quality of training labels. As mentioned earlier, recent studies use strong branching scores or selections as training labels, which provides

the local optimal solution, but is not guaranteed to be the global optimal solution. In general, we do not have access to optimal labels for most combinatorial optimization problems, and thus the supervised learning paradigm is not suitable in most cases [13]. Another approach is to learn through the interactions with an uncertain environment and provide some reward feedbacks to a learning algorithm. This is also known as the reinforcement learning (RL) paradigm. The RL algorithm makes a sequence of decisions and learns the decision-making policy through trial and error to maximize the long-term reward. Previous studies have shown that the combinatorial optimization problem can be solved using RL algorithms, such as the travelling salesman problem [13], [26], [27], maximum cut problem [28], [29], [30], and more. This study proposes a deep reinforcement learning framework to learn the global optimal variable selection strategy. We adopt the structure of GCNN as the design of our policy and value network.

## III. BACKGROUND

In this section, we describe the fundamental concepts related to the paper, and provide formal definitions to various terms.

### A. Mixed integer linear program (MILP)

A MILP is a mathematical optimization problem that has a set of linear constraints, a linear objective function, and several decision variables that are continuous or integral with the form:

$$\arg \min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}, \quad \text{s.t.} \quad \mathbf{A} \mathbf{x} \leq \mathbf{b}, \quad \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}, \\ x_i \in \mathbb{Z} \text{ where } i \in \mathcal{I}, \quad |\mathcal{I}| \leq n,$$

where  $\mathbf{c}$  is the objective coefficient matrix,  $\mathbf{x}$  is the variable vector,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  denotes the constraint coefficient matrix,  $\mathbf{b} \in \mathbb{R}^m$  represents the constraint constant term vector,  $\mathbf{l} \in (\mathbb{R} \cup \{-\infty\})^n$  and  $\mathbf{u} \in (\mathbb{R} \cup \{\infty\})^n$  indicate the lower and upper variable bound vectors, respectively. Here  $n$ ,  $m$ , and  $\mathcal{I}$  respectively denote the number of variables, number of constraints, and index set of integer variables where  $|\mathcal{I}| \leq n$ . If a variable has no lower or upper bound, then we set the associated  $l$  and  $u$  to infinite values respectively. A *candidate* solution is any assignment of  $\mathbf{x}$  that satisfy the variable bounds. A *feasible* solution is a candidate solution that satisfies all constraints in the MILP instance, and an *optimal* solution is a feasible solution that minimize the objective function.

A MILP can be relaxed to a linear program (LP) by ignoring the integer constraints in the MILP; this is also called *LP relaxation*. LP is convex and therefore can be solved efficiently using various algorithms, such as the simplex algorithm. Since removing the integer constraints expands the feasible set, the optimal solution for LP is then used as the lower bound for the corresponding MILP, namely the *dual bound*.

### B. Branch-and-bound (B&B) algorithm

The B&B algorithm constructs a search tree recursively. Each node in the search tree is a MILP. The B&B algorithm can be described as follows. The original MILP is treated as the root node in the search tree. The algorithm then recursively picks a node from the search tree by a given node selection

rule, picks a variable to decompose the selected node, and adds two children to the selected node that are produced by the decomposition. The dual bound of these two children are then being used to update the dual bound of the root node, and the algorithm selects the next node to expand. To decompose a MILP on variable  $x_i$ , we first find the optimal solution  $\mathbf{x}^*$  to the LP relaxation. Then, if  $\mathbf{x}_i^*$  does not meet the integrality constraint, we can decompose the MILP into two sub-problems with additional constraints  $x_i \leq \lfloor x_i^* \rfloor$  and  $x_i \geq \lceil x_i^* \rceil$ . The variable  $x_i$  is called the *branching variable*, and all variables that can be selected are called *branching candidates*.

1) *Strong branching (SB)*: SB is one of the most powerful state-of-the-art variable selection rules. The idea of SB is to test which branching candidate can provides the best improvement measured in children nodes. This method is a greedy method that selects the locally best variable to branch on, which usually works well in terms of the number of nodes visited to solve the problem. However, it requires to branch on every branching candidates to calculate the score, which is computationally expensive. Moreover, this greedy approach cannot guarantee to provide the global optimal selection.

### C. Markov decision process (MDP) formulation

We can formulate the sequential decision making of variable selection as a MDP. Each node in the search tree can be encoded as a state. The agent exerts a branching variable from all branching candidates to decompose the current node. This action causes a transition to a child node. Through interactions with the MDP, the algorithm learns an optimal *policy*  $\pi$ , that is a sequence of control actions starting from the root node.

1) *State*: The state  $s_t$  at node  $t$ , can be represented as:

$$s_t = \{(X, E, C)_t, J_t\},$$

where the first tuple is the bipartite graph representation  $(X, E, C)_t$  of the current node MILP, as done in [12], and index set  $J_t$  is the index set of branching candidates. Two sets of nodes in the bipartite graph correspond to the  $n$  variable to be optimized and  $m$  constraints to meet. The edge  $e_{i,j}$  is added if the variable  $x_i$  has a non-zero coefficient  $A_{i,j}$  in the constraint  $c_j$ , where  $d_e$  features form the constraints constant term.  $E \in \mathbb{R}^{m \times n \times d_e}$  represents the sparse edge feature matrix.  $X \in \mathbb{R}^{n \times d_x}$  is a feature matrix for all variable nodes, including the features extracted from the objective function and variable constraints. Similarly,  $C \in \mathbb{R}^{m \times d_c}$  represents the feature matrix for all constraint nodes, where each constraint is encoded into  $d_c$  features. Figure 1 illustrates the bipartite representation of a general MILP instance. We calculate the optimal solution of the current node's LP relaxation and mark variables with integer constraints and having a non-integer solution as the branching variable to get  $J_t$ .

2) *Action and transition*: The action at node  $t$ , denoted by  $a_t$ , determines the branching variable from the branching candidates:  $a_t \in J_t$ . After an action is performed, the search tree will add two children nodes to the current node and then prunes the search tree if needed, as described in Section III-B. All children share the same  $p(s_{t+1}|s_t, a_t)$  in this study.

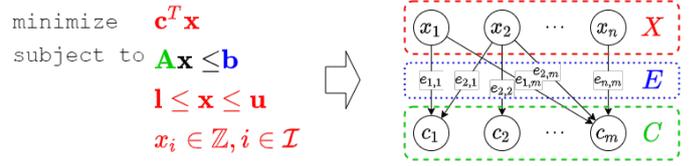


Fig. 1: Bipartite graph representation  $(X, E, C)$  of a MILP.

3) *Reward*: The reward function is designed to encourage the agent increase the dual bound quickly with as less branching operations as possible. Because we do not control the selection of state, and the global dual bound is highly related to the search tree constructed based on the selection of branching node at each step, using the improvement of global dual bound is not valid as the selected branching node might not be able to improve the global dual bound by any action. Therefore, we calculate the reward based on the improvement of the local dual bound:

$$r(s_t, a_t) = \min\{c^T \mathbf{x}_{\lfloor a_t \rfloor_{LP}}^*, c^T \mathbf{x}_{\lceil a_t \rceil_{LP}}^*\} - c^T \mathbf{x}_{s_t}^*,$$

where  $\mathbf{x}_{s_t}^*$  is the dual bound of the current node  $s_t$ ,  $\mathbf{x}_{\lfloor a_t \rfloor_{LP}}^*$  and  $\mathbf{x}_{\lceil a_t \rceil_{LP}}^*$  are respectively the dual bound of children nodes after adding constraint  $x_{a_t} \leq \lfloor x_{a_t}^* \rfloor$  and  $x_{a_t} \geq \lceil x_{a_t}^* \rceil$  to  $s_t$ .

## IV. METHODOLOGY

In this section, we discuss the design of the RL agent, techniques to address the cold-start problem, the training algorithm, and how to exploit the knowledge of a trained RL agent to select branching variables from a given set of branching candidates. Figure 2 shows an overview of our approach, which entails (1) designing the RL agent; (2) using imitation learning to pre-train the RL agent; (3) training the RL agent with PPO; (4) finally, selecting reliable branching variables for test environments using RL agent based on the search result of Monte-Carlo tree search (MCTS); We describe each of these tasks below.

### A. Designing the RL agent

Reinforcement learning methods can find a policy that maximizes the total reward, especially when the MDP is identified. In this study, we use a policy gradient-based method called proximal policy optimization (PPO) to find the optimal policy  $\pi$  using the actor-critic framework. PPO has shown a strong performance in nearly all reinforcement learning tasks, thanks to the clipping method that limits the update of the behaviour policy within a trust region.

To evaluate the step size of the policy gradient method, PPO keeps tracking two policies, current policy  $\pi_\theta$  and old policy  $\pi_{\theta_{old}}$ . Each policy contains two networks: a policy network that estimates the action distribution of a given state and a value network that estimates the state value. The state value  $V(s_t)$  in this study is defined as follow:

$$V(s_t) = \sum_a p(a|s_t) \left( r(s_t, a) + \gamma \frac{V(\lfloor s_t' \rfloor) + V(\lceil s_t' \rceil)}{2} \right),$$

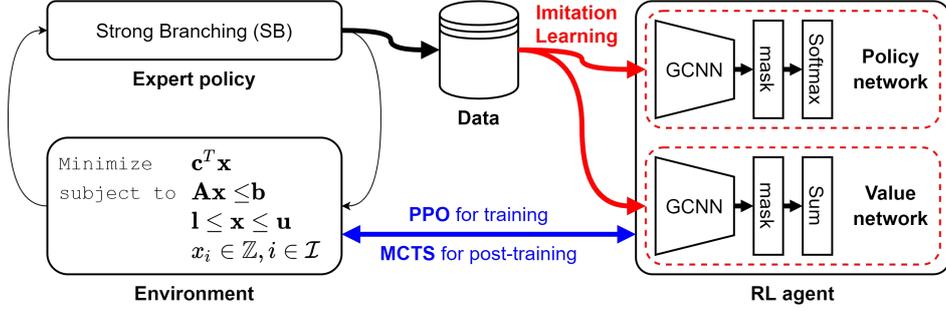


Fig. 2: The flow-diagram of the proposed approach.

where  $\gamma$  is a discount factor with value of 0.99 to encourage immediate reward, and states  $|s'_t|$  and  $|s_t|$  are two children after branching on state  $s_t$  with variable  $x_a$ . This is the different from the next state we obtained through the interaction with the environment. That is being said, state  $s_t$  and  $s_{t+1}$  might not have an edge in the search tree, because the node selection rule pick the next state from all leaf nodes in the search tree. In the calculation of the state value, the next state must be the child of the current state, to therefore correctly represent the state value in the search tree. If  $s_t$  is a leaf node in the search tree, then the state value  $V(s_t)$  is set to 0.

Because the state consists of a bipartite graph, we use graph convolutional neural network (GCNN) as our policy and value network. Previous studies also proved that GCNNs can effectively capture structural characteristics of combinatorial optimization problems. We adopt the similar GCNN design from [12], which use two successive passes to perform a single graph convolution. These passes are

$$c'_p \leftarrow f_c \left( c_p, \sum_{(p,q) \in E} g_c(\text{emb}_x(x_q), e_{p,q}, \text{emb}_c(c_p)) \right),$$

$$x'_q \leftarrow f_x \left( x_q, \sum_p g_x(\text{emb}_x(x_q), e_{p,q}, c'_p) \right),$$

for all  $p \in C, q \in X$ . Next, the value of  $x'$  is sent to a 2-layer perceptron. For the policy network, we apply masked softmax activation to estimate the action distribution. For the value network, we compute masked sum to predict the state value.

### B. Imitating the Strong Branching (SB)

Theoretically, the RL agent can find the optimal policy  $\pi$  from scratch after training for enough episodes. However, as the search tree is huge, with a branching factor usually over 1,000, training an RL agent from scratch becomes time-consuming and therefore impractical. To avoid the initial aimless exploration of the RL agent, we use the imitation learning approach to pretrain the RL agent policy and value network, paving the way for learning sophisticated policy. We select SB as our expert policy to generate offline training data, including the state, corresponding SB score for each branching candidate, as well as the reward. Then we reconstruct the state value  $V(s_t)$

from the offline data and pretrain the policy and value network by minimizing the following loss:

$$L^{policy}(\theta) = -\frac{1}{N} \sum_{s,a \in \mathcal{D}} \log \pi_\theta(a|s)$$

$$L^{value}(\theta) = \frac{1}{N} \sum_{s \in \mathcal{D}} (V_\theta(s) - V(s))^2$$

### C. Training the RL agent

Once the RL agent is pretrained using offline data, it is necessary to learn an advanced policy by interacting with the environment directly. To update the policy parameter  $\theta$  with some trajectories generated through the interaction with the environment, we first save the parameters  $\theta$  into  $\theta_{old}$ , and then calculate the loss as follows:

$$A_t = V(s_t) - V_\theta(s_t), \quad r = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)},$$

$$L_t^{policy}(\theta) = \mathbb{E}_t [\min(rA_t, \text{clip}(r, 1 - \epsilon, 1 + \epsilon)A_t)],$$

$$L_t(\theta) = \mathbb{E}_t \left[ L_t^{policy}(\theta) - c_1 A_t^2 - c_2 \sum_a \pi_\theta(a|s_t) \log \pi_\theta(a|s_t) \right],$$

where  $V(s_t)$  is the state value calculated from the experiences,  $V_\theta(s_t)$  is the value network estimated state value for  $s_t$ , and  $\epsilon, c_1, c_2$  are hyperparameters of the model. In this study, we use  $\epsilon = 0.1, c_1 = 0.5, c_2 = 0.01$ .

### D. Enhancing policy with Monte-Carlo Tree Search (MCTS)

After we obtain the stable policy  $\pi_\theta^*$  and  $V_\theta^*$ , it is essential to make reliable selections for a given state  $s_t$ . One common and straightforward approach is to take action with the highest estimated action probability,  $\arg \max_a \pi_\theta^*(a|s_t)$ . However, this result could be biased when the policy is not optimized or has a significant variance. Since we have a tree-like search space, it is possible to adopt the idea of MCTS to update the policy  $\pi_\theta^*$  further. In MCTS, we generate multiple trajectories starting from the current node. Then, we expand trajectories by taking action based on some probability distribution until a certain number of steps or the final state is reached. Finally, we pick the best action based on these trajectories. This is similar to the SB, except MCTS does not explore all branching candidates.

To run MCTS efficiently, we incorporate the knowledge learned by our RL agent. In the action selection step for MCTS,

we use a modified version of upper confidence bound (UCB), which selects the action that maximizes the following equation:

$$\arg \max_{a \in \mathcal{A}(s)} (s, a) + c\pi_{\theta}^*(a|s) \sqrt{\frac{\log(1 + \sum_a N(s, a))}{N(s, a) + 1}}.$$

Here, the  $Q(s, a)$  represents the action value based on the trajectories done previously, and  $N(s, a)$  keeps tracking the number of times  $a$  has been selected on state  $s$ . We introduce trained policy  $\pi_{\theta}^*$  to encourage the algorithm to search for promising directions. To minimize the size of the search tree, we further limit the branching candidates on each state  $s$  to  $\mathcal{A}(s)$ , which only contains the top  $k$  actions based on the  $\pi_{\theta}^*(s)$ . In this study, we use  $k = 10$ .

Also, to reduce the simulation time, we do not perform the branching operation when we run MCTS. Instead, we directly modify the constraint feature matrix and edge feature matrix to simulate the next state  $s'$  based on the action, with half chance to reach the left child and half chance to reach the right child. We then set the reward of all actions to 0 and use the value network  $V_{\theta}^*$  trained by the RL agent to calculate the value of  $V_{\theta}^*(s')$ , and use it to calculate the  $Q(s, a)$ . We initialize the action value and the visit count as follow:

$$Q(s, a) = \gamma V_{\theta}^*(s'_t), \quad N(s, a) = 1.$$

The  $Q(s, a)$  and  $N(s, a)$  are then updated when the agent reaches the leaf of the search tree or the maximum number of steps is reached. We apply the following update rule for each state  $\{s_t, \dots, s_0\}$ :

$$Q(s_{\tau}, a_{\tau}) \leftarrow Q(s_{\tau}, a_{\tau}) + \frac{-Q(s_{\tau}, a_{\tau}) + \sum_{t'=\tau}^t \gamma^{t-t'} V_{\theta}^*(s'_{t'})}{N(s_{\tau}, a_{\tau}) + 1}.$$

After all MCTS simulations are finished, we identify the action  $a$  with the highest  $Q(s, a)$  as the best branching variable for each state  $s$  that has been visited at least ten times and use this to train the policy network by minimizing the cross-entropy loss. In this study, we limit the maximum depth to 3 and run 1,000 simulations of the MCTS for each state.

## V. EVALUATION

In this section, we study the efficacy of different variable selection strategies. We adopt the average solving time, average number of resulting B&B nodes, and average dual integral as our evaluation metrics. All experiments are repeated five times with different random seeds to eliminate randomness. All numbers are the averaged value across all five runs.

### A. Data sets

To test the generalizability of our framework, we evaluate our approach on four different types of NP-hard problems. The first problem is called set covering problem proposed in [31]. Our instances contain 1,000 columns and 500 rows per instance. The second problem is generated following the arbitrary relationships procedure of [32]. This problem is also known as the combinatorial auction problem. In our experiment, we generate instances with 100 items for 500 bids. Our third data set is called capacitated facility location described in [33].

We collect instances with 100 facilities and 100 customers. The last data set we used in this study is proposed in [34], which is called the maximum independent set problem. The affinity is set to 4, and the graph size is set to 500 in this study.

These problems are selected based on the previous works and the hardness of the problem itself. According to [12], these problems are the most representative of the types of integer programming problems encountered in practice. We use SCIP 7.0.3 [35] as the backend solver throughout the study, with ecole 0.7.3 [36] as the environment interface. All SCIP parameters are kept to default in this study.

### B. Baselines

In the rest of this paper, we use PPO-MCTS to refer to our proposed reinforcement learning framework. We compare our approach with three different variable selection baseline strategies. The first naive baseline strategy is the pure random strategy, in which we select the branching variable from a set of candidate variables uniformly. We use the full strong branching (FSB) strategy as our second baseline, which we use the default parameter defined in SCIP in this study. We also re-implemented the GCNN model from [12] as our third baseline. Based on the ML4CO NeurIPS 2021 competition result, the GCNN model yields the best performance among all other competing methods [37], [38], [39]. The performance of PPO agents that have no MCTS learning afterwards (PPO) is also reported for ablation study.

### C. Evaluation metrics

We evaluate the performance of each approach using three metrics, including the average solving time for each problem instance, the average number of B&B search tree nodes visited before the problem is solved, and a reward score that takes into account both the solving time and the improvement of the dual bound. Solving time and the number of nodes visited measure the computational cost of each algorithm. Solving time is evaluated based on the wall clock time, including feature extraction time, model inference time, branching time, and more. Therefore, a shorter solving time does not guarantee to optimize the number of nodes visited during the B&B method. To optimize the branching variable selection strategy, we expect to minimize the number of nodes visited during the branching and the total solving time to select optimal branching variables with minimum computational cost. The score is calculated by:

$$-T\mathbf{c}^T \mathbf{x}^* + \int_{t=0}^T \mathbf{z}_t^* \partial t, \quad (1)$$

where  $\mathbf{x}^*$  is the optimal solution of the MILP instance,  $T$  is the time budget to solve the problem, and  $\mathbf{z}_t^*$  is the best dual bound at time  $t$ . This score is to be maximized, representing a fast improvement of the dual bound. This reward metric was first introduced in the ML4CO NeurIPS 2021 competition [37] and is expected to be adopted further by the community.

TABLE I: Number of resulting B&amp;B nodes on the test data sets

	Set Covering	Independent Set	Combinatorial Auction	Capacitated Facility Location
Random	2225.20	257.60	25543.26	2292.24
FSB	47.42	103.85	<b>193.83</b>	<b>47.9</b>
GCNN	44.07	<b>88.66</b>	201.82	886.66
PPO-MCTS	<b>43.90</b>	90.23	194.25	863.21

TABLE II: MILP instance solving time (in seconds) on the test data sets

	Set Covering	Independent Set	Combinatorial Auction	Capacitated Facility Location
Random	19.2	15.36	99.08	92.28
FSB	95.8	240.65	113.58	864.75
GCNN	3.28	<b>6.54</b>	4.15	80.68
PPO-MCTS	<b>3.13</b>	6.60	<b>3.87</b>	<b>77.24</b>

TABLE III: Evaluation score on the test data sets

	Set Covering	Independent Set	Combinatorial Auction	Capacitated Facility Location
FSB	149930	-191876	-7093620	16119158
GCNN	<b>150654</b>	<b>-191123</b>	-7077028	16159789
PPO-MCTS	150652	-191139	<b>-7076023</b>	<b>16160324</b>

TABLE IV: Performance comparison between PPO and PPO-MCTS

	Nodes visit		Solving time		Evaluation score	
	PPO	PPO-MCTS	PPO	PPO-MCTS	PPO	PPO-MCTS
Set Covering	57.68	<b>43.90</b>	3.52	<b>3.13</b>	150651	<b>150652</b>
Independent Set	<b>88.18</b>	90.23	8.34	<b>6.60</b>	-203328	<b>-191139</b>
Combinatorial Auction	270.97	<b>194.25</b>	5.28	<b>3.87</b>	-7077229	<b>-7076023</b>
Capacitated Facility Location	2202.46	<b>863.21</b>	139.65	<b>77.24</b>	16155103	<b>16160324</b>

#### D. Experiment result

Table I shows the average number of nodes visited before solving the instances for each approach. We noticed both GCNN and our PPO-MCTS have more nodes visited in complex problems, namely the combinatorial auction and capacitated facility location problems, compared to FSB. We conclude this to the fact that the number of branching candidates in these two problems are more significant than the other two problem types, and therefore leads to the approximate function getting more complex, which lowers the performance of the GCNN model. Similarly, as the search tree branching factor grows, RL agents become more challenging to learn the environment thoroughly. The agent can potentially struggle in local optimal as the maximum depth is set to three for our PPO-MCTS agent. It is worth noting that number of nodes on its own is not enough of a measure to judge different approaches with. There reason is that a method may visit a larger number of nodes, but may in fact be faster on each visit and result a better overall reward value for convergence.

On the other hand, it is readily seen from Table II that FSB takes a significant amount of time to select a variable for each node in these two problems, and therefore the total solving time for FSB is the longest compared to all other approaches in all four problems. The GCNN and PPO-MCTS are having similar inference time as the network designs are similar. As the PPO-MCTS has a lower average number of visited nodes in all but independent set problems, our method provides the shortest solving time in all problems except the independent set problem. However, the performance differences between GCNN and PPO-MCTS on independent set problem in all three metrics are negligible, which proves the effectiveness of our proposed framework. In addition, when the problem is easy,

such as the set covering and independent set problems, both GCNN and PPO-MCTS can find better branching variables with fewer nodes visited than FSB. In general, PPO-MCTS has a slightly better performance across different data sets and metrics than GCNN, with a trade-off on the training expenses.

The average evaluation score for each approach is shown in Table III. The GCNN and PPO-MCTS approaches keep dominating the score in all problems, whereas the PPO-MCTS has a higher score on challenging problems, and GCNN performs better with easy problems.

#### E. Ablation study

We present an ablation study of the proposed PPO-MCTS model to evaluate the importance of having post MCTS retraining. Table IV demonstrates the performance of PPO without post MCTS retraining on all metrics across all data sets, denoted by PPO, comparing with the proposed PPO-MCTS. It is observed that the PPO-MCTS perform better than PPO in all cases, except the number of nodes visited in the independent set problem. This empirical result suggests that the post MCTS retraining offers a better performing RL agent.

## VI. CONCLUSION

We proposed a reinforcement learning framework to learn the variable selection policy for the B&B method. We formulated a reward function that helps the agent learn optimal policies without generating labels. We used imitation learning and MCTS to deal with sample inadequacy challenges by initializing the policy to a relatively good policy and enhancing it with multiple steps look-ahead. We demonstrated the performance of the proposed framework with three baseline approaches on four NP-hard problems and showed that the proposed method yields a strong performance in most problems.

## REFERENCES

- [1] João CN Clímaco, Marta MB Pascoal, José MF Craveirinha, and M Eugénia V Captivo. Internet packet routing: Application of a k-quickest path algorithm. *European Journal of Operational Research*, 181(3):1045–1054, 2007.
- [2] Yves Crama. Combinatorial optimization models for production scheduling in automated manufacturing systems. *European Journal of Operational Research*, 99(1):136–153, 1997.
- [3] Mehul M Khimasia and Peter V Coveney. Protein structure prediction as a hard optimization problem: the genetic algorithm approach. *Molecular Simulation*, 19(4):205–226, 1997.
- [4] Tuomas Sandholm. Algorithm for optimal winner determination in combinatorial auctions. *Artificial intelligence*, 135(1-2):1–54, 2002.
- [5] Pedro Larranaga, Cindy M. H. Kuijpers, Roberto H. Murga, Inaki Inza, and Sejla Dizdarevic. Genetic algorithms for the travelling salesman problem: A review of representations and operators. *Artificial intelligence review*, 13(2):129–170, 1999.
- [6] Moustapha Diaby. Linear programming formulation of the vertex colouring problem. *International Journal of Mathematics in Operational Research*, 2(3):259–289, 2010.
- [7] Paolo Toth and Daniele Vigo. *The vehicle routing problem*. SIAM, 2002.
- [8] Bruce L Golden, Subramanian Raghavan, and Edward A Wasil. *The vehicle routing problem: latest advances and new challenges*, volume 43. Springer Science & Business Media, 2008.
- [9] Ailsa H Land and Alison G Doig. An automatic method of solving discrete programming problems. *Econometrica*, 28(3):497–520, 1960.
- [10] Eugene L Lawler and David E Wood. Branch-and-bound methods: A survey. *Operations research*, 14(4):699–719, 1966.
- [11] Ambros Gleixner, Gregor Hendel, Gerald Gamrath, Tobias Achterberg, Michael Bastubbe, Timo Berthold, Philipp Christophel, Kati Jarck, Thorsten Koch, Jeff Linderoth, et al. Miplib 2017: data-driven compilation of the 6th mixed-integer programming library. *Mathematical Programming Computation*, pages 1–48, 2021.
- [12] Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. Exact combinatorial optimization with graph convolutional neural networks. *arXiv preprint arXiv:1906.01629*, 2019.
- [13] Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- [14] Nina Mazyavkina, Sergey Sviridov, Sergei Ivanov, and Evgeny Burnaev. Reinforcement learning for combinatorial optimization: A survey. *Computers & Operations Research*, page 105400, 2021.
- [15] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [16] Jeff T Linderoth and Martin WP Savelsbergh. A computational study of search strategies for mixed integer programming. *INFORMS Journal on Computing*, 11(2):173–187, 1999.
- [17] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [18] Stephen Boyd and Jacob Mattingley. Branch and bound methods. *Notes for EE364b, Stanford University*, pages 2006–07, 2007.
- [19] Tobias Achterberg, Thorsten Koch, and Alexander Martin. Branching rules revisited. *Operations Research Letters*, 33(1):42–54, 2005.
- [20] Michel Bénéichou, Jean-Michel Gauthier, Paul Girodet, Gerard Hentges, Gerard Ribière, and O Vincent. Experiments in mixed-integer linear programming. *Mathematical Programming*, 1(1):76–94, 1971.
- [21] Jan Karel Lenstra and David Shmoys. The traveling salesman problem: A computational study, 2009.
- [22] Alejandro Marcos Alvarez, Quentin Louveaux, and Louis Wehenkel. A machine learning-based approximation of strong branching. *INFORMS Journal on Computing*, 29(1):185–195, 2017.
- [23] Elias Khalil, Pierre Le Bodic, Le Song, George Nemhauser, and Bistra Dilkina. Learning to branch in mixed integer programming. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.
- [24] Christoph Hansknecht, Imke Joormann, and Sebastian Stiller. Cuts, primal heuristics, and learning to branch for the time-dependent traveling salesman problem. *arXiv preprint arXiv:1805.01415*, 2018.
- [25] Hanjun Dai, Elias B Khalil, Yuyu Zhang, Bistra Dilkina, and Le Song. Learning combinatorial optimization algorithms over graphs. *arXiv preprint arXiv:1704.01665*, 2017.
- [26] Hanjun Dai, Bo Dai, and Le Song. Discriminative embeddings of latent variable models for structured data. In *International conference on machine learning*, pages 2702–2711. PMLR, 2016.
- [27] Hao Lu, Xingwen Zhang, and Shuang Yang. A learning-based iterative method for solving vehicle routing problems. In *International Conference on Learning Representations*, 2019.
- [28] Thomas Barrett, William Clements, Jakob Foerster, and Alex Lvovsky. Exploratory combinatorial optimization with reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3243–3250, 2020.
- [29] Quentin Cappart, Emmanuel Goutierre, David Bergman, and Louis-Martin Rousseau. Improving optimization bounds using machine learning: Decision diagrams meet deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 1443–1451, 2019.
- [30] Kenshin Abe, Zijian Xu, Issei Sato, and Masashi Sugiyama. Solving np-hard problems on graphs with extended alphago zero. *arXiv preprint arXiv:1905.11623*, 2019.
- [31] Egon Balas and Andrew Ho. Set covering algorithms using cutting planes, heuristics, and subgradient optimization: a computational study. In *Combinatorial Optimization*, pages 37–60. Springer, 1980.
- [32] Kevin Leyton-Brown, Mark Pearson, and Yoav Shoham. Towards a universal test suite for combinatorial auction algorithms. In *Proceedings of the 2nd ACM conference on Electronic commerce*, pages 66–76, 2000.
- [33] Gérard Cornuéjols, Ranjani Sridharan, and Jean-Michel Thizy. A comparison of heuristics and relaxations for the capacitated plant location problem. *European journal of operational research*, 50(3):280–297, 1991.
- [34] David Bergman, Andre A Cire, Willem-Jan Van Hoes, and John Hooker. *Decision diagrams for optimization*, volume 1. Springer, 2016.
- [35] Gerald Gamrath, Daniel Anderson, Ksenia Bestuzheva, Wei-Kun Chen, Leon Eifler, Maxime Gasse, Patrick Gemander, Ambros Gleixner, Leona Gottwald, Katrin Halbig, et al. The scip optimization suite 7.0. 2020.
- [36] Antoine Prouvost, Justin Dumouchelle, Lara Scavuzzo, Maxime Gasse, Didier Chételat, and Andrea Lodi. Ecole: A gym-like library for machine learning in combinatorial optimization solvers. In *Learning Meets Combinatorial Algorithms at NeurIPS2020*, 2020.
- [37] ecole.ai. ML4CO: 2021 neurips competition on machine learning for combinatorial optimization. <https://www.ecole.ai/2021/ml4co-competition/>.
- [38] Zixuan Cao, Yang Xu, Zhewei Huang, and Shuchang Zhou. MI4co-kida: Knowledge inheritance in dataset aggregation. *arXiv preprint arXiv:2201.10328*, 2022.
- [39] Amin Banitalebi-Dehkordi and Yong Zhang. MI4co: Is gcnn all you need? graph convolutional neural networks produce strong baselines for combinatorial optimization problems, if tuned and trained properly, on appropriate data. *arXiv preprint arXiv:2112.12251*, 2021.