



CIRRELT

Centre interuniversitaire de recherche
sur les réseaux d'entreprise, la logistique et le transport

Interuniversity Research Centre
on Enterprise Networks, Logistics and Transportation

RLBS: An Adaptive Backtracking Strategy Based on Reinforcement Learning for Combinatorial Optimization

Ilyess Bachiri
Jonathan Gaudreault
Brahim Chaib-draa
Claude-Guy Quimper

February 2015

CIRRELT-2015-07

Bureaux de Montréal :
Université de Montréal
Pavillon André-Aisenstadt
C.P. 6128, succursale Centre-ville
Montréal (Québec)
Canada H3C 3J7
Téléphone : 514 343-7575
Télécopie : 514 343-7121

Bureaux de Québec :
Université Laval
Pavillon Palasis-Prince
2325, de la Terrasse, bureau 2642
Québec (Québec)
Canada G1V 0A6
Téléphone : 418 656-2073
Télécopie : 418 656-2624

www.cirrelt.ca

RLBS: An Adaptive Backtracking Strategy Based on Reinforcement Learning for Combinatorial Optimization

Ilyess Bachiri^{1,2,*}, Jonathan Gaudreault^{1,2}, Brahim Chaib-draa²,
Claude-Guy Quimper²

¹ Interuniversity Research Centre on Enterprise Networks, Logistics and Transportation (CIRRELT)

² Département d'informatique et de génie logiciel, 1065, avenue de la Médecine, Université Laval, Québec, Canada G1V 0A6

Abstract. Combinatorial optimization problems are often very difficult to solve and the choice of a search strategy has a tremendous influence over the solver's performance. A search strategy is said to be adaptive when it dynamically adapts to the structure of the problem instance and identifies the areas of the search space that contain good solutions. We introduce an algorithm (RLBS) that learns to efficiently backtrack when searching non-binary trees. Branching can be carried on using any variable/value selection strategy. However, when backtracking is needed, the selection of the target node involves reinforcement learning. As the trees are non-binary, we have the opportunity to backtrack many times to each node during the search, which allows learning which nodes generally lead to the best rewards (that is, to the most interesting leaves). RLBS is evaluated for a scheduling problem using real industrial data. It outperforms classic (non-adaptive) search strategies (DFS, LDS) as well as an adaptive branching strategy (IBS).

Keywords: Search, optimization, learning, backtracking.

Acknowledgements. This work has been supported by the FORAC Research Consortium industrial partners.

Results and views expressed in this publication are the sole responsibility of the authors and do not necessarily reflect those of CIRRELT.

Les résultats et opinions contenus dans cette publication ne reflètent pas nécessairement la position du CIRRELT et n'engagent pas sa responsabilité.

* Corresponding author: Ilyess.Bachiri@cirrelt.ca

1 Introduction

Combinatorial optimization problems are often very difficult to solve and the choice of a search strategy has a tremendous influence over the solver’s performance. To solve a problem using search, one needs to choose a variable selection strategy (defining the order in which variables will be instantiated), a value selection strategy (defining the sequence in which we will try the variable possible values) and a backtracking strategy (that determines to which node we should backtrack/backjump, when a leaf is reached or a dead-end is encountered). Some backtracking policies are encoded into full deterministic algorithms (e.g. Depth-First Search, DFS) while others rely on more dynamic node evaluator mechanisms (e.g. Best-First Search). Others (e.g. Limited Discrepancy Search [9]) can be implemented as a deterministic iterative algorithm or as a node evaluator [3].

A strategy is said to be *adaptive* when it dynamically adapts to the structure of the problem and identifies the areas of the search space that contain good solutions. Some have proposed adaptive branching strategies (e.g. Impact-based Search (IBS) [17]) or a backtracking strategy (e.g. Adaptive Discrepancy Search [7], proposed for distributed optimization problems).

In this paper, we consider a machine learning approach which improves the performance of the solver. More specifically, we use Reinforcement Learning (RL) to identify the areas of the search space that contain good solutions. The approach was developed for optimization problems for which the search space is encoded as a non-binary tree. As the trees are non-binary, we have the opportunity to backtrack multiple times to each node during the search. This allows learning which nodes generally lead to the best rewards (that is, to the most interesting leaves).

Section 2 reviews some preliminary concepts regarding adaptive search and reinforcement learning. Section 3 explains how backtracking can be encoded as a reinforcement learning task and introduces the proposed algorithm (*Reinforcement Learning Backtracking Search*, or RLBS). Section 4 presents results for a complex industrial problem that combines planning and scheduling. RLBS is compared to more classic (non-adaptive) search strategies (DFS, LDS) as well as an other adaptive branching strategy (IBS). Section 5 concludes the paper.

2 Background

2.1 Learning Variable/Value Selection Heuristics

Some algorithms learn during the search which variables are the most difficult to instantiate, in order to dynamically change the order of the variables (e.g. YIELDS [10]). In [4] and [8], each time a constraint causes a failure, the priority of the variables involved in this constraint is increased.

In Impact Based Search (IBS) [17], the impact of the variables is measured by observing how their instantiation reduces the size of the search space. Since IBS picks the variable to assign and the value to try all at once, it can be considered learning a combination of a variable and value ordering strategies.

2.2 Learning to Backtrack

Approaches where the system learns to evaluate the quality of the nodes are of particular interest for backtracking strategies. Ruml [18] makes an interesting proposal regarding this. While a basic LDS policy gives the same importance to any discrepancy, Best Leaf First Search (BLFS) dynamically attributes different weights to discrepancies according to their depth. BLFS uses a linear regression in order to establish the value of the weights. The model was not really used in order to define a backtracking strategy. Instead, the search algorithm proceeds by a series of successive descents in the tree. Ruml has achieved very good results with this algorithm (see [19]). It was the inspiration for the following algorithm.

Adaptive Discrepancy Search (ADS) [7] is an algorithm that was proposed for distributed optimization but it could be used in a classic COP context. During the search, it dynamically learns which nodes it pays the most to backtrack to (in order to concentrate on those areas of the tree first). For each node, it tries learning a function $Improvement(i)$ predicting how good would be the first leaf reached after backtracking to this node for the i -th time, in comparison to previous backtracks to the same node. The drawback of this method is that a function needs to be learned for each open node, and updated each time it leads to a new solution [13] (although an approximation can be computed using regression).

The algorithm in Section 3 introduces a simplified learning mechanism based on a basic reinforcement learning technique.

2.3 Reinforcement Learning

The fundamental idea of Reinforcement Learning (RL) is to figure out a way to map actions to situations in order to maximize the total reward. The learner is not told which actions to take, it must discover by itself which actions lead to the highest reward (at long-term). Actions may affect not only the immediate reward but also the next situation and, through all, all subsequent rewards [2]. Moreover, the actions may not lead to the expected result due to the uncertainty of the environment.

RL uses a formal framework defining the interaction between the learner and the environment in terms of states, actions, and rewards. The environment that supports RL is typically formulated as a finite-state Markov Decision Process (MDP). In each state $s \in S$, a set of actions $a \in A$ are available to the learner, among which it has to pick the one that maximizes the cumulative reward. The evaluation of actions is entirely based on the learner's experience, built through its interactions with the environment. The goal of the learner is to find, through its interactions with the environment, an optimal policy $\pi : S \rightarrow A$ maximizing the cumulative reward. The cumulative reward is either expressed as a sum of all the rewards $R = r_0 + r_1 + \dots + r_n$ or as a discounted sum of the rewards $R = \sum_t \gamma^t r_t$. The discount factor $0 \leq \gamma \leq 1$ is applied to promote the recent rewards. The discounted sum representation of the cumulative reward is mostly used for an MDP with no terminal state.

A central intuition underlying reinforcement learning is that actions that lead to large rewards should be made more likely to recur.

In a RL task, each action a , in each state s , is associated with a numeric value $Q(s, a)$ that represents the desirability to take the action a in the state s . These values are called Q -Values. The higher the Q -Value, the more likely the action is going to lead to a good solution, according to the learner's judgment. Every time a reward is returned to the learner, the learner must update the Q -Value of the action that has led to this reward. However, the older Q -Value should not be completely forgotten, otherwise the learner would be acting based on the very last experience every single time. To do so, we keep a part of the old Q -Value and we update it with a part of the new experience. Also, we assume that the learner is going to act optimally afterward. Moreover, the expected future rewards need to be discounted to express the idea of the sooner a reward is received, the better.

Let s be the current state, s' the next state, a an action, r the returned reward after having taken the action a , α the learning rate, and γ the discount factor. The update formula for the Q -Values is as follows:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a')] \quad (1)$$

This update formula comes in handy when the learner has to learn an action-value representation, like in Q -Learning [20].

2.4 Reinforcement Learning and Search

The idea of using RL in solving combinatorial problems is supported by many publications [14, 16, 21]. Some researches tried to apply RL to solve optimization problems and some have considered solving Constraint Satisfaction Problems (CSP) using RL techniques.

For instance, Xu et al. [21] proposed a formulation of a CSP as a RL task. A set of different variable ordering heuristics is provided to the algorithm that learns which one to use, and when to use it, in order to solve a CSP in a shorter amount of time. The learning process is accomplished in an offline manner and applied on different instances of the same CSP. The states are the instances or sub-instances of the CSP and the actions are defined as the variable ordering heuristics. A reward is assigned each time an instance is solved. This approach relies on Q -learning to learn the optimal variable ordering heuristic at each decision point of the search tree, for a given (sub)-instance of the CSP.

Moreover, Loth et al. [11] have proposed the Bandit Search for Constraint Programming (BASCOP) algorithm that guides the exploration in the neighborhood of the previous best solution, based on statistical estimates gathered across multiple restarts. BASCOP has been applied on a job shop problem in [12] and has been shown to match the CP-based state of the art.

A local search technique using RL is also proposed in [14]. This approach aims at solving COPs based on a population of RL agents. The pairs (variable, value) are considered as the RL task states, and the branching strategies as the actions.

Each RL agent is assigned a specific area of the search space where it has to learn and find good local solutions. The expertise of the entire population of RL agents is used. A new solution is produced by taking a part of the locally-best solution found by one agent, and complete the remaining assignments using the expertise of another agent.

According to [16], local search can be seen as a policy of a Markov Decision Process (MDP) where states represent solutions, and actions define neighboring solutions. Reinforcement learning techniques can then be used to learn a cost function in order to improve local search. One way to do so is to learn a new cost function over multiple search trajectories of the same problem instance. Boyan and Moore’s STAGE algorithm [5] follows this approach and alternates between using the learned and the original cost function. By enhancing the predictive accuracy of the learned cost function, the guidance of the heuristics improves as the search goes on.

Another approach that uses reinforcement learning to improve local search in the context of combinatorial optimization is to learn a cost function off-line, and then use it on new instances of the same problem. Zhang and Dietterich’s work [22] falls into this category.

3 RLBS: Backtracking as a Reinforcement Learning Task

This section introduces *Reinforcement Learning Backtracking Search* (RLBS). Branching is performed according to any variable/value selection heuristic. Each time we reach a leaf/solution, we need to select the node to backtrack to. To each available candidate (node with at least one unvisited child) corresponds a possible *action* (“backtracking to this node”). Once we select a node, the search continues from that point until we reach a new leaf/solution. The difference between the quality of this new solution and the best solution so far is the *reward* we get for performing the previous action. As our trees are non-binary, we backtrack multiple times to each node during the search. This is an opportunity to identify the actions that pay the most (that is, nodes that are more likely to lead to interesting leaves/solutions).

This situation reminds the k-armed-bandit problem [1]. It is a single-state reinforcement learning problem. Many actions are possible (pulling one of the arms/levels of the slot machine). Each action may lead to a reward (which is stochastic) and we need balancing between exploration and exploitation. In our specific backtracking situation, performing an action makes us discover new nodes/actions, in addition to giving us a reward (which is stochastic and non-stationary).

3.1 Learning

As in classic reinforcement learning, the valuation (Q -value) of an action a is updated each time we get a reward after performing the action. As we are in

a single-state environment, the discount factor γ is equal to 0 and equation (1) reduces to equation (2):

$$Q(a) \leftarrow Q(a) + \alpha[r(a) - Q(a)] \quad (2)$$

where $r(a)$ is the reward and α is the learning rate.

The next action to perform is selected based on those valuations. A node that paid well at first but never got good solutions afterward will see its Q -value decrease over time, until it becomes less interesting than other nodes/actions.

3.2 Initialization of the Algorithm

At the beginning of the search, we descend to the first leaf/solution of the tree using a DFS. We then backtrack once to each open node (this is similar to performing the first 2 iterations of LDS), which allows computing their Q -Values. Then, we start using the Q -Values in order to choose the next node to backtrack to. Each time a new node is visited for the first time, its Q -Value is initialized using the value of its parent.

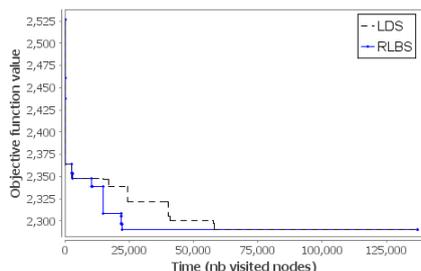


Fig. 1. Solution objective function value according to computation time of LDS and RLBS for case #1

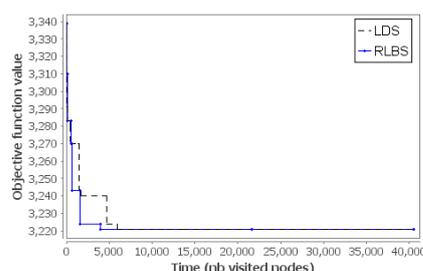


Fig. 2. Solution objective function value according to computation time of LDS and RLBS for case #2

4 Experimentation using Industrial Data

We carried out experiments for a combined planning and scheduling problem from the forest-products industry (lumber planning and scheduling problem). The problem is difficult as it involves *divergent processes with coproduction*: a single process produces many different products at the same time, from a single type of raw material. Moreover, many alternative processes can produce the same product. Finally, it involves complex setup rules. The objective is to minimize orders lateness.

The problem is fully described in [6] which provides a good variable/value selection heuristic specific for it. In [7], this heuristic was used to guide the

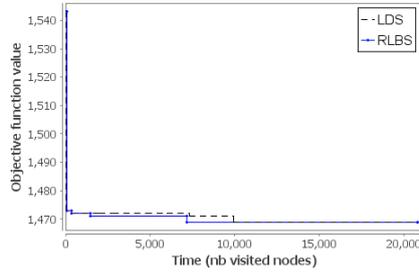


Fig. 3. Solution objective function value according to computation time of LDS and RLBS for case #3

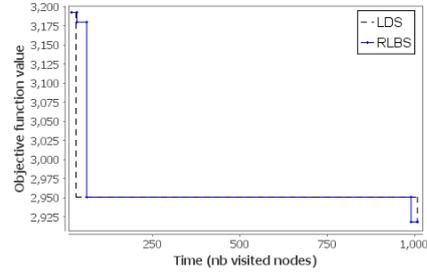


Fig. 4. Solution objective function value according to computation time of LDS and RLBS for case #4

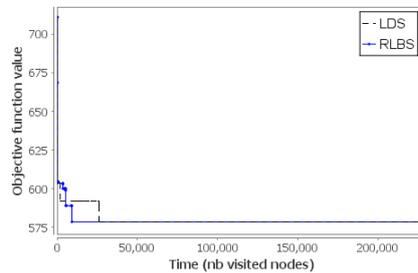


Fig. 5. Solution objective function value according to computation time of LDS and RLBS for case #5

search in a constraint programming model. Provided with this heuristic, LDS outperformed DFS as well as a mathematical programming approach. In [15], parallelization was used to improve performance: however, the visiting order of the nodes is the same as the centralized version, so it implements the same strategy.

Table 1. Computation time needed to get the best solution (RLBS vs. LDS)

	Case 1	Case 2	Case 3	Case 4	Case 5	Average
LDS	57926	5940	9922	1008	26166	20192.4
RLBS	22164	3949	7172	990	9545	8764
Time reduction	↓ 61.73%	↓ 33.52%	↓ 27.72%	↓ 1.79%	↓ 63.52%	↓ 37.66%

We used the same variable/value selection heuristic as in previous work. We also used the same industrial data provided by a Canadian forest-products company. However, in order to be able to compare the algorithms according to

Table 2. Average computation time to get a solution of a given quality (RLBS vs. LDS)

	Case 1	Case 2	Case 3	Case 4	Case 5	Average
LDS	8777.82	1243.86	386.5	127.34	2776.05	2662.31
RLBS	4254.81	606.79	264.24	152.17	1320.15	1319.63
Time reduction	↓ 51.53%	↓ 51.22%	↓ 31.63%	↑ 19.5%	↓ 52.44%	↓ 33.46%

the time needed to get optimal solutions, we reduced the size of the problems (5 periods instead of 44 periods).

RLBS was evaluated using a learning rate $\alpha = 0.5$. We compared the algorithm to an LDS-based policy (selecting the node showing the least discrepancies).

Figures 1 to 5 present the results for five different cases. Table 1 shows the reduction of computation time (measured as the number of visited nodes) needed to get an optimal solution. RLBS reduced computation time for each case (on average by 37.66%).

As in industrial context we usually do not have time to wait for the optimal solution, we also wanted to consider the time needed to get solutions of intermediate qualities. Table 2 shows, for each case, the average time needed to get a solution of any given quality. The last column shows that on average, for all problems and all needed solution qualities, the expected improvement of computation time is 33.46%.

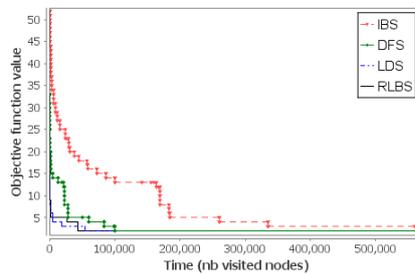


Fig. 6. Solution objective function value according to computation time of RLBS, LDS, IBS and DFS for a toy problem

Finally, we also tried using IBS. It performs adaptive variable/value selection, which prevents us from using our specific variable/value selection heuristic. We were not able to find good solutions in reasonable computation time (over 150 hours using Choco v2.1.5). Therefore, we generated really small toy problems (Fig. 6) in order to compare RLBS, DFS, LDS and IBS. IBS showed the worst

result, presumably because it cannot make use of the specific branching strategy known to be really efficient for this problem. DFS was also outperformed by LDS, as it was already reported in the literature for this problem.

5 Conclusion

We proposed a simple learning mechanism based on reinforcement learning which allows a solver to dynamically learn how to backtrack. It was evaluated for a difficult industrial planning and scheduling problem which is only efficiently solved when using specific branching heuristics. The proposed adaptive strategy greatly improved the performance in comparison with a LDS policy. This is made possible as the mechanism allows identifying which nodes are the most profitable to backtrack to and, thus, focusing on them first.

Using real industrial data showed the value of this approach. However, there are still open questions regarding how the algorithm should perform with problems for which we do not know good branching heuristics. In this situation, is it worth trying to identify which node we should backtrack to?

The combination of the adaptive backtracking strategy and adaptive branching strategies would be another interesting research opportunity.

Acknowledgments

This work has been supported by the FORAC Research Consortium industrial partners.

References

1. Ethem Alpaydin. *Introduction to machine learning*. MIT press, 2004.
2. Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
3. J Christopher Beck and Laurent Perron. Discrepancy-bounded depth first search. In *Proceedings of the Second International Workshop on Integration of AI and OR Technologies for Combinatorial Optimization Problems (CPAIOR), Germany, Paderborn*, pages 7–17, 2000.
4. Frédéric Boussemart, Fred Hemery, Christophe Lecoutre, and Lakhdar Sais. Boosting systematic search by weighting constraints. In *ECAI*, volume 16, page 146, 2004.
5. Justin A Boyan and Andrew W Moore. Using prediction to improve combinatorial optimization search. In *Sixth International Workshop on Artificial Intelligence and Statistics*, 1997.
6. Jonathan Gaudreault, Pascal Forget, Jean-Marc Frayret, Alain Rousseau, Sebastien Lemieux, and Sophie D’Amours. Distributed operations planning in the softwood lumber supply chain: models and coordination. *International Journal of Industrial Engineering: Theory Applications and Practice*, 17:168–189, 2010.
7. Jonathan Gaudreault, Gilles Pesant, Jean-Marc Frayret, and Sophie D’Amours. Supply chain coordination using an adaptive distributed search strategy. *Journal of Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 42(6):1424–1438, 2012.

8. Diarmuid Grimes and Richard J Wallace. Learning from failure in constraint satisfaction search. In *Learning for Search: Papers from the 2006 AAAI Workshop*, pages 24–31, 2006.
9. William D Harvey and Matthew L Ginsberg. Limited discrepancy search. In *Proceedings of International Joint Conference on Artificial Intelligence (1)*, pages 607–615, 1995.
10. Wafa Karoui, Marie-José Huguet, Pierre Lopez, and Wady Naanaa. YIELDS: A yet improved limited discrepancy search for CSPs. In *Proceedings of Integration of AI and OR Techniques in Constraint Programming for Combinatorial Optimization Problems*, pages 99–111. Springer, 2007.
11. Manuel Loth, Michele Sebag, Youssef Hamadi, and Marc Schoenauer. Bandit-based search for constraint programming. In *Principles and Practice of Constraint Programming*, pages 464–480. Springer, 2013.
12. Manuel Loth, Michele Sebag, Youssef Hamadi, Marc Schoenauer, and Christian Schulte. Hybridizing constraint programming and monte-carlo tree search: Application to the job shop problem. In *Learning and Intelligent Optimization*, pages 315–320. Springer, 2013.
13. Donald W Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *Journal of the Society for Industrial & Applied Mathematics*, 11(2):431–441, 1963.
14. Victor V Miagkikh and William F Punch III. Global search in combinatorial optimization using reinforcement learning algorithms. In *Proceedings of Evolutionary Computation, 1999. CEC 99. Proceedings of the 1999 Congress on*, volume 1. IEEE, 1999.
15. Thierry Moisan, Jonathan Gaudreault, and Claude-Guy Quimper. Parallel discrepancy-based search. In *Principles and Practice of Constraint Programming*, pages 30–46. Springer, 2013.
16. Robert Moll, Andrew G Barto, Theodore J Perkins, and Richard S Sutton. Learning instance-independent value functions to enhance local search. In *Advances in Neural Information Processing Systems*. Citeseer, 1998.
17. Philippe Refalo. Impact-based search strategies for constraint programming. In *Proceedings of Principles and Practice of Constraint Programming–CP 2004*, pages 557–571. Springer, 2004.
18. Wheeler Ruml. *Adaptive tree search*. PhD thesis, Citeseer, 2002.
19. Wheeler Ruml. Heuristic search in bounded-depth trees: Best-leaf-first search. In *Working Notes of the AAAI-02 Workshop on Probabilistic Approaches in Search*, 2002.
20. Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
21. Yuehua Xu, David Stern, and Horst Samulowitz. Learning adaptation to solve constraint satisfaction problems. In *Proceedings of Learning and Intelligent Optimization (LION)*, 2009.
22. Wei Zhang and Thomas G Dietterich. A reinforcement learning approach to job-shop scheduling. In *Proceedings of International Joint Conferences on Artificial Intelligence*, volume 95, pages 1114–1120. Citeseer, 1995.