

University of Groningen

## Model-Based Multi-Objective Reinforcement Learning

Wiering, Marco; Withagen, Maikel; Drugan, Madalina

*Published in:*

IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*

Final author's version (accepted by publisher, after peer review)

*Publication date:*

2014

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

Wiering, M., Withagen, M., & Drugan, M. (2014). Model-Based Multi-Objective Reinforcement Learning. In IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning Orlando: IEEE (The Institute of Electrical and Electronics Engineers).

### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

# Model-Based Multi-Objective Reinforcement Learning

Marco A. Wiering (*IEEE Member*)

Institute of Artificial Intelligence, University of Groningen, The Netherlands, Email: m.a.wiering@rug.nl

Maikel Withagen

Institute of Artificial Intelligence, University of Groningen, The Netherlands, Email: maikel.withagen@gmail.com

Mădălina M. Drugan

Artificial Intelligence Lab, Vrije Universiteit Brussel, Belgium, Email: mdrugan@vub.ac.be

**Abstract**—This paper describes a novel multi-objective reinforcement learning algorithm. The proposed algorithm first learns a model of the multi-objective sequential decision making problem, after which this learned model is used by a multi-objective dynamic programming method to compute Pareto optimal policies. The advantage of this model-based multi-objective reinforcement learning method is that once an accurate model has been estimated from the experiences of an agent in some environment, the dynamic programming method will compute all Pareto optimal policies. Therefore it is important that the agent explores the environment in an intelligent way by using a good exploration strategy. In this paper we have supplied the agent with two different exploration strategies and compare their effectiveness in estimating accurate models within a reasonable amount of time. The experimental results show that our method with the best exploration strategy is able to quickly learn all Pareto optimal policies for the Deep Sea Treasure problem.

## I. INTRODUCTION

Reinforcement learning (RL) [1], [2] enables an autonomous agent to learn from its interactions with a particular environment that emits reward signals to the agent. The objective of the agent is to learn a policy that obtains the highest possible discounted cumulative reward intake. In this paper we consider value-based reinforcement learning, where the agent estimates a value function denoting the future reward intake and uses this value function to select actions. Many value-based reinforcement learning algorithms have been proposed [2]. These algorithms can be divided into model-free and model-based reinforcement learning algorithms. Model-free methods such as Q-learning [3] update the Q-value function after each interaction with the environment without estimating a model. Model-based RL methods first learn to estimate a model of the environment and then use a dynamic programming algorithm to compute the policy. The advantage of model-based RL methods is that experiences of the agent are used more effectively, leading to faster convergence to optimal policies.

Although traditionally reinforcement learning algorithms have been applied solely to single objective decision problems, during the last decade the amount of research on multi-objective problems has considerably increased [4], [5]. In multi-objective reinforcement learning (MORL), the reward function emits a reward vector instead of a single scalar reward, and the goal is to learn all Pareto optimal policies.

In this paper we describe a novel model-based reinforcement learning algorithm for solving multi-objective reinforcement learning problems. Similar to [6], our multi-objective reinforcement learning algorithm considers the Pareto dominance relation to order the policies. The Pareto dominance relation [7] considers that two policies are incomparable if one policy is better in one objective and worse in another objective than the second policy. Intuitively, a policy is better, or dominates, another policy if it is better (or equal) in all objectives, and a policy is worse than another policy if it is worse (or equal) in all objectives. The set of Pareto optimal policies is the set of policies for which there is no policy that dominates it.

**Related Work.** Most state-of-the-art multi-objective reinforcement learning algorithms are model-free value-based reinforcement learning algorithms [4], [5]. The MORL algorithms often use scalarization functions [8] to transform the reward vectors into scalar reward values in order to use standard reinforcement learning algorithms. Since a single scalarization function results in a single Pareto optimal policy, the scalarized MORL methods use multiple scalarization functions that will converge to a set of Pareto optimal policies. The scalarized MORL algorithms have difficulties in identifying the complete set of Pareto optimal policies for non-convex Pareto fronts [5], [9], but they are efficient for convex Pareto fronts [4], [10].

There are also multi-objective multi-armed bandits algorithms, which are considered reinforcement learning algorithms with a single state used to study the theoretical properties of reinforcement learning algorithms. The well known exploration/exploitation trade-off plays a very important role in this problem [11]. These methods have also used both the Pareto dominance relation [12] and the scalarization functions [13] to identify the Pareto front. Some multi-objective reinforcement learning algorithms use the lexicographical order relation [14] that assumes that one objective is more important than another objective. Finally, the hypervolume unary indicator has been used [15], [16], which is another function to transform a set of reward vectors into a single reward value.

**Novel Contributions.** Current MORL methods do not make use of model-based methods that first estimate the model of the environment and then solve this model. Model-based RL

methods for a single objective have been researched quite extensively [17], [18] and have been shown to be able to make effective use of the agent’s experiences. The transition from standard model-based RL methods to model-based MORL methods involves different dynamic programming-like methods that can solve the problem.

Furthermore, MORL is much more challenging than standard RL, since many policies have to be learned at the same time during the interaction with the environment. For some problems the amount of Pareto optimal policies may become very large, which can make the computational time huge. For this reason, we consider in this paper deterministic multi-objective sequential decision making problems, which we solve with our model-based MORL method. Our method uses the previously introduced CON-MODP algorithm [19], which is a multi-objective dynamic programming method that can solve deterministic multi-objective Markov Decision Processes (MOMDPs) [19], [20]. Instead of using CON-MODP to solve an MOMDP, in this paper a model is learned from the interaction with the environment and CON-MODP uses the estimated model to compute the Pareto optimal policies.

Although the CON-MODP method is optimal in the sense that it computes all Pareto optimal policies for deterministic MOMDPs, the model that is estimated may not be completely accurate. This is because the agent has to visit all state-action pairs at least a single time in order to learn the optimal model underlying the environment. Therefore, we have supplied the agent with two different exploration strategies which will be compared on their effectiveness in learning a model that is useful for computing all Pareto optimal policies. We developed a least-visited exploration strategy that updates how often each action in each state has been executed, and then always chooses the action in a state that has been executed the least amount of times. This exploration strategy is compared to a fully randomized exploration strategy that always selects a random action in each state. Both methods will finally learn to estimate a perfect model for a deterministic environment, but it is interesting to see which method learns it faster.

The results show that the proposed model-based MORL method with the least-visited exploration strategy is very effective in quickly learning all Pareto optimal policies for the Deep Sea Treasure problem [5]. Also the use of the random algorithm leads to finding almost all Pareto optimal policies within 2000 epochs. The Deep Sea Treasure problem is a deterministic problem involving two objectives. Furthermore, it has a discrete state space containing only 110 states. Therefore, our model-based MORL technique is very effective for this problem and performs in a near-optimal way.

**Outline.** In Section II we describe model-based reinforcement learning. Section III explains our model-based MORL method. In Section IV experimental results are presented for the Deep Sea Treasure problem. Finally, Section V concludes with our main findings and some possible future work directions.

## II. MODEL-BASED REINFORCEMENT LEARNING

Reinforcement learning algorithms are very useful to let an agent learn from its interaction with an environment. Typically, the RL agent is in some particular state, uses its policy to select an action, and after executing this action, the agent makes a transition to a next state and receives a scalar reward signal. Most often the RL agent is learning without any a-priori information about the environment, and therefore it has to learn from trial-and-error which policy can be used to obtain the highest future cumulative reward intake. In this section we will explain model-based RL methods that estimate a model of the environment while interacting with it and use a dynamic programming-like method to compute the best policy for its current estimated model.

### A. Markov Decision Processes

In reinforcement learning the environment is typically modeled as a Markov Decision Process (MDP). A finite MDP consists of the following information:

- A set of environmental states  $S$ , where  $s_t \in S$  is the state of the environment at time-step  $t$ .
- A set of actions  $A$ , where  $a_t \in A$  is the action executed by the agent at time-step  $t$ .
- A transition function  $T(s'|s, a)$  denoting the probability that the agent finds itself in each possible next state  $s'$  after executing action  $a$  in state  $s$ .
- A reward function  $R(s, a)$  denoting the expected immediate reward obtained by executing action  $a$  in state  $s$ .
- A discount factor  $\gamma$ , with  $0 \leq \gamma < 1$ , which gives more importance to immediate rewards compared to rewards obtained in the future.

The goal is to calculate a policy  $\pi(\cdot)$  mapping states to actions that will maximize the expected sum of discounted rewards  $J^\pi$  defined by:

$$J^\pi \equiv E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t))\right]$$

### B. Dynamic Programming

Each MDP has one or more optimal policies, with optimal being a maximized expected sum of (discounted) rewards. To compute an optimal policy, dynamic programming (DP) methods use value functions to compute the expected utility of a certain action given a certain state. The value of a state  $V^\pi(s)$  denotes the expected cumulative discounted future reward when the agent starts in state  $s$  and follows policy  $\pi$ :

$$V^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, \pi\right]$$

Another useful function is the Q-function, which stores the expected value of a state-action pair.  $Q^\pi(s, a)$  denotes the expected cumulative discounted future reward when the agent starts in state  $s$ , takes action  $a$ , and then follows policy  $\pi$ :

$$Q^\pi(s, a) = E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, a_0 = a, \pi\right]$$

If the optimal Q-function  $Q^*$  is known, the agent can select optimal actions by always taking the action with the largest Q-value in each state. It can be easily shown that the optimal state value equals the highest state-action value in each state.

Bellman [21] has shown that the optimal state-action value of a state relates to the optimal state values of the states that can be reached in a single step:

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s'|s, a) V^*(s'),$$

where  $V^*(s') = \max_{a'} Q^*(s', a')$ . A number of efficient DP techniques have been developed that solve this set of non-linear equations. Value iteration for example uses the Bellman equation iteratively on all state-action pairs in order to compute an optimal Q-function:

$$Q(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s'|s, a) V(s'),$$

where  $V(s') = \max_{a'} Q(s', a')$ . After a finite amount of such updates for all states and actions, the optimal Q-function  $Q^*$  is computed, resulting in an optimal policy for the MDP.

### C. Model-based Reinforcement Learning

One of the biggest disadvantages of dynamic programming is that the transition and reward functions should be known a-priori. This means that a human has to design the MDP for a particular problem, which can be time-consuming, erroneous, or infeasible for particular problems.

Model-based RL methods learn to estimate the model from the experiences of the agent interacting with the environment. Although model-learning can be hard for particular environments with continuous or huge state spaces, for moderately sized problems it can be done efficiently. Furthermore, some comparisons have shown model-based RL to be much more effective than model-free methods such as Q-learning [18]. This is why we want to use model-based RL for solving multi-objective MDPs which are usually difficult to solve due to the existence of many Pareto optimal policies.

Model-based RL methods learn the transition and reward models of the environment by making use of counters that are used in a maximum-likelihood way to compute approximate transition probabilities and average rewards. Each time the agent selects action  $a$  in state  $s$  and makes a transition to state  $s'$ , the transition model's counter values  $C(s, a)$  and  $C(s, a, s')$  are increased by one. In a similar fashion, the obtained reward  $r$  is added to the value  $R^T(s, a)$  which computes the sum of all rewards obtained by selecting action  $a$  in state  $s$ . Finally, the maximum likelihood model of the MDP is computed as:

$$T'(s'|s, a) = \frac{C(s, a, s')}{C(s, a)} \quad \text{and} \quad R'(s, a) = \frac{R^T(s, a)}{C(s, a)}$$

At each moment in time, the estimated model can be used by a dynamic programming method to compute a new policy. Some algorithms such as prioritized sweeping have been proposed [17] which use a smart strategy to re-compute the Q-function with much less computational effort.

## III. MODEL-BASED MULTI-OBJECTIVE REINFORCEMENT LEARNING

In multi-objective Markov decision processes (MOMDPs), instead of the usual scalar reward function  $R(s, a)$ , a reward vector  $\vec{R}(s, a)$  is used. The vector  $\vec{R}(s, a)$  consists of  $l$  dimensions or components representing the different objectives. This means that a reward function  $\vec{R}(s, a) = (R_1(s, a), \dots, R_l(s, a))$  is given in an MOMDP that returns the expected reward vector for each state-action pair.

The reward vector in MOMDPs has multiple components, and therefore conflicts can arise between them. For example, it may be possible that an agent is able to compute a policy for a very nice sight-seeing tour, but that the cost of this tour is larger than the cost of other possible (less nice) tours. Therefore, the algorithms have to deal with several trade-offs and compute all policies that are not dominated by another policy. This set of policies is called the *Pareto efficient (optimal) set*. When given the set of Pareto optimal policies, we can let a user or autonomous agent select one of them given their preferences at that time.

To solve MOMDPs we can use particular multi-objective dynamic programming (MODP) algorithms that compute the Pareto front of optimal policies and corresponding value functions. These MODP methods are based on value iteration [20] or policy iteration [22], [23] and extend conventional dynamic programming algorithms by computing sets of policies and value functions.

### A. Multi-objective Dynamic Programming

The solution of most multi-objective DP methods is to keep track of all non-dominated (or Pareto efficient) value functions and policies. The multi-objective value function has different cumulative reward components or values and this is denoted by a value function  $V^i(s) = (V_1^i(s), V_2^i(s), \dots, V_l^i(s))$ , where  $V_x^i$  denotes the discounted cumulative reward intake of reward component  $x$  of policy  $i$ .

The dominance function  $\succ$  works on two value vectors for a state  $s$  as follows:

$$V^i(s) \succ V^j(s) \Leftrightarrow \exists x; V_x^i(s) > V_x^j(s) \wedge \neg \exists y; V_y^i(s) < V_y^j(s)$$

So a value vector of policy  $i$  dominates a value vector of policy  $j$  if policy  $i$  has a higher value on some component  $x$  and does not have a lower value on any other component.

We will use  $V^O$  to denote the set of value functions (and policies) that are not dominated:

$$V^O(s) = \{V^i(s) | V^i(s) \text{ is not dominated by a policy in } s\}$$

Furthermore,  $V^D$  is used for denoting the set of value functions computed at some given moment that may include dominated ones.

The same is done for the Q-function, thus the algorithm keeps track of a set  $Q^O$  that denotes the set of non-dominated Q-functions. Here the Q-functions should not be dominated by another Q-vector for the same state-action pair. Thus:

$$Q^O(s, a) = \{Q^i(s, a) | Q^i(s, a) \text{ is not dominated in } s, a\}$$

The non-dominated operator  $ND$  tells whether a policy  $i$  is not dominated in the state  $s$  by any value function vector of the set of value functions  $V^D(s)$ :

$$ND(V^i(s), V^D(s)) \Leftrightarrow \neg \exists V^j(s) \in V^D(s); V^j(s) \succ V^i(s)$$

Analogous definitions hold for Q-vector functions.

### B. CON-MODP

For our novel model-based MORL algorithm, we combine model-building methods with a multi-objective dynamic programming method. In this paper we will make use of the CON-MODP algorithm [19] that focuses on only computing stationary deterministic Pareto optimal policies. This makes this method more effective than previous dynamic programming methods for solving MOMDPs, although the drawback is that CON-MODP assumes deterministic problems.

With the previous formulations it is now possible to enhance dynamic programming to compute non-dominated value function sets. For simplicity we restrict ourselves to deterministic MDPs. We define the Pareto optimal operator  $PO$  as:

$$PO(Q^D(s, a)) = \{Q^i(s, a) | Q^i(s, a) \in Q^D(s, a) \wedge ND(Q^i(s, a), Q^D(s, a))\}$$

Furthermore, we construct the dynamic programming operator as follows where  $\oplus$  denotes an addition operator working on sets (and vectors) and  $s'$  is the next state:

$$DP(Q^D(s, a)) = (\vec{R}(s, a) \oplus \gamma V^O(s')) | P(s, a, s') = 1.0$$

Here  $V^O(s)$  is computed as:

$$V^O(s) = PO(\cup_a Q^D(s, a))$$

Note that this dynamic programming operator is defined for deterministic environments (therefore  $P(s, a, s') = 1$  for some  $s'$ ). For stochastic environments, this operator should be changed, but this could lead to a huge increase in the number of Pareto optimal policies.

It is known that in infinite horizon discounted Markov decision processes, there is always a single optimal value function and one or more stationary deterministic policies belonging to it. The CON-MODP algorithm exploits this fact and uses a consistency operator that eliminates most non-stationary policies and reevaluates non-stationary (or inconsistent) policies that are only inconsistent in a single state that is being expanded. Dynamic programming methods usually expand each state and this could sometimes lead to temporally inconsistent policies where different actions are selected in the same state at different time-steps.

CON-MODP detects when a policy is made inconsistent due to the last lookahead update step, and then changes it to a consistent policy by forcing the last action in the current state that is evaluated all the times this state will be visited. In this way, the policy is consistent again.

Because the sequence of actions from which the value function has been computed is changed, policy evaluation is used by CON-MODP to compute the true value of the consistent

policy. The CON-MODP algorithm uses the operators  $CON$ ,  $PO$ , and  $DP$ . The  $CON$  operator may make the policy  $\pi$  consistent (in case it was not) and recomputes its Q-value vector as:

$$\begin{aligned} CON(Q^\pi(s, a)) &= \pi \text{ with } Q^\pi(s, a), \text{ if } \pi(s) = a \\ &= \pi' \text{ with } Eval(\pi'), \text{ if } \pi(s) \neq a \text{ is the only} \\ &\quad \text{inconsistency, and } ND(Q^\pi(s, a), Q^O(s, a)) \\ &\quad \text{where } \pi'(s) = a \text{ and } \pi'(s') = \pi(s') \forall s' \neq s \\ &= \emptyset, \text{ otherwise} \end{aligned}$$

Here  $Eval(\pi')$  means that policy  $\pi'$  is evaluated using policy evaluation for the same number of steps as the original policy  $\pi$  is computed. Evaluation is only done if the original inconsistent policy was not already dominated. If the previous policy was dominated, the policy can be discarded anyway, since its values will never be larger than those of the previous policy.

We also let  $CON$  work immediately on sets of Q-functions and policies. The CON-MODP algorithm is now defined as:

$$Q^{O*} = (PO(CON(DP(Q_0))))^*$$

Where  $X(\cdot)^*$  means that operator  $X$  is repeated until convergence, and  $Q_0$  is the initial Q-vector containing for example only zero-values. This algorithm is able to compute all stationary Pareto optimal policies and corresponding value functions for deterministic finite MOMDPs.

### C. Exploration Strategies

We have developed two different exploration policies to obtain experiences by interacting with the environment from which the model is learned. The simplest exploration strategy is the Random-exploration approach, where actions are taken randomly and the findings are used to update the model. This will finally converge to a perfect environmental model, but faster methods exist. The other exploration strategy we used, is a Least-Visited approach, which chooses the actions to take according to the times it has explored a certain state-action combination.

1) *Least-Visited exploration:* As stated before, Least-Visited exploration chooses the actions to take according to the times it has taken those actions before. More specifically, it keeps track of the actions in each state and counts the times it has explored them. This gives a value for each possible action in each possible state. Action selection is then done by choosing the action with the lowest counter. After an action is taken, its corresponding state-action counter is incremented. If there are multiple actions with the lowest state-action counter value, the last equal action is chosen.

2) *Random-exploration:* In Random-exploration, actions are chosen completely randomly. The agent starts at its starting state and takes randomly chosen actions until it lands in a final state, while observing and updating rewards and transition probabilities in its model.

#### IV. EXPERIMENTAL RESULTS

In order to test our approach, we have used a multi-objective reinforcement learning problem known as the ‘‘Deep Sea Treasure’’ [5]. The Deep Sea Treasure environment states an episodic task where an agent has to explore the sea bottom for treasures. The world consists of a  $10 \times 11$  grid with 10 goal states. The goal state value is increased as the location is further away from the starting state. There are four deterministic actions possible; Go Up, Go Down, Go Left, Go Right. The two objectives are to minimize the number of steps taken before reaching a goal state, and maximizing the goal reward. The agent receives a reward of -1 for every step taken (the first objective). If a goal state is reached, the goal reward is the value of the treasure (the second objective). The fastest path to each goal state is part of the Pareto optimal set, so the Pareto front exists of 10 policies. The distribution of the goal state values causes an entirely concave Pareto front. Figure 1 shows an illustration of the problem.

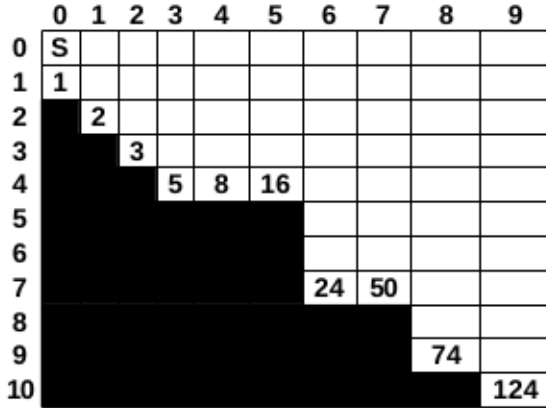


Fig. 1. A visual representation of the Deep Sea Treasure world. The starting state is shown as *S* and the numbers in the cells represent the goal state values.

**Set-up** We will compare our model-based MORL method with the two exploration strategies to Scalarized Q-Learning [9]. This method uses scalarization functions in order to transform the multi-objective Q-values into a single Q-value and then selects actions with an  $\epsilon$ -greedy exploration strategy. Because this approach works with linear scalarization functions, it can only find Pareto optimal policies in a convex Pareto front. Therefore, in the Deep Sea Treasure world, it can only find the  $(-1,1)$  and the  $(-19,124)$  solutions. We have supplied this method with scalarization weights of  $(1, 0)$  and  $(0, 1)$  for the two objectives.

Our results are collected and averaged over 10 trials of each 2000 epochs. For the Scalarized Q-learning algorithm, the learning rate  $\alpha = 0.1$ ,  $\gamma = 0.9$  and  $\epsilon = 0.1$ . For our model-based MORL methods we set  $\gamma = 0.9$  as well. The maximum amount of steps to find a solution was 1000.

##### A. Results

We first present the cardinality results, which are the number of Pareto optimal policies found after a specific number of

epochs. For the three methods we use a greedy exploration strategy to compute the Pareto optimal policies after each 10 epochs. As mentioned before, the maximal number of Pareto optimal policies is 10 for the Deep Sea Treasure problem. Table I shows results after 200, 500, 1000, and 2000 epochs. The table clearly shows that our model-based method with the random or least-visited exploration strategy significantly outperforms the scalarized Q-learning algorithm. Furthermore, the results also show that our model-based MORL method with the least-visited exploration strategy finds all Pareto optimal solutions within 2000 epochs.

TABLE I

The number of different Pareto optimal policies found during the learning process and the standard deviations.

Epochs	Q-learning	Random	Least-visited
200	$1.1 \pm 0.3$	$6.6 \pm 0.8$	$7.8 \pm 0.8$
500	$1.2 \pm 0.4$	$8.3 \pm 0.8$	$9.4 \pm 0.5$
1000	$1.8 \pm 0.4$	$9.1 \pm 0.7$	$9.4 \pm 0.5$
2000	$2.0 \pm 0.0$	$9.6 \pm 0.5$	$10.0 \pm 0.0$

In Figure 2 we also show the entire learning performance of the three methods. The model-based MORL methods learn to find a number of Pareto optimal policies much faster and will finally converge to having learned all 10 of them.

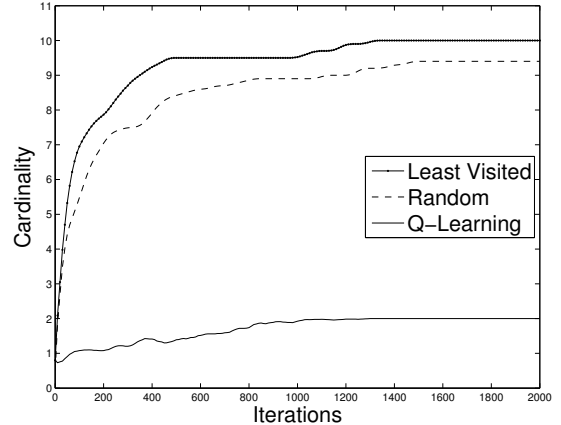


Fig. 2. Learning performance on the Deep Sea Treasure problem measured in the number of found Pareto optimal policies.

We also show the results of the three methods using the hypervolume assessment metric. The hypervolume gives the volume between the summed rewards obtained by the Pareto optimal policies and the reference point  $(-25,0)$ . Table II shows that the hypervolume is maximized after 2000 epochs by the model-based MORL method with the least-visited exploration strategy.

The hypervolume has as maximal value 1155. It is interesting to see that the hypervolume of our best method after 200 epochs is already larger than the one of the scalarized Q-learning algorithm.

We finally show the learning performance as measured by the hypervolume assessment function in Figure 3.

TABLE II

The values of the Hypervolume unary indicator during the learning process and the standard deviations.

Epochs	Q-learning	Random	Least-visited
200	98 ± 2331	686 ± 264	852 ± 164
500	172 ± 311	890 ± 211	1101 ± 90
1000	614 ± 311	971 ± 183	1101 ± 90
2000	762 ± 0	1055 ± 141	1155 ± 0

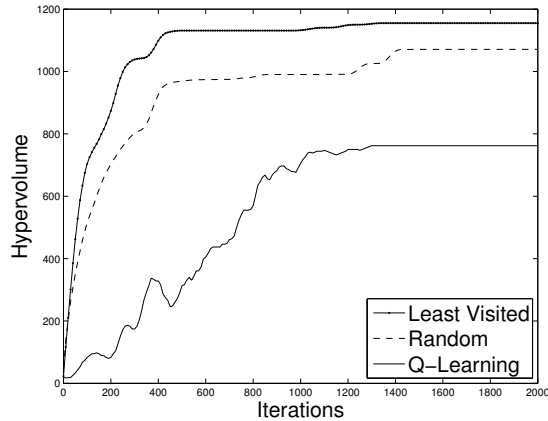


Fig. 3. Learning performance on the Deep Sea Treasure problem measured using the hypervolume unary indicator.

## B. Discussion

The two policies that the scalarized Q-learning algorithm finds are the two ends of the Pareto front. This confirms the results obtained in [9]. The model-based MORL methods will both converge to learning all Pareto optimal policies. If every state-action pair is visited at least one time, the CON-MODP algorithm is able to compute all Pareto optimal solutions. Some other methods in MORL literature have been used for the Deep Sea Treasure problem, but none of them was as effective as our method. Finally, we want to note that a single call of CON-MODP does not cost more than one second for this problem.

## V. CONCLUSION AND FUTURE WORK

This paper described a model-based reinforcement learning algorithm for solving deterministic multi-objective reinforcement learning problems. The method learns a model of the environment while the agent is interacting with it, and uses a previously introduced multi-objective dynamic programming method to compute the set of Pareto optimal policies. The results on the Deep Sea Treasure problem have shown that the proposed method is near-optimal for this problem. It finds all Pareto optimal policies within a short learning period. The only way to improve on our method would be to develop a smarter exploration method.

Although the current results are promising, in future work we intend to use multi-objective dynamic programming methods that can handle stochastic environments. Although this may lead to a slow algorithm, the structure of the problem may be exploited to make such algorithms faster.

## REFERENCES

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. The MIT press, Cambridge MA, A Bradford Book, 1998.
- [2] M. Wiering and M. van Otterlo, *Reinforcement Learning: State-of-the-art*. Springer, 2012.
- [3] C. J. C. H. Watkins, "Learning from delayed rewards," Ph.D. dissertation, King's College, Cambridge, England, 1989.
- [4] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley, "A survey of multi-objective sequential decision-making," *J. Artif. Intell. Res. (JAIR)*, vol. 48, pp. 67–113, 2013.
- [5] P. Vamplew, R. Dazeley, A. Berry, R. Issabekov, and E. Dekker, "Empirical evaluation methods for multiobjective reinforcement learning algorithms," *Machine Learning*, vol. 84, no. 1-2, pp. 51–80, 2011.
- [6] K. van Moffaert, M. M. Drugan, and A. Nowe, "Learning sets of Pareto optimal policies," in *Thirteenth International Conference on Autonomous Agents and Multiagent Systems - Adaptive Learning Agents Workshop (ALA)*, 2014.
- [7] E. Zitzler, L. Thiele, M. Laumanns, C. M. Fonseca, and V. da Fonseca, "Performance assessment of multiobjective optimizers: An analysis and review," *IEEE T. on Evol. Comput.*, vol. 7, pp. 117–132, 2003.
- [8] G. Eichfelder, *Adaptive Scalarization Methods in Multiobjective Optimization*. Springer, 2008.
- [9] K. V. Moffaert, M. M. Drugan, and A. Nowe, "Scalarized multi-objective reinforcement learning: Novel design techniques," in *Proceedings of Adaptive Dynamic Programming and Reinforcement Learning (ADPRL)*, 2013, pp. 191–199.
- [10] D. J. Lizotte, M. Bowling, and S. A. Murphy, "Efficient reinforcement learning with multiple reward functions for randomized clinical trial analysis," in *Proceedings of the Twenty-Seventh International Conference on Machine Learning (ICML)*, 2010.
- [11] S. Q. Yahyaa, M. M. Drugan, and B. Manderick, "Exploration vs exploitation in the multi-objective multi-armed bandit problem," in *Proceedings of International Joint Conference of Neural Networks (IJCNN)*, 2014.
- [12] M. M. Drugan and A. Nowe, "Designing multi-objective multi-armed bandits: a study," in *Proceedings of International Joint Conference of Neural Networks (IJCNN)*, 2013.
- [13] —, "Scalarization based Pareto optimal set of arms identification algorithms," in *Proc of International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2014.
- [14] Z. Gabor, Z. Kalmar, and C. Szepesvari, "Multi-criteria reinforcement learning," in *Proceedings of the Fifteenth International Conference on Machine Learning*, M. K. P. Inc., Ed., 1998.
- [15] K. van Moffaert, M. M. Drugan, and A. Nowe, "Hypervolume-based multi-objective reinforcement learning," in *Proc of Evolutionary Multi-objective Optimization (EMO)*. Springer, 2013.
- [16] W. Wang and M. Sebag, "Hypervolume indicator and dominance reward based multi-objective Monte-Carlo tree search," *Machine Learning*, vol. 92, no. 2-3, pp. 403–429, 2013.
- [17] A. W. Moore and C. G. Atkeson, "Prioritized sweeping: Reinforcement learning with less data and less time," *Machine Learning*, vol. 13, pp. 103–130, 1993.
- [18] M. A. Wiering, "Explorations in efficient reinforcement learning," Ph.D. dissertation, University of Amsterdam, February 1999.
- [19] M. A. Wiering and E. D. de Jong, "Computing optimal stationary policies for multi-objective markov decision processes," in *Proc of Approximate Dynamic Programming and Reinforcement Learning (ADPRL)*. IEEE, 2007, pp. 158–165.
- [20] D. White, "Multi-objective infinite-horizon discounted Markov decision processes," *Journal of Mathematical Analysis and Applications*, vol. 89, pp. 639–647, 1982.
- [21] R. Bellman, *Dynamic Programming*. Princeton University Press, 1957.
- [22] L. Thomas, "Constrained Markov decision processes as multi-objective problems," in *Multi-Objective Decision Making*, S. French, L. Thomas, R. Hartley, and D. White, Eds. Academic Press, 1983, pp. 77–94.
- [23] N. Furukawa, "Vector valued Markovian decision processes within countable state space," in *Recent Developments in Markov Decision Processes*, R. Hartley, L. Thomas, and D. White, Eds. Academic Press, New York, 1980, pp. 205–223.