

Choosing Search Heuristics by Non-Stationary Reinforcement Learning

Alexander Nareyek

GMD FIRST, Kekuléstr. 7, D - 12489 Berlin, Germany

November 29, 2001

Abstract. Search decisions are often made using heuristic methods as real-world applications can rarely be tackled without any heuristics. In many cases, multiple heuristics can potentially be chosen, and it is not clear a priori which would perform best. In this article, we propose a procedure that learns, during the search process, how to select promising heuristics. The learning is based on weight adaptation and can even switch between different heuristics during search. Different variants of the approach are evaluated within a constraint programming environment.

Keywords: Non-Stationary Reinforcement Learning, Optimization, Local Search, Constraint Programming

1. Introduction

All kinds of search techniques include choice points at which decisions must be made between various alternatives. For example, in refinement search, an extension step from a partial solution toward a complete solution must be chosen. In local search methods, it must be decided how a complete but suboptimal/infeasible solution is to be changed toward an optimal/feasible solution.

However, for large and complex real-world problems, decisions can rarely be made in an optimal way. Especially for local search techniques, this is a very critical issue because they do not normally incorporate backtracking mechanisms. Many different meta-heuristic techniques have therefore been developed to handle the complications involved when choosing an alternative.

Figure 1 shows a choice point, representing the current state/solution of local search, and multiple alternatives, representing the so-called *neighbor states* that can be reached within an iteration.

Nearly all local search methods evaluate all neighbor states in a kind of look-ahead step in order to choose the most beneficial alternative. However, complex real-world problems – such as action planning including time, resources and optimization – often have utility functions whose computation requires a great deal of computing power. Analyzing large neighborhoods is mostly out of the question, and even smaller neighborhoods are difficult to check. Techniques like simulated



© 2001 Kluwer Academic Publishers. Printed in the Netherlands.

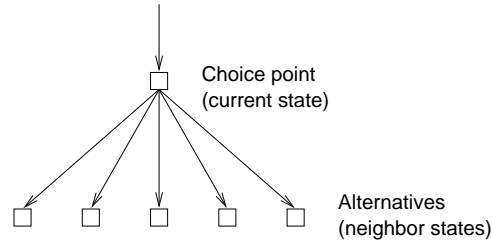


Figure 1. A decision point in local search

annealing (Kirkpatrick et al., 1983) are highly suitable for these tasks because *only one neighbor is analyzed* for a choice decision (though more neighbors may be analyzed if the current neighbor appears to be unsuitable).

For the purposes of this paper, we go one step further, not analyzing any neighbor, but choosing a neighbor according to learned utility values. In addition, we do not choose a specific neighbor state but a *transformation heuristic* that will be applied to create the new state. Unlike other reinforcement learning approaches for learning which heuristics perform well, our approach allows the search to switch between different heuristics during search in order to adapt to specific regions of the search space.

Section 2 introduces the constraint programming environment that is applied in our experiments, and details the use of heuristics. Weights and their adaptation are presented in Sec. 3. The scheme is evaluated in Sec. 4. Conclusions and related work are discussed in Sec. 5.

2. Search Decisions

As an example of local search, we give a brief description of the search method applied in the **DragonBreath** engine. The underlying paradigm is presented in detail in (Nareyek, 2001 (a)).

The problem is specified as a so-called constraint satisfaction problem (CSP). A CSP consists of

- a set of variables $x = \{x_1, \dots, x_n\}$
- where each variable is associated with a domain d_1, \dots, d_n
- and a set of constraints $c = \{c_1, \dots, c_m\}$ over these variables.

The domains can be symbols as well as numbers, continuous or discrete (e.g., “door”, “13”, “6.5”). Constraints are relations between

variables (e.g., “ x_a is a friend of x_b ”, “ $x_a < x_b \times x_c$ ”) that restrict the possible value assignments. Constraint satisfaction is the search for a variable assignment that satisfies the given constraints. Constraint optimization requires an additional function that assigns a quality value to a solution and tries to find a solution that maximizes this value.

In our local search approach, a specific cost function is specified for every constraint (so-called *global constraints*), which returns a value that represents the constraint’s current inconsistency/optimality with respect to the connected variables. For example, a simple **Sum** constraint with two variables a and b to be added and an s variable for the sum could specify its costs as $\text{Sum}_{costs} = |a + b - s|$.

In addition, a constraint has a number of heuristics to improve its cost function. For example, a heuristic for the **Sum** constraint could randomly choose one of the related variables and change it such that there are no more costs. Another heuristic might resolve the inconsistency by distributing the necessary change such that all variables are changed by the same (minimal) amount. The constraint must make the choice as to which heuristic to apply on its own.

On top of all constraints is a *global search control* which selects, in each iteration of local search, one of the constraints which is to perform a change, i.e., the transition to a neighbor state. Figure 2 shows the control flow.

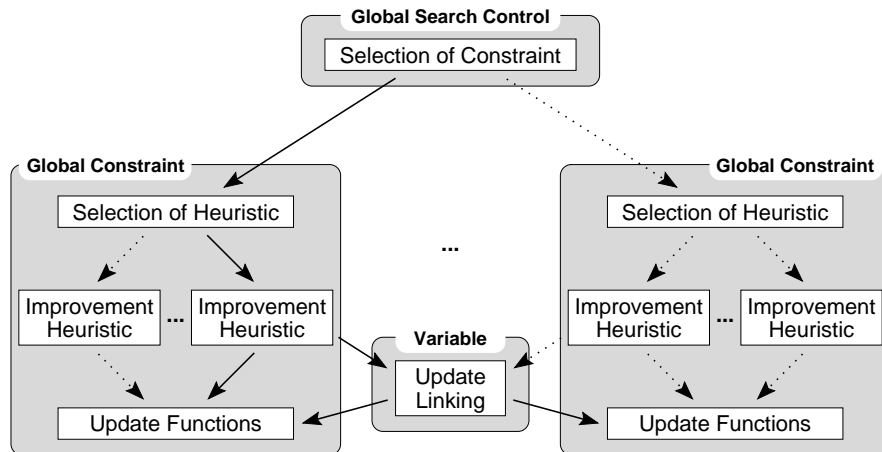


Figure 2. Using global constraints for local search

The global search control possesses qualitative and quantitative information from the constraints’ cost functions to decide which constraint to choose (e.g., the constraint with the maximal costs), but a constraint itself has little guidance as to which of its heuristics to

choose. This choice point — for choosing one of the constraint’s heuristics — is investigated below.

3. Utility Weight

For a choice point, a *utility value* $\omega_a \geq 1$ is computed/maintained for every alternative a (an *alternative* stands for a *heuristic* here) that expresses the expected benefit of choosing this alternative.

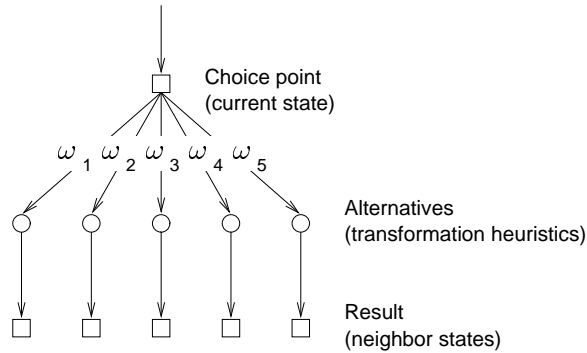


Figure 3. A decision point in our approach

The utility values are subject to learning schemes, which change the values based on past experiences with choosing this alternative. In many cases, an appropriate balance of the utility values will depend on the area of the search space that the search is currently in. We will therefore focus on schemes that dynamically adapt during search and not only after a complete run.

3.1. SELECTION FUNCTION

Selection between possible alternatives is done based on the alternative’s utility values. Here, we look at two simple ones. The first is a *fair random choice* (a so-called *softmax* kind of selection rule), referred to below as M:O, which selects an alternative a from the choice point’s alternatives \mathcal{A} with a choice probability p_a in proportion to the alternative’s utility value ω_a :

$$\text{M:O} : \quad p_a = \frac{\omega_a}{\sum_{i \in \mathcal{A}} \omega_i}$$

Another possibility is to make a random choice between the alternatives with *maximal* utility values:

$$\text{M:1} : p_a = \frac{\begin{cases} 0 & : \exists i \in \mathcal{A} : \omega_a < \omega_i \\ 1 & : \forall i \in \mathcal{A} : \omega_a \geq \omega_i \end{cases}}{\sum_{i \in \mathcal{A} | \forall j \in \mathcal{A} : \omega_i \geq \omega_j} 1}$$

3.2. WEIGHT ADAPTATION

All utility weights have integer domains and are initially set to 1. If the choice point is selected, the utility weight of the alternative is changed that was chosen by the choice point when it was called last time. The kind of change depends on the relation of the current objective function value o_{now} to the objective function value when the choice point was called last time o_{before} (i.e., if there is a positive or negative reinforcement). The update schemes below can be combined to give many different strategies, e.g., a simple P:1-N:1 strategy.

o_{now} **better-than** o_{before} : (positive reinforcement)

P:1 (Additive Adaptation): $\omega_a \leftarrow \omega_a + 1$

P:2 (Escalating Additive Adaptation): $\omega_a \leftarrow \omega_a + m_{promotion}$

P:3 (Multiplicative Adaptation): $\omega_a \leftarrow \omega_a \times 2$

P:4 (Escalating Multiplicative Adaptation): $\omega_a \leftarrow \omega_a \times m_{promotion}$

P:5 (Power Adaptation): $\omega_a \leftarrow \begin{cases} \omega_a \times \omega_a & : \omega_a > 1 \\ 2 & : \omega_a = 1 \end{cases}$

o_{now} **worse-than-or-equal-to** o_{before} : (negative reinforcement)

N:1 (Subtractive Adaptation): $\omega_a \leftarrow \omega_a - 1$

N:2 (Escalating Subtractive Adaptation): $\omega_a \leftarrow \omega_a - m_{demotion}$

N:3 (Divisional Adaptation): $\omega_a \leftarrow \frac{\omega_a}{2}$

N:4 (Escalating Divisional Adaptation): $\omega_a \leftarrow \frac{\omega_a}{m_{demotion}}$

N:5 (Root Adaptation): $\omega_a \leftarrow \sqrt{\omega_a}$

If a utility value falls below 1, it is reset to 1; if a utility value exceeds a certain max_ω , it is reset to max_ω ; if a utility value is assigned a non-integer value, it is rounded down. In the case of an escalating adaptation, each time there is a consecutive improvement/deterioration, the

$m_{promotion}/m_{demotion}$ value is doubled. Otherwise, it is reset to 1 (for P:2 and N:2) or 2 (for P:4 and N:4).

3.3. INVALID ALTERNATIVES

For some choice points, more than one alternative must be tested. For example, an alternative may turn out to be infeasible. An *applicability flag* f with a value of 0 or 1 is introduced for every alternative, indicating whether the alternative is still a valid option:

$$\omega_a \leftarrow f_a \times \omega_a$$

By the option of setting an applicability flag to 0, alternatives can often be ruled out a priori by simple feasibility tests.

However, in some cases, the infeasibility of an alternative will only become apparent during the state-transformation process of the chosen heuristic, i.e., after the choice has been made. In such a case, all changes in the current state that were made after the choice point are reversed, the corresponding applicability flag is set to 0 and the choice process is repeated. If no alternative remains applicable, the constraint improvement fails.

If the choice point's selection is subject to the learning scheme, applicability flags are not set to 0 if an alternative fails. The failure may be caused by a bad random decision during the alternative's computations and the alternative may not be *fully* inapplicable. The learning process can handle this situation more appropriately than in a non-learning case, skipping the usual update of the utility weights and *temporarily* dividing the failed alternative's utility weight by two (though no weight may fall below one). If one of the alternatives has been successfully applied, all adaptations of the utility weights that were done for the restarts are undone.

4. Empirical Evaluation

Two optimization problems — the Orc Quest problem and the Logistics Domain — are evaluated with different learning/selection schemes. The concrete problems and solving heuristics are not explained here because they are not relevant to the techniques applied. A detailed presentation of the problems can be found in (Nareyek, 2001 (b)). The Orc Quest problem's solving process involves only three constraints with six heuristics each. For all of these, the learning scheme is applied. The Logistics Domain's solving process includes a much greater and varying

number of constraints. The learning scheme is applied to all constraint of a particular type¹, which includes five alternative heuristics.

4.1. RESULTS

A strategy is denoted by P-N-M, $P \in \{1..5\}$ indicating the adaptation scheme that is applied in the case of an improvement, $N \in \{1..5\}$ the adaptation scheme for non-improvement, and $M \in \{0, 1\}$ if the fair random choice is applied or a maximal value is chosen.

The results for some strategies for the Orc Quest problem are shown in Fig. 4 as the percentage of test runs (100 % = 100,000 test runs) that found the optimal solution after a specific number of iterations. The problem from the Logistics Domain is much harder, so only the best solution (minimal duration) found after 100,000 iterations is shown in Fig. 5 (100 % test runs = 1,000 test runs).

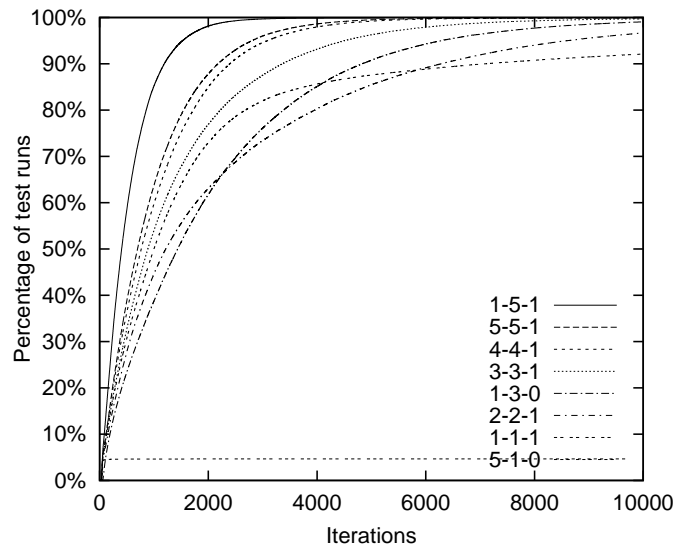


Figure 4. Sample strategies for the Orc Quest problem

A more detailed analysis is given in Figs. 6 and 7. For specific percentages, it is shown after how many iterations this percentage of test runs found the optimum (for Fig. 6), and the lowest duration that was found by this percentage of test runs after 100,000 iterations (for Fig. 7) The strategies are sorted according to which strategy resulted

¹ For the Logistics Domain, the duration minimization of Problem 6-1a is analyzed, applying the weight adaptation for the STATE RESOURCE CONSTRAINT's selection of an improvement heuristic.

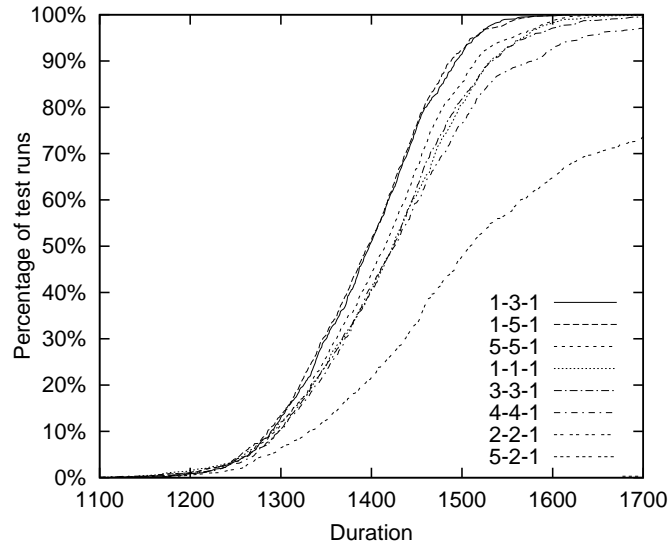


Figure 5. Sample strategies for the Logistics Domain

in the least iterations/duration for a maximal percentage of test runs (not considering the 100 % rate).

The general trend is that a low (e.g., additive/P:1) rate of adaptation is good in the case of an improvement, a strong (e.g., root/N:5) rate of adaptation is good in the case of a deterioration, and a choice of a maximal weight is often better than a fair random choice. The explorative feature of the fair random choice may not be that important because there are very often cases of negative reinforcement that quickly change the weight situation.

Because the Orc Quest problem involves only three constraints, we can easily visualize some further properties of the search process for it. One interesting property is the ratio of positive to negative reinforcements shown in Fig. 8. However, it is not very surprising that this ratio deteriorates according to the strategy ordering shown in Fig. 6.

Figure 9 shows how many times a constraint's highest weight changes, i.e., how many times a weight is assigned a value above a certain percentage of the total values of the choice point's weights, and the last time this percentage was reached, it was reached by another weight. Strategies that perform many changes in the configuration appear to perform better. This might be an indication of why strategies with a low rate of adaptation in the case of an improvement and a strong rate of adaptation in the case of a deterioration are likely to perform better, because such strategies facilitate a reconfiguration of the weight situation.

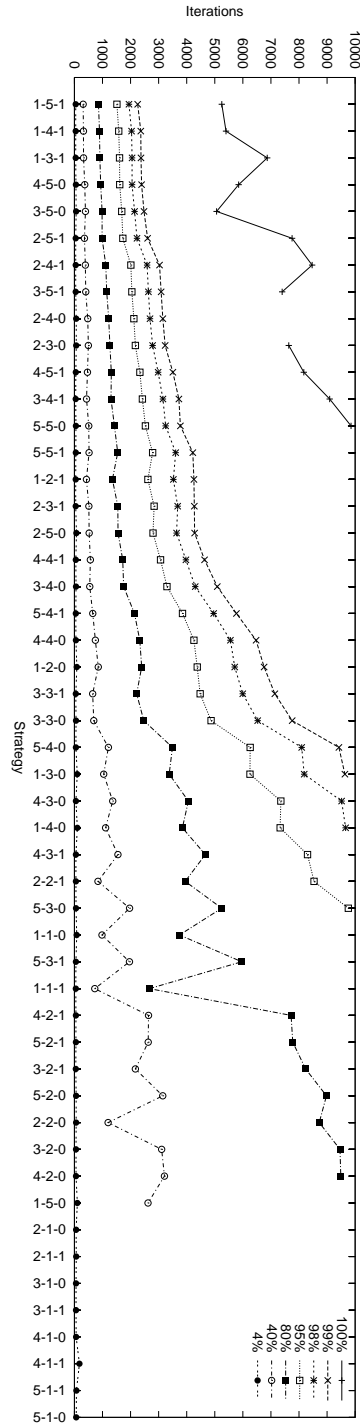


Figure 6. Weight-adaptation results for the Orc Quest problem

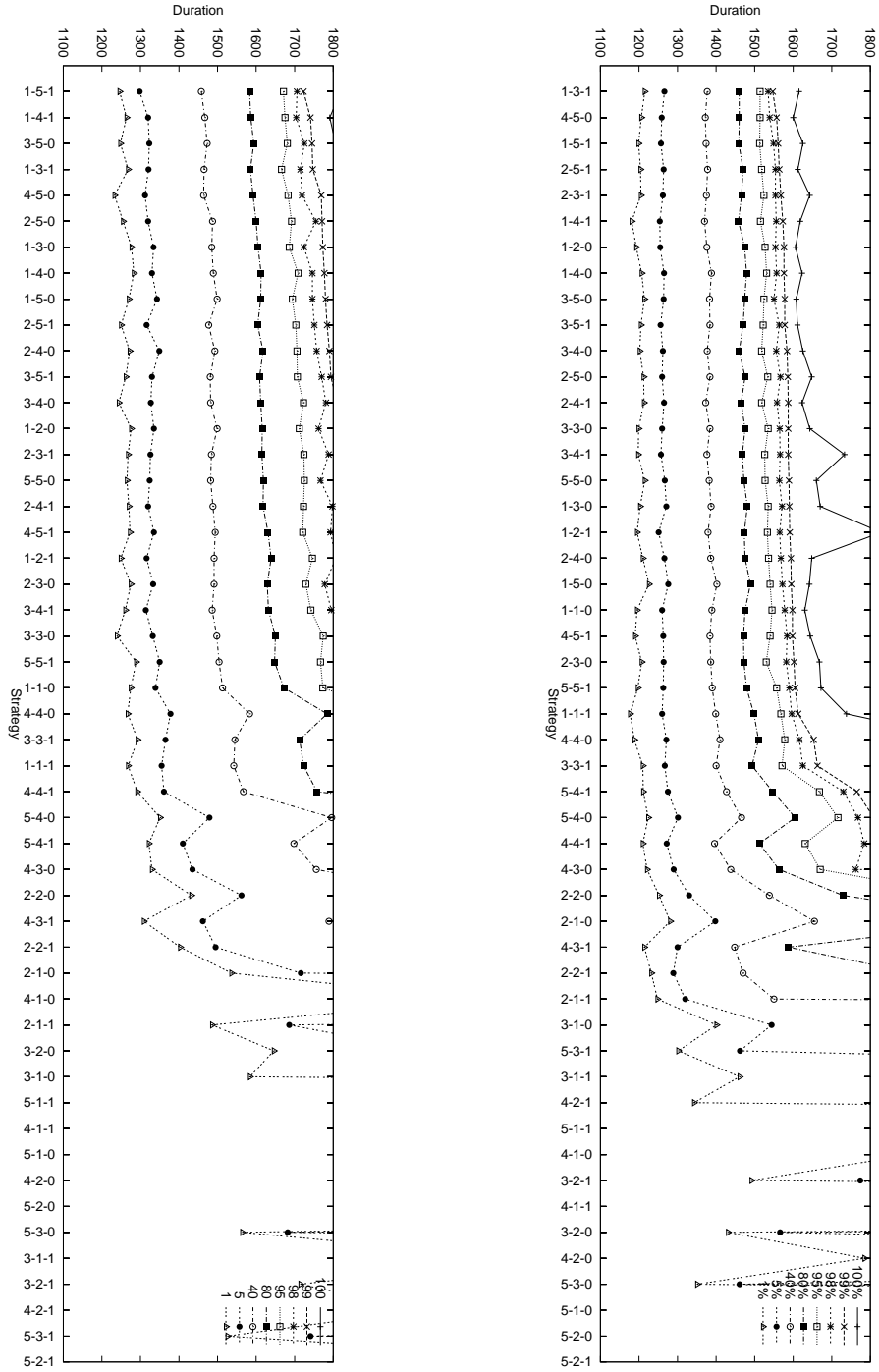


Figure 7. Weight-adaptation results for the Logistics Domain after 25,000 (left) and 100,000 (right) iterations

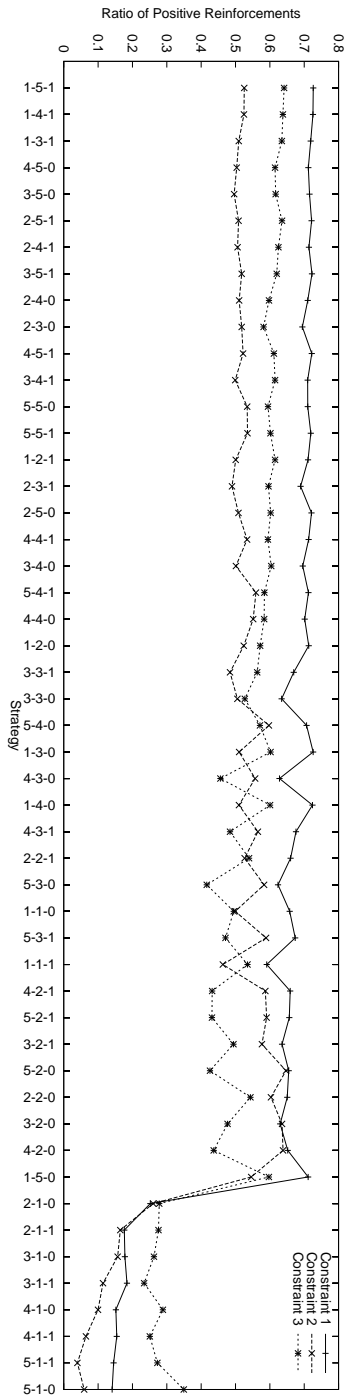


Figure 8. Ratio of positive reinforcements for the Orc Quest problem

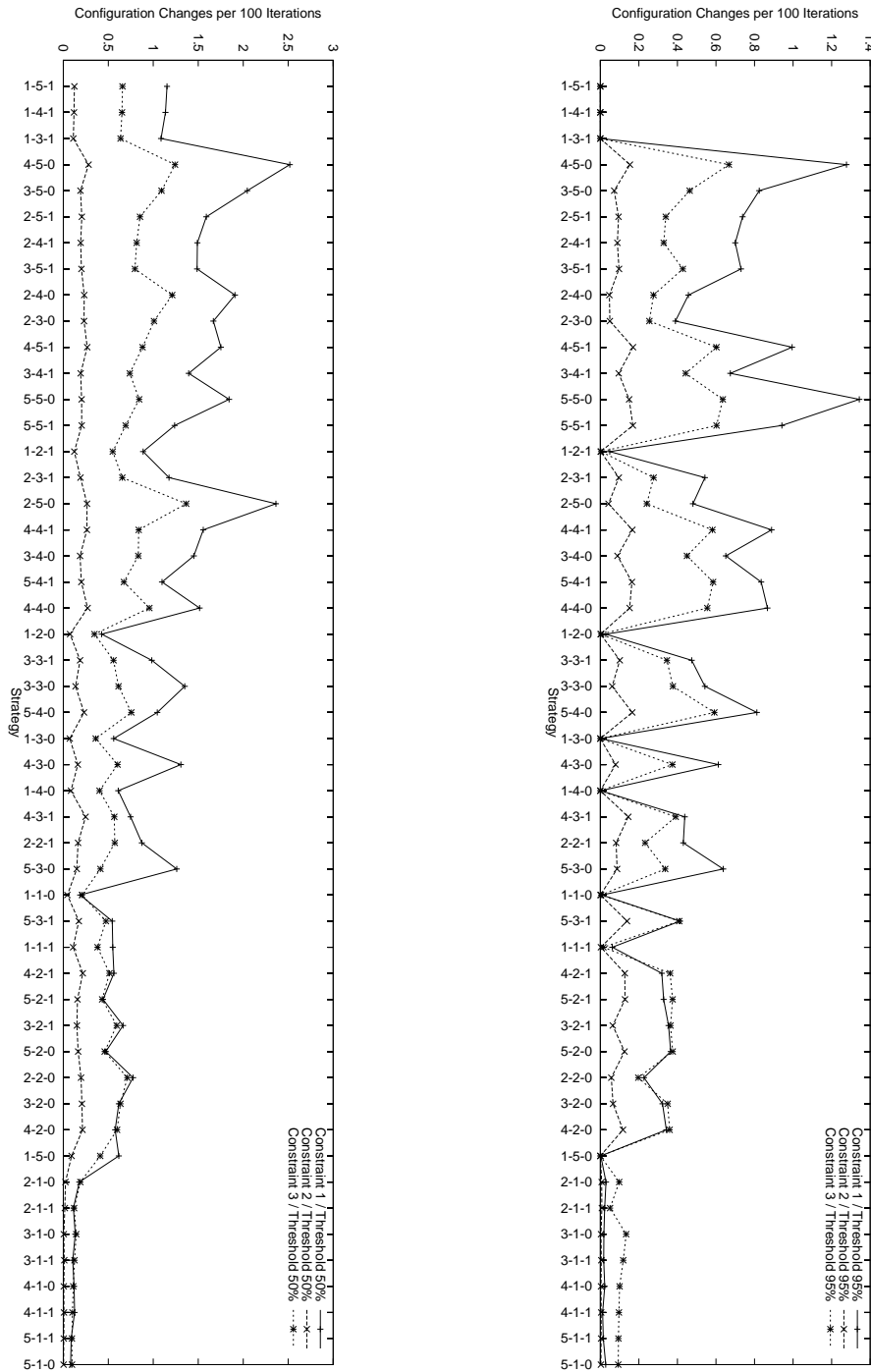


Figure 9. Number of configuration changes for the Orc Quest problem

An exception here are strategies with a very low positive reaction (P:1). Because of the slow growth of the weights, a smaller number of clear reconfigurations are performed. However, this simmering situation would appear to have its advantages as well.

Figure 10 shows how many times the highest weight is re-established, i.e., how many times a weight is assigned a value above a certain percentage of the total values of the choice point's weights, and the last time this percentage was reached, it was reached by the *same* weight. In general, one would expect strategies that re-establish old configurations to do needless work and thus, possibly perform worse. However, the figures do not show many differences here. The reason for this is probably that the better strategies perform a lot of configuration changes in general, and, are thus also more likely to re-establish configurations more often. But, of course, the ratio of re-established configurations to all reconfigurations is much better here.

4.2. EXTENDED EXPERIMENTS

Following the observed trend, we can extend our experiments by more extreme options in this direction:

P:0 (No Adaptation): $\omega_a \leftarrow \omega_a$

To enable negative adaptations for this option, in the case of a negative change the decrease of ω_a is distributed as an increase to all $\omega_{i \neq a}$ (starting with high initial weights).

N:6 (Total Loss Adaptation): $\omega_a \leftarrow 1$

Figure 11 shows that these options do not improve performance for the Orc Quest problem. However, as shown in Fig. 12, strategies with an N:6 option appear to work well for the early phase of search, i.e., for less constrained problems.

4.3. STATIONARY REINFORCEMENT LEARNING

So far, we have looked at different methods to adapt the weights during search, assuming that different areas of the search space can be handled more efficiently using different search strategies. Although this assumption seems to be intuitively correct, it remains to be shown to be true. This section, then, compares adaptive non-stationary learning with stationary approaches.

Previous approaches adapted learning parameters *after a complete run* or *when a local minimum was reached*. Of these two options, only an adaptation after a complete run (with an upper bound of a specific

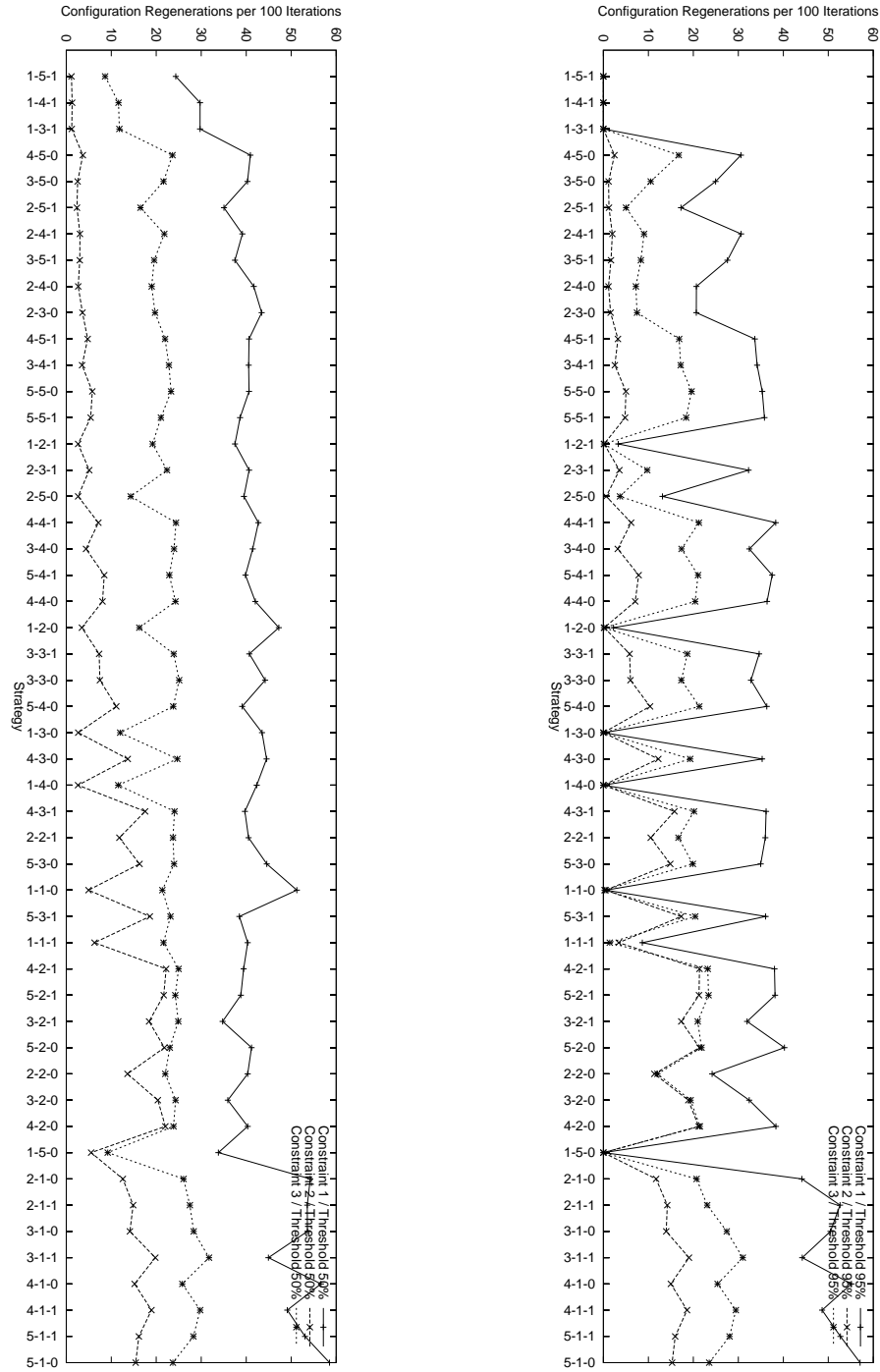


Figure 10. Number of configuration regenerations for the Orc Quest problem

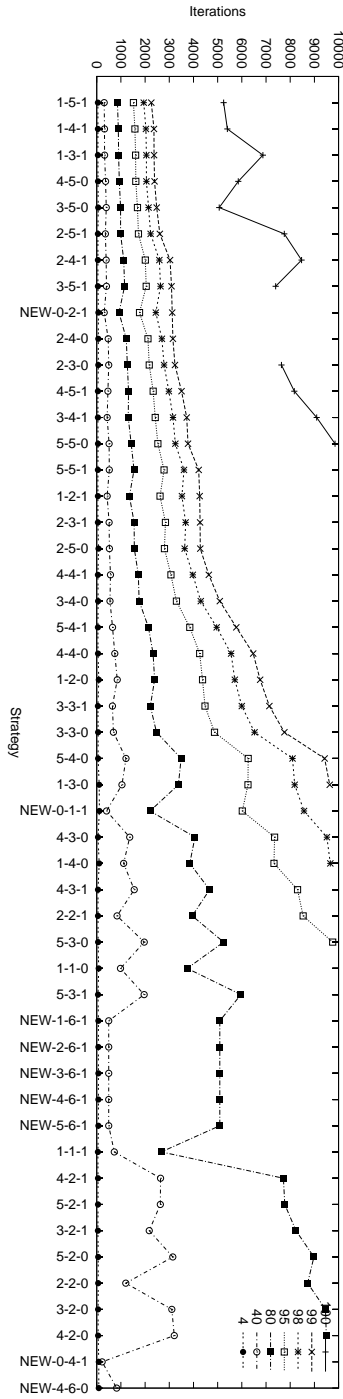


Figure 11. Extended weight-adaptation results for the Orc Quest problem; showing only the 50 best strategies

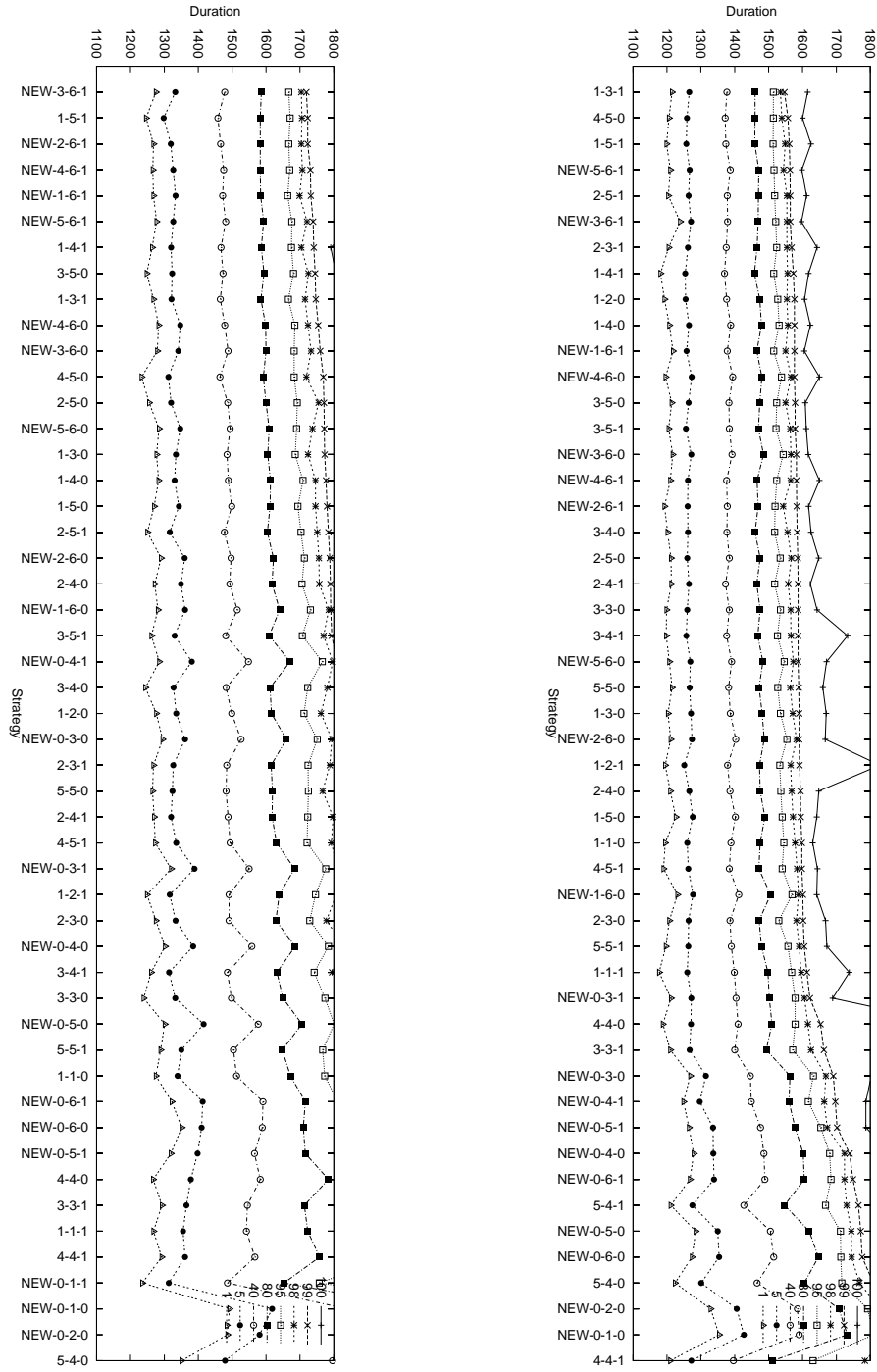


Figure 12. Extended weight-adaptation results for the Logistics Domain after 25,000 (left) and 100,000 (right) iterations; showing only the 50 best strategies

number of iterations) is applicable here because we do not evaluate the whole neighborhood and cannot therefore tell if we are in a local minimum. The time taken by the learning process to find an optimal stationary weight distribution is not measured because the results may be very different for different learning techniques. Thus, the adaptive learning strategies are compared with the optimal stationary distribution.

For the Orc Quest problem, we can actually find an optimal static weight distribution such that all test runs find the optimum in about **35** iterations. With the most simple, non-stationary 1-1-0 strategy, 50% of the test runs found the optimum after **1,305** iterations, and after **405** iterations for strategy 1-5-1. Thus, an adaptive strategy would seem to perform very poorly for the simple Orc Quest problem. This is not completely true, however, given the time that would be required to learn the optimal stationary distribution. For example, using a simple static distribution such that every heuristic is chosen equally often, **none** of the 100,000 test runs found the optimum within 100,000 iterations. We conclude that, if the problem (or very “similar” problems) is solved very often, a stationary reinforcement learning approach will ultimately perform much better; but for a short time-frame, the non-stationary approach is probably much superior.

For the more complex Logistics Domain, our findings are different. The performance of even the most simple, non-stationary 1-1-0 strategy is similar to that a carefully hand-tailored static weight distribution, i.e., a static distribution does not work well even disregarding the learning time (see Fig. 13). Our assumption that it is useful to switch between different heuristics during search in order to adapt to specific regions of the search space proves valid for this more complicated problem.

5. Conclusion

The use of a neighborhood of repair heuristics is a promising way to implement a local search — especially for complex real-world problems in which the computation of a state’s objective function value is often very costly. Using the repair heuristics, domain-dependent knowledge to guide the search can easily be incorporated into the search process. The approach used here is based on (Nareyek, 2001 (a)). Similar techniques were applied in (Rabideau et al., 1999; Smith, 1994; Zweben et al., 1994).

However, finding an appropriate strategy that guides when to apply which heuristics is not easy. This article has presented an approach to learn a selection strategy by modifying weights. Other approaches

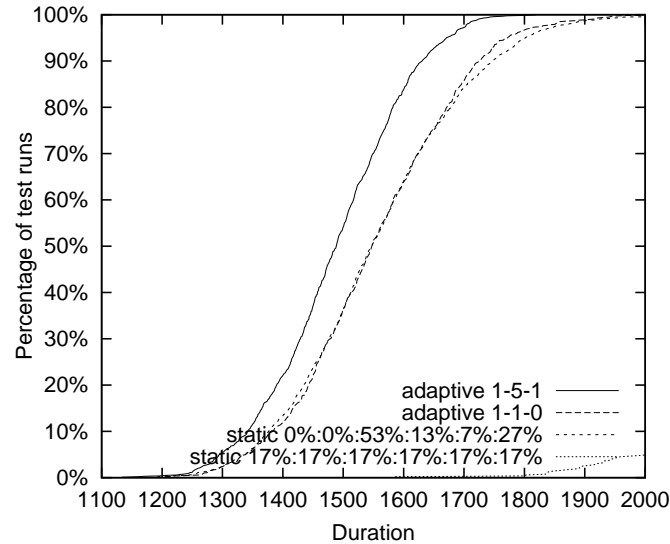


Figure 13. Adaptive vs. static strategies for the Logistics Domain (after 25,000 iterations), disregarding the learning time for the static distribution

that make use of weights for re-structuring the neighborhood include (Boyan and Moore, 1997; Frank, 1997; Schuurmans and Southey, 2000; Voudouris and Tsang, 1995). Unlike these approaches, we do not only change the weights when a local optimum or solution is reached. Search usually undergoes different phases (i.e., search-space regions) in which different search heuristics work best. Thus, even during search, the heuristics' configuration is constantly updated. For this purpose, a less-carrot-more-stick strategy seems to be appropriate, allowing for quick configuration changes and preventing the old configuration from being re-established too quickly.

In reinforcement learning, non-stationary environments (such as the search-space region) are only rarely considered. Examples include approaches based on supervised techniques (Schmidhuber, 1990), evolutionary learning (Littman and Ackley, 1991) and model-based learning (Michaud and Mataric, 1998). Unlike these approaches, we have used a modification of standard action-value methods (Sutton and Barto, 1998), applying functional updates instead of cumulative value additions in order to influence the impact of the already learned reinforcements. This simple method enables the search to compute weight updates very quickly – which is very important for a local search environment because a single iteration should consume only very little computing power.

Adaptive weights are not restricted to local search; they can also be used for (esp. restart-based) refinement search. Examples include the pheromone trails in ant colony optimization (Dorigo et al., 1999), the use of domain-specific prioritizers (Joslin and Clements, 1999) and action costs in adaptive probing (Ruml, 2001). The results obtained in this study may be transferred to these areas, and techniques like pheromone evaporation are worth studying for neighborhoods of heuristics as well.

So far, we have not considered quantitative cost-function effects of decisions. Improvement or non-improvement was the only criterion for learning. However, “good” heuristics may not be equally good on the quantitative level and incorporating mechanisms to exploit the quantitative differences is a promising idea for future work.

The presented techniques are integrated into the **DragonBreath** engine, which is a free optimization engine based on constraint programming and local search. It can be obtained via:

<http://www.ai-center.com/projects/dragonbreath/>

Acknowledgements

Thanks to Michael Littman for his feedback.

References

- Boyan, J. A., and Moore, A. W. Using Prediction to Improve Combinatorial Optimization Search. In Proceedings of the Sixth International Workshop on Artificial Intelligence and Statistics (AISTATS-97), 1997.
- Dorigo, M.; Di Caro, G.; and Gambardella, L. M. Ant Algorithms for Discrete Optimization. *Artificial Life* 5(3): 137–172, 1999.
- Frank, J. Learning Short-Term Weights for GSAT. In Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence (IJCAI-97), 384–391, 1997.
- Joslin, D. E., and Clements, D. P. Squeaky Wheel Optimization. *Journal of Artificial Intelligence Research* 10: 353–373, 1999.
- Kirkpatrick, S.; Gelatt, C. D.; and Vecchi, M. P. Optimization by Simulated Annealing. *Science* 220(4598): 671–680, 1983.
- Littman, M. L., and Ackley, D. H. Adaptation in constant utility non-stationary environments. In Proceedings of the Fourth International Conference on Genetic Algorithms, 136–142, 1991.
- Michaud, F., and Matarić, M. J. Learning from History for Behavior-Based Mobile Robots in Non-Stationary Environments. *Machine Learning* 31, Joint Special Issue on Learning in Autonomous Robots, 141–167, 1998.
- Nareyek, A. (a) Using Global Constraints for Local Search. In Freuder, E. C., and Wallace, R. J. (eds.), *Constraint Programming and Large Scale Discrete*

- Optimization*, American Mathematical Society Publications, DIMACS Volume 57, 9–28, 2001.
- Nareyek, A. (b) *Constraint-Based Agents – An Architecture for Constraint-Based Modeling and Local-Search-Based Reasoning for Planning and Scheduling in Open and Dynamic Worlds*. Reading, Springer LNAI 2062, 2001.
- Rabideau, G.; Knight, R.; Chien, S.; Fukunaga, A.; and Govindjee, A. Iterative Repair Planning for Spacecraft Operations in the ASPEN System. International Symposium on Artificial Intelligence Robotics and Automation in Space (iSAIRAS 99), 1999.
- Ruml, W. Incomplete Tree Search using Adaptive Probing. In Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence (IJCAI-01), 235–241, 2001.
- Schmidhuber, J. Making the World Differentiable: On Using Self-Supervised Fully Recurrent Neural Networks for Dynamic Reinforcement Learning and Planning in Non-Stationary Environments. Technical Report, TR FKI-126-90, Department of Computer Science, Technical University of Munich, 1990.
- Schuermans, D., and Southey, F. Local search characteristics of incomplete SAT procedures. In Proceedings of the Seventeenth National Conference on Artificial Intelligence (AAAI-2000), 297-302, 2000.
- Smith, S. F. OPIS: A Methodology and Architecture for Reactive Scheduling. In Zweben, M., and Fox, M. S. (eds.), *Intelligent Scheduling*, Morgan Kaufmann, 29–66, 1994.
- Sutton, R. S., and Barto, A. G. *Reinforcement Learning: An Introduction*. Reading, MIT Press, 1998.
- Voudouris, C., and Tsang, E. Guided Local Search. Technical Report CSM-247, University of Essex, Department of Computer Science, Colchester, United Kingdom, 1995.
- Zweben, M.; Daun, B.; Davis, E.; and Deale, M. Scheduling and Rescheduling with Iterative Repair. In Zweben, M., and Fox, M. S. (eds.), *Intelligent Scheduling*, Morgan Kaufmann, 241–255, 1994.