

Theoretical Results on Reinforcement Learning with Temporally Abstract Behaviors

Doina Precup¹, Richard S. Sutton¹, and Satinder Singh²

¹ Department of Computer Science
University of Massachusetts
Amherst, MA 01003-4610
<http://www.cs.umass.edu/~dprecup|~rich>

² Department of Computer Science
University of Colorado
Boulder, CO 80309-0430
<http://www.cs.colorado.edu/~baveja>

Abstract. We present new theoretical results on planning within the framework of temporally abstract reinforcement learning (Precup & Sutton, 1997; Sutton, 1995). Temporal abstraction is a key step in any decision making system that involves planning and prediction. In temporally abstract reinforcement learning, the agent is allowed to choose among "behaviors", whole courses of action that may be temporally extended, stochastic, and contingent on previous events. Examples of behaviors include closed-loop policies such as picking up an object, as well as primitive actions such as joint torques. Knowledge about the consequences of behaviors is represented by special structures called multi-time models. In this paper we focus on the theory of planning with multi-time models. We define new Bellman equations that are satisfied for sets of multi-time models. As a consequence, multi-time models can be used interchangeably with models of primitive actions in a variety of well-known planning methods including value iteration, policy improvement and policy iteration.

1 Introduction

Model-based reinforcement learning offers a possible solution to the problem of integrating planning with real-time learning and decision-making [19]. However, conventional model-based reinforcement learning uses one-step models [11, 13, 17], that cannot represent common-sense, higher-level actions, such as picking an object or traveling to a specified location.

Several researchers have proposed extending reinforcement learning to a higher level by treating entire closed-loop policies as actions [3–6, 9, 10, 12, 16]. In order to use such actions in planning, an agent needs the ability to create and handle models at a variety of different, interrelated levels of temporal abstraction. Sutton [18] introduced an approach to modeling at different time scales, based on prior work by Singh [16], Dayan [2] and by Sutton and Pinette [20].

This approach enables models of the environment at different temporal scales to be intermixed, producing temporally abstract models. In previous work [14], we generalized this approach from the prediction case, to the full control case. In this paper, we summarize the framework of temporally abstract reinforcement learning and present new theoretical results on planning with general *behaviors* and temporally abstract models of behaviors.

Behaviors are similar to AI’s classical “macro operators” [7, 8, 15], in that they can take control for some period of time, determining the actions during that time, and in that one can choose among behaviors much as one originally chose among primitive actions. However, classical macro operators are only a fixed sequence of actions, whereas behaviors incorporate a general (possibly non-Markov) closed-loop policy and completion criterion. These generalizations are required when the environment is stochastic and uncertain with general goals, as in reinforcement learning and Markov decision processes (MDP).

The predictive knowledge needed in order to plan using behaviors can be represented through *multi-time models* [14]. Such models summarize several time scales and have the ability to predict events that can happen at various unknown moments. In this paper, we focus on the theoretical properties of multi-time models. We show formally that such models can be used interchangeably with models of primitive actions in a variety of well-known dynamic programming methods, while preserving the same guarantees of convergence to correct solutions. The benefit of using such temporally extended models is a significant improvement in the convergence rates of these algorithms.

2 Reinforcement Learning (MDP) Framework

First we briefly summarize the mathematical framework of the reinforcement learning problem that we use in the paper. In this framework, a learning *agent* interacts with an *environment* at some discrete, lowest-level time scale $t = 0, 1, 2, \dots$. At each time step, the agent perceives the state of the environment, \mathbf{s}_t , and on that basis chooses a primitive action, \mathbf{a}_t . In response to each primitive action, \mathbf{a}_t , the environment produces one step later a numerical reward, r_{t+1} , and a next state, \mathbf{s}_{t+1} .

The agent’s objective is to learn a policy π , which is a mapping from states to probabilities of taking each action, that maximizes the expected discounted future reward from each state \mathbf{s} :

$$V^\pi(\mathbf{s}) = E \left\{ r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \mid s_0 = \mathbf{s}, \pi \right\},$$

where $\gamma \in [0, 1)$ is a *discount-rate* parameter. The quantity $V^\pi(\mathbf{s})$ is called the *value* of state \mathbf{s} under policy π , and V^π is called the value function for policy π . The optimal value of a state is denoted

$$V^*(\mathbf{s}) = \max_{\pi} V^\pi(\mathbf{s})$$

The environment is henceforth assumed to be a stationary, finite Markov decision process. We assume that the states are discrete and form a finite set,

$s_t \in \{1, 2, \dots, n\}$. The latter assumption is a temporary theoretical convenience; it is not a limitation of the ideas we present.

3 Behaviors

In order to achieve faster planning, the agent should be able to predict what happens if it follows a certain course of action over a period of time. By a “course of action” we mean any way of behaving, i.e. any way of mapping representations of states to primitive actions. A behavior is a way of choosing actions that is initiated, takes control for some period of time, and then eventually ends. Behaviors are defined by three elements:

- the set of states \mathcal{I} in which the behavior applies
- a decision rule μ which specifies what actions are executed by the behavior
- a completion function β which specifies the probability of completing the behavior on every time step.

μ is of a slightly more general form than the policies used in conventional reinforcement learning. The first generalization is based on the observation that primitive actions qualify as behaviors: they are initiated in a state, take control for a while (one time step), and then end. Therefore, we allow μ to choose among behaviors, rather than only among primitive actions.

The conventional reinforcement learning setting also requires that a policy’s decision probabilities at time t should be a function only of the current state s_t . This type of policy is called *Markov*. For the policy of a behavior, μ , we relax this assumption, and we allow the probabilities of selecting a sub-behavior to depend on all the states and actions from time t_0 , when the behavior began executing, up through the current time, t . We call policies of this more general class *semi-Markov*.

The completion function β can also be semi-Markov. This property enables us to describe various kinds of completion. Perhaps the simplest case is that of a behavior that completes after some fixed number of time steps (e.g. after 1 step, as in the case of primitive actions, or after n steps, as in the case of a classical macro operator consisting of a sequence of n actions). A slightly more general case is that in which the behavior completes during a certain time period, e.g. 10 to 15 time steps later. In this case, the completion function β should specify the probability of completion for each of these time steps. The case which is probably the most useful in practice is the completion of a behavior with the occurrence of a critical state, often a state that we think of as a subgoal. For instance, the behavior **pick-up-the-object** could complete when the object is in the hand. This event occurs at a very specific moment in time, but this moment is indefinite, not known in advance. In this case, the completion function depends on the state history of the system, rather than explicitly on time. Lastly, if β is always 0, the behavior does not complete. This is the case of usual policies used in reinforcement learning.

Behaviors are typically executed in a *call-and-return* fashion. Each behavior can be viewed as a “subroutine” which calls sub-behaviors, according to its internal policy μ . When a sub-behavior b' is selected, it takes control of the action choices until it completes. Upon completion, b' transfers the control back to b . b also inherits the current time t' and the current state $s_{t'}$, and has to decide if it should terminate or pick a new sub-behavior.

Two behaviors, a and b , can be *composed* to yield a new behavior, denoted ab , that first follows a until it terminates and then follows b until it terminates, also terminating the composed behavior ab .

4 Models of Behaviors

Planning in reinforcement learning refers to the use of models of the effects of actions to compute value functions, particularly V^* . We use the term *model* for any structure that generates predictions based on the representation of the state of the system and on the course of action that the system is following.

In order to plan at the level of behaviors, we need to have a knowledge representation form that predicts their consequences. The *multi-time model* of a behavior characterizes the states that result upon the behavior’s completion and the truncated return received along the way when the behavior is executed in various states [14]. Let \mathbf{p} be an n -vector and g a scalar. The pair \mathbf{p}, g is an *accurate prediction* for the execution of semi-Markov behavior b in state s if and only if

$$g = E \left\{ r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{T-1} r_T \mid s_t = s, b \right\} \quad (1)$$

and

$$\mathbf{p} = E \left\{ \gamma^T \mathbf{s}_T \mid s_t = s, b \right\}, \quad (2)$$

where T is the random variable denoting the time at which b terminates when executed in $s_t = s$, and \mathbf{s}_T denotes the unit basis n -vector corresponding to s_T . \mathbf{p} is called the *state prediction vector* and g is called the *reward prediction* for state s and behavior b .

We will use the notation “ \cdot ”, as in $\mathbf{x} \cdot \mathbf{y}$, to represent the inner or dot product between vectors. We will refer to the i -th element of a vector \mathbf{x} as $\mathbf{x}(i)$. The transpose of a vector \mathbf{x} will be denoted by \mathbf{x}^T .

The predictions corresponding to the states $s \notin \mathcal{I}$ are always 0. If the behavior never terminates, the reward prediction g is equal to the value function of the internal policy of the behavior, and the elements of \mathbf{p} are all 0. In general, for any behavior and any starting state, there exists a constant ϵ_b such that $\|\mathbf{p}\|_1 \leq \epsilon_b < 1$, where $\|\cdot\|_1$ denotes the l_1 -norm of a vector: $\|\mathbf{x}\|_1 = \sum_i \mathbf{x}(i)$.

The state prediction vectors for the same behavior in all the states are often grouped as the rows of an $n \times n$ *state prediction matrix*, \mathbf{P} , and the reward predictions are often grouped as the components of an n -component *reward prediction vector*, \mathbf{g} . \mathbf{g}, \mathbf{P} form the *accurate model* of the behavior.

A key theoretical result refers to the way in which the model of a composed behavior can be obtained from the models of its components.

Theorem 1 (Composition or Sequencing). *Given an accurate prediction g_a, \mathbf{p}_a for some behavior a applied in state s and an accurate model g_b, \mathbf{P}_b for some behavior b , the prediction:*

$$g = g_a + \mathbf{p}_a \cdot \mathbf{g}_b \quad \text{and} \quad \mathbf{p} = \mathbf{p}_a^T \mathbf{P}_b \quad (3)$$

is an accurate prediction for the composed behavior ab when it is executed in s .

Proof. Let k be the random variable denoting the time at which a completes, and T be the random variable denoting the time at which ab completes. Then we have:

$$\begin{aligned} g &= E\{r_{t+1} + \dots + \gamma^{T-1} r_T \mid s_t = s, ab\} \\ &= E\{r_{t+1} + \dots + \gamma^{k-1} r_k \mid s_t = s, a\} + E\{\gamma^k r_{k+1} + \dots + \gamma^{T-1} r_T \mid s_t = s, ab\} \\ &= g_a + \sum_{i=1}^{\infty} P\{k=i \mid s_t = s, a\} \gamma^i \sum_{s'} P\{s_k = s' \mid s_t = s, k=i\} \\ &\quad E\{r_{k+1} + \dots + \gamma^{T-k-1} r_T \mid s_k = s', b\} \\ &= g_a + \sum_{s'} \sum_{i=1}^{\infty} P\{k=i \mid s_t = s, a\} \gamma^i P\{s_k = s' \mid s_t = s, k=i\} g_b(s') \\ &= g_a + \sum_{s'} p_a(s') g_b(s') \\ &= g_a + \mathbf{p}_a \cdot \mathbf{g}_b \end{aligned}$$

The equation for \mathbf{p} follows similarly. □

A simple combination of models can also be used to predict the effect of probabilistic choice among behaviors:

Theorem 2 (Averaging or Choice). *Let g_i, \mathbf{p}_i be a set of accurate predictions for the behaviors b_i executed in state s , and let $w_i > 0$ be a set of numbers such that $\sum_i w_i = 1$. Then the prediction defined by:*

$$g = \sum_i w_i g_i \quad \text{and} \quad \mathbf{p} = \sum_i w_i \mathbf{p}_i \quad (4)$$

is accurate for the behavior b , which chooses in state s among sub-behaviors b_i with probabilities w_i and then follows the chosen b_i until completion.

Proof.

$$g = E\{r_{t+1} + \dots + \gamma^{T-1} r_T \mid s_t = s, b\} = \sum_i w_i E\{r_{t+1} + \dots + \gamma^{T-1} r_T \mid s_t = s, b_i\} = \sum_i w_i g_i$$

The equation for \mathbf{p} follows similarly. □

5 Planning with Models of Behaviors

In this section, we extend the theoretical results of dynamic programming for the case in which the agent is allowed to use an arbitrary set of behaviors, \mathcal{B} . If \mathcal{B} is exactly the set of primitive actions, then our results degenerate to the conventional case. We assume, for the sake of simplicity, that for any state s , the set of behaviors that apply in s , denoted \mathcal{B}_s , is always non-empty. However, the theory that we present extends, with some additional complexity, to the case in which no behaviors apply in certain states.

For every state s , the value of a policy π that chooses among the behaviors from \mathcal{B} is defined similarly to the conventional case, as the expected discounted reward if the policy is applied starting in s :

$$v^\pi(s) = E \left\{ r_{t+1} + \gamma r_{t+2} + \dots \mid s_t = s, \pi \right\}. \quad (5)$$

The *optimal value function, given the set \mathcal{B}* , can be defined as

$$v_{\mathcal{B}}^*(s) = \sup_{\pi \in \Pi_{\mathcal{B}}} v^\pi(s), \quad (6)$$

for all s , where $\Pi_{\mathcal{B}}$ is the set of policies that can be defined using the behaviors from \mathcal{B} . The value functions are sometimes represented as n -vectors, \mathbf{v}^π and $\mathbf{v}_{\mathcal{B}}^*$, with each component representing the value of a different state.

The value function of any Markov policy $\pi \in \Pi_{\mathcal{B}}$ satisfies the Bellman evaluation equations:

$$v^\pi(s) = \sum_{b \in \mathcal{B}_s} \pi(s, b) (g_b + \mathbf{p}_b \cdot \mathbf{v}^\pi), \quad \text{for all } s \in S, \quad (7)$$

where $\pi(s, b)$ is the probability of choosing behavior b in state s when acting according to π and g_b, \mathbf{p}_b is the accurate prediction for sub-behavior b in state s . We will show that, similarly to the conventional case, \mathbf{v}^π is the unique solution to this system of equations.

Similarly, the optimal value function given a set of behaviors \mathcal{B} satisfies the Bellman optimality equations:

$$v_{\mathcal{B}}^*(s) = \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*), \quad \text{for all } s \in S, \quad (8)$$

As in the conventional case, $\mathbf{v}_{\mathcal{B}}^*$ is the unique solution of this set of equations, and there exists at least one deterministic Markov policy $\pi_{\mathcal{B}}^* \in \Pi_{\mathcal{B}}$ that is optimal, i.e., for which $\mathbf{v}^{\pi_{\mathcal{B}}^*} = \mathbf{v}_{\mathcal{B}}^*$.

We will now present in detail the proofs leading to these theoretical results. The practical consequence is that all the usual update rules used in reinforcement learning can be used to compute \mathbf{v}^π and $\mathbf{v}_{\mathcal{B}}^*$ when the agent makes choices among behaviors. We focus first on the Bellman policy evaluation equations.

Theorem 3 (Value Functions for Composed Policies). *Let π be a policy that, when starting in state s , follows behavior b and, when the behavior completes, follows policy π' . Then*

$$v^\pi(s) = g_b + \mathbf{p}_b \cdot \mathbf{v}^{\pi'}. \quad (9)$$

Proof. According to the theorem statement, $\pi = b\pi'$. The conclusion follows immediately from the composition theorem. \square

Theorem 4 (Bellman Policy Evaluation Equation). *The value function of any Markov policy satisfies the Bellman policy evaluation equations (7).*

Proof. Using the averaging rule and the fact that π is Markov, we have:

$$v^\pi(s) = \sum_{b \in \mathcal{B}_s} \pi(s, b) v^{b\pi}(s)$$

We can expand $v^{b\pi}(s)$ using (9), to obtain the desired result:

$$v^\pi(s) = \sum_{b \in \mathcal{B}_s} \pi(s, b) (g_b + \mathbf{p}_b \cdot \mathbf{v}^\pi)$$

\square

We now prove that the Bellman evaluation equations have a unique solution, and that this solution can be computed using the well-known algorithm of *policy evaluation with successive approximations*:

- start with arbitrary initial values $v_0(s)$ for all states
- iterate for all s

$$v_{k+1}(s) \leftarrow g_{\pi(s)} + \mathbf{P}_{\pi(s)} \cdot \mathbf{v}_k,$$

where $\pi(s)$ is the behavior suggested by π in state s .

Proof. For any behavior b , we show that the operator $T_b(\mathbf{v}) = g_b + \mathbf{p}_b \cdot \mathbf{v}$ is a contraction with constant $\epsilon_b = \|\mathbf{p}_b\|_1$. For arbitrary vectors \mathbf{v} and \mathbf{v}' , we have:

$$T_b(\mathbf{v}) - T_b(\mathbf{v}') = \sum_{s'} p_b(s') (v(s') - v'(s')) \leq \sum_{s'} p_b(s') \|\mathbf{v} - \mathbf{v}'\|_\infty \leq \epsilon_b \|\mathbf{v} - \mathbf{v}'\|_\infty,$$

where $\|\cdot\|_\infty$ denotes the l_∞ -norm of a vector: $\|\mathbf{x}\|_\infty = \max_i x(i)$. The result follows from the contraction mapping theorem [1]. \square

So far we have established that the value functions of Markov behavior policies have similar properties with the Markov policies that use only primitive actions. Now we establish similar results for the optimal value function that can be obtained when planning with a set of behaviors.

Theorem 5 (Bellman Optimality Equation). *For any set of behaviors \mathcal{B} , the optimal value function $\mathbf{v}_\mathcal{B}^*$ satisfies the Bellman optimality equations (8).*

Proof. Let π be an arbitrary policy which, at time step 0, in state $s_0 = s$, chooses among the available behaviors with probabilities w_b .

$$v^\pi(s) = \sum_{b \in \mathcal{B}_s} w_b (g_b + \sum_{s'} p_b(s') W_\pi(s'))$$

where $W_\pi(s')$ is the expected discounted return from state s' on. Since $p_b(s') \geq 0, \forall s'$, using the definition of $v_{\mathcal{B}}^*(s')$, we have:

$$v^\pi(s) \leq \sum_{b \in \mathcal{B}_s} w_b(g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*) \leq \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*)$$

Since π is arbitrary, due to the definition of $v_{\mathcal{B}}^*(s)$,

$$v_{\mathcal{B}}^*(s) \leq \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*)$$

On the other hand, let $b_0 = \arg \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*)$. Let π be the policy that chooses b_0 at time step 0 and, after b_0 ends, in state s' , switches to a policy $\pi_{s'}$ such that $v^{\pi_{s'}}(s') \geq v_{\mathcal{B}}^*(s') - \epsilon$. $\pi_{s'}$ exists because of the way in which $v_{\mathcal{B}}^*$ was defined. From (9), we have:

$$v^\pi(s) = g_{b_0} + \sum_{s'} p_{b_0}(s') v^{\pi_{s'}}(s') \geq g_{b_0} + \sum_{s'} p_{b_0}(s') (v_{\mathcal{B}}^*(s') - \epsilon) \geq g_{b_0} + \mathbf{p}_{b_0} \cdot \mathbf{v}_{\mathcal{B}}^* - \epsilon$$

Since $v_{\mathcal{B}}^*(s) \geq v^\pi(s)$, we have:

$$v_{\mathcal{B}}^*(s) \geq \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*) - \epsilon$$

Since ϵ is arbitrary, it follows that

$$v_{\mathcal{B}}^*(s) = \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*)$$

The following step is to show that the solution of the Bellman optimality equations (8) is bounded and unique, and that it can be computed through a *value iteration algorithm*:

- start with arbitrary initial values $v_0(s)$
- iterate the update:

$$v_{k+1}(s) \leftarrow \max_{b \in \mathcal{B}_s} g_b + \mathbf{p}_b \cdot \mathbf{v}_k, \quad \forall s$$

Proof. For any set of actions \mathcal{B} let us consider the operator $T_{\mathcal{B}}$ which, in any state s , performs the following transformation: $T_{\mathcal{B}_s}(\mathbf{v}) = \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b \cdot \mathbf{v})$. Let \mathbf{v} and \mathbf{v}' be two arbitrary vectors. Then we have:

$$\begin{aligned} T_{\mathcal{B}_s}(\mathbf{v}) &= \max_{b \in \mathcal{B}_s} (g_b + \mathbf{p}_b(\mathbf{v}' + \mathbf{v} - \mathbf{v}')) = \max_{b \in \mathcal{B}_s} (T_b(\mathbf{v}') + \mathbf{p}_b \cdot (\mathbf{v} - \mathbf{v}')) \\ &\leq T_{\mathcal{B}_s}(\mathbf{v}') + \max_{b \in \mathcal{B}_s} \sum_{s'} p_b(s') \|\mathbf{v} - \mathbf{v}'\|_\infty = T_{\mathcal{B}_s}(\mathbf{v}') + \epsilon_{\mathcal{B}_s} \|\mathbf{v} - \mathbf{v}'\|_\infty, \end{aligned}$$

where $\epsilon_{\mathcal{B}_s} = \max_{b \in \mathcal{B}_s} \epsilon_b$. Similarly,

$$T_{\mathcal{B}_s}(\mathbf{v}') \leq T_{\mathcal{B}_s}(\mathbf{v}) + \epsilon_{\mathcal{B}_s} \|\mathbf{v} - \mathbf{v}'\|_\infty$$

Therefore, $\forall s$,

$$\|T_{\mathcal{B}}(\mathbf{v}) - T_{\mathcal{B}}(\mathbf{v}')\|_\infty \leq \epsilon_{\mathcal{B}} \|\mathbf{v} - \mathbf{v}'\|_\infty,$$

where $\epsilon_{\mathcal{B}} = \max_s \epsilon_{\mathcal{B}_s}$. Therefore $T_{\mathcal{B}}$ is a contraction with constant $\epsilon_{\mathcal{B}}$. The results follow from the contraction mapping theorem [1].

So far we have shown that the optimal value function $\mathbf{v}_{\mathcal{B}}^*$ is the unique bounded solution of the Bellman optimality equations. Now we will show that this value function can be achieved by a deterministic Markov policy:

Theorem 6 (Value Achievement). *The policy $\pi_{\mathcal{B}}^*$ defined as*

$$\pi_{\mathcal{B}}^*(s) = \arg \max_{b \in \mathcal{B}_s} g_b + \mathbf{p}_b \cdot \mathbf{v}_{\mathcal{B}}^*$$

achieves $\mathbf{v}_{\mathcal{B}}^$.*

Proof. For any arbitrary state s , we have:

$$0 \leq v_{\mathcal{B}}^*(s) - v^{\pi_{\mathcal{B}}^*}(s) = \mathbf{p}_{\pi_{\mathcal{B}}^*(s)} \cdot (\mathbf{v}_{\mathcal{B}}^* - \mathbf{v}^{\pi_{\mathcal{B}}^*}) \leq \epsilon_{\mathcal{B}_s} \|\mathbf{v}_{\mathcal{B}}^* - \mathbf{v}^{\pi_{\mathcal{B}}^*}\|_{\infty}.$$

Since $\epsilon_{\mathcal{B}} < 1$, $v_{\mathcal{B}}^*(s) = v^{\pi_{\mathcal{B}}^*}(s)$. □

In order to find the optimal policy given a set of behaviors, one can simply compute $\mathbf{v}_{\mathcal{B}}^*$ and then use it to pick actions greedily. Another popular planning method for computing optimal policies is *policy iteration*, which alternates steps of policy evaluation and policy improvement. We have already investigated policy evaluation, so we turn now to policy improvement:

Theorem 7 (Policy Improvement Theorem). *For any Markov policy π defined using actions from a set \mathcal{B} , let π' be a new policy which, for some state s , chooses greedily among the available behaviors, and then follows π*

$$\pi'(s) = \arg \max_{b \in \mathcal{B}_s} g_b + \mathbf{p}_b \cdot \mathbf{v}^{\pi},$$

Then $v^{\pi'}(s) \geq v^{\pi}(s)$.

Proof. Let b_0 be the behavior chosen by π in state s . Then, from (7),

$$v^{\pi}(s) = g_{b_0} + \mathbf{p}_{b_0} \cdot \mathbf{v}^{\pi} \leq \max_{b \in \mathcal{B}_s} g_b + \mathbf{p}_b \cdot \mathbf{v}^{\pi} = v^{\pi'}(s).$$

□

Given a set of behaviors \mathcal{B} such that \mathcal{B}_s is finite for all s , the policy iteration algorithm, which interleaves policy evaluation and policy improvement, converges to $\pi_{\mathcal{B}}^*$ in a finite number of steps.

The final result relates the models of behaviors to the optimal value function of the environment, V^* .

Theorem 8. *If g_b, \mathbf{p}_b is an accurate prediction for some behavior b in state s , then*

$$g_b + \mathbf{p}_b \cdot \mathbf{V}^* \leq V^*(s)$$

We say that accurate models are non-overpromising, i.e. they never promise more than the agent can actually achieve.

Proof. Assume that the agent can use b and all the primitive actions in the environment. Let $\pi = b\pi^*$, where π^* is the optimal policy of the environment. Then we have:

$$V^*(s) \geq v^{\pi}(s) = g_b + \mathbf{p}_b \cdot \mathbf{V}^*$$

□

6 Illustrations

The theoretical results presented so far show that accurate models of behaviors can be used in all the planning algorithms typically employed for solving MDPs, with the same guarantees of convergence to correct plans as in the case of primitive actions. We will now illustrate the speedup that can be obtained when using these methods in two simple gridworld learning tasks.

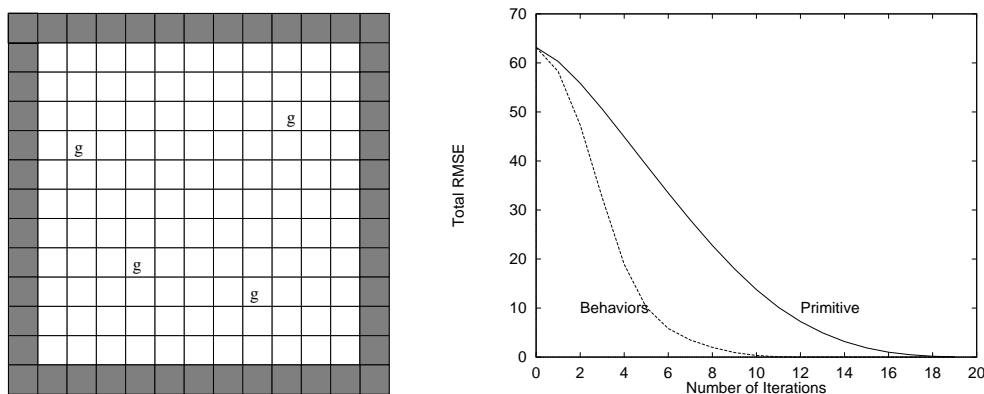


Fig. 1. Empty gridworld task. Behaviors allow the error in the value function estimation to decrease more quickly

The first task is depicted on the left panel of figure 1. The cells of the grid correspond to the states of the environment. From any state the agent can perform one of four primitive actions, **up**, **down**, **left** or **right**. With probability $2/3$, the actions cause the agent to move one cell in the corresponding direction (unless this would take the agent into a wall, in which case it stays in the same state). With probability $1/3$, the agent moves instead in one of the other three directions (unless this takes it into a wall, of course). There is no penalty for bumping into walls. In addition to these primitive actions, the agent can use four additional higher-level behaviors, to travel to each of the marked locations. These locations have been chosen randomly inside the environment. Accurate models for all the behaviors are also available. Both the behaviors and their models have been learned during a prior random walk in the environment, using Q-learning [21] and the β -model learning algorithm [18].

The agent is repeatedly given new goal positions and it needs to compute optimal paths to these positions as quickly as possible. In this experiment, we considered all possible goal positions. In each case, the value of the goal state is 1, there are no rewards along the way, and the discounting factor is $\gamma = 0.9$. We performed planning according to the standard value iteration method, where the starting values are $v_0(\mathbf{s}) = 0$ for all the states except the goal state, for which $v_0(\mathit{goal}) = 1$. In the first experiment, the agent was only allowed to use

primitive actions, while in the second case, it used both the primitive actions and the higher-level behaviors.

The right panel in figure 1 shows the root mean squared error in the estimate of the optimal value function over the whole environment. The use of higher-level behaviors introduces a significant speedup in convergence, even though the behaviors have been chosen arbitrarily. Note that an iteration using all the behaviors is slightly more expensive than an iteration using only primitive actions. This aspect can be improved by using more sophisticated methods of ordering the behaviors before doing the update. However, such methods are beyond the scope of this paper.

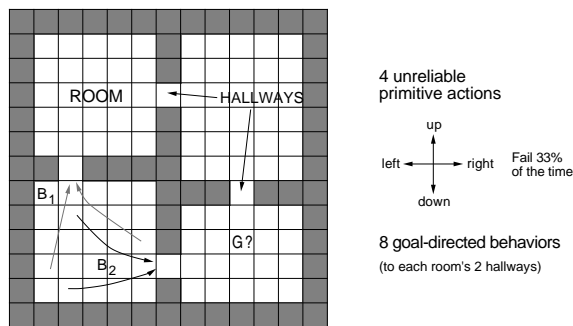


Fig. 2. Example Task. The natural behaviors are to move from room to room

In order to analyze in more detail the effect of behaviors, let us consider a second environment. In this case, the gridworld has four “rooms”. The basic dynamics of the environment are the same as in the previous case. For each state in a room, two higher-level behaviors are available, which can take the agent to each of the hallways adjacent to the room. Each of these behaviors has two outcome states: the target hallway, which corresponds to a successful outcome, and the state adjacent to the other hallway, which corresponds to failure (the agent has wandered out of the room). The completion function β is therefore 0 for all the states except these outcome states, where it is 1. The policy π underlying the behavior is the optimal policy for reaching the target hallway.

The goal state can have an arbitrary position in any of the rooms, but for this illustration let us suppose that the goal is two steps down from the right hallway. The value of the goal state is 1, there are no rewards along the way, and the discounting factor is $\gamma = 0.9$. We performed planning again, according to the standard value iteration method, with the same setting as in the previous task.

When using only primitive actions, the values are propagated one step on each iteration. After six iterations, for instance, only the states that are within six steps of the goal are attributed non-zero values. Figure 3 shows the value

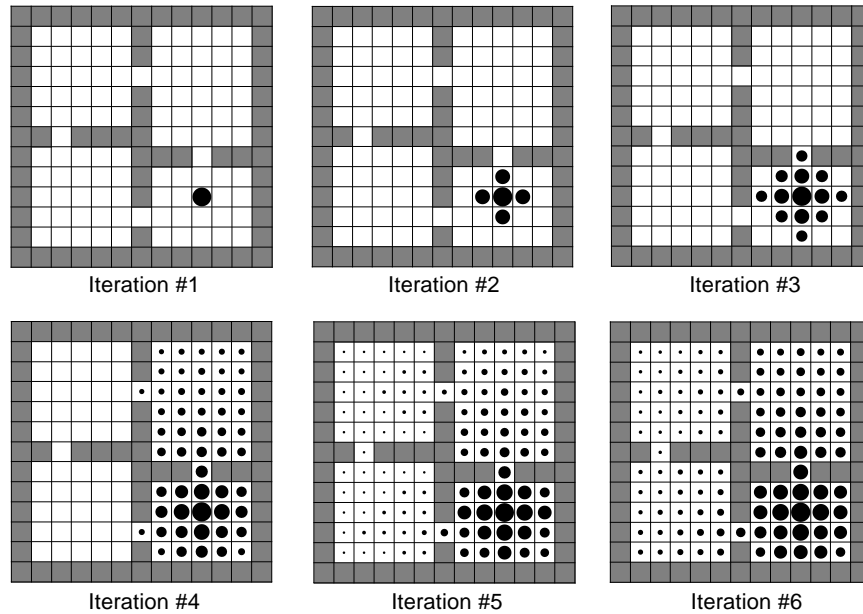


Fig. 3. Value iteration using primitive actions and higher-level behaviors

function after each iteration, using all available behaviors. The area of the circle drawn in each state is proportional to the value attributed to the state. The first three iterations are identical to the case in which only primitive actions are used. However, once the values are propagated to the first hallway, all the states in the rooms adjacent to the hallway receive values as well. For the states in the room containing the goal, these values correspond to performing the behavior of getting into the right hallway, and then following the optimal primitive actions to get to the goal. At this point, a path to the goal is known from each state in the right half of the environment, even if the path is not optimal for all the states. After six iterations, an optimal policy is known for all the states in the environment.

7 Discussion

Planning with multi-time models converges to correct solutions significantly faster than planning at the level of primitive actions. There are two intuitive reasons that justify this result. First, the temporal abstraction achieved by the models enables the agent to reason at a higher level. Second, the knowledge captured in the models only depends only on the MDP underlying the environment and on the behavior itself. The behaviors and multi-time models can therefore be seen as an efficient means of transferring knowledge across different

reinforcement learning tasks, as long as the dynamics of the environment are preserved.

Theoretical results similar to the ones presented in this paper are also available for planning in optimal stopping tasks [1], as well as for a different regime of executing behaviors, which allows early termination. Further research will be devoted to integrating temporal and state abstraction, and to the issue of discovering useful behaviors.

Acknowledgments

The authors thank Amy McGovern, Andy Fagg, Leo Zelevinsky, Manfred Huber and Ron Parr for helpful discussions and comments contributing to this paper. This research was supported in part by NSF grant ECS-9511805 to Andrew G. Barto and Richard S. Sutton, and by AFOSR grant AFOSR-F49620-96-1-0254 to Andrew G. Barto and Richard S. Sutton. Doina Precup also acknowledges the support of the Fulbright foundation.

References

1. Dimitri P. Bertsekas. *Dynamic Programming: Deterministic and Stochastic Models*. Prentice Hall, Englewood Cliffs, NJ, 1987.
2. Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5:613–624, 1993.
3. Peter Dayan and Geoff E. Hinton. Feudal reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 5, pages 271–278, Cambridge, MA, 1993. MIT Press.
4. Thomas G. Dietterich. Hierarchical reinforcement learning with maxq value function decomposition. Technical report, Computer Science Department, Oregon State University, 1997.
5. Manfred Huber and Roderic A. Grupen. Learning to coordinate controllers - reinforcement learning on a control basis. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence, IJCAI-97*, San Francisco, CA, 1997. Morgan Kaufmann.
6. Leslie P. Kaelbling. Hierarchical learning in stochastic domains: Preliminary results. In *Proceedings of the Tenth International Conference on Machine Learning ICML'93*, pages 167–173, San Mateo, CA, 1993. Morgan Kaufmann.
7. Richard E. Korf. *Learning to Solve Problems by Searching for Macro-Operators*. Pitman Publishing Ltd, London, 1985.
8. John E. Laird, Paul S. Rosenbloom, and Allan Newell. Chunking in SOAR: The anatomy of a general learning mechanism. *Machine Learning*, 1:11–46, 1986.
9. Sridhar Mahadevan and Jonathan Connell. Automatic programming of behavior-based robots using reinforcement learning. *Artificial Intelligence*, 55(2-3):311–365, 1992.
10. Amy McGovern, Richard S. Sutton, and Andrew H. Fagg. Roles of macro-actions in accelerating reinforcement learning. In *Grace Hopper Celebration of Women in Computing*, pages 13–18, 1997.
11. Andrew W. Moore and Chris G. Atkeson. Prioritized sweeping: Reinforcement learning with less data and less real time. *Machine Learning*, 13:103–130, 1993.

12. Ronald Parr and Stuart Russell. Reinforcement learning with hierarchies of machines. In *Advances in Neural Information Processing Systems* volume 10, Cambridge, MA, 1997. MIT Press. To appear.
13. Jing Peng and John Williams. Efficient learning and planning within the Dyna framework. *Adaptive Behavior*, 4:323–334, 1993.
14. Doina Precup and Richard S. Sutton. Multi-Time models for temporally abstract planning. In *Advances in Neural Information Processing Systems* volume 10, Cambridge, MA, 1997. MIT Press. To appear.
15. Earl D. Sacerdoti. *A Structure for Plans and Behavior*. Elsevier, North-Holland, NY, 1977.
16. Satinder P. Singh. Scaling reinforcement learning by learning variable temporal resolution models. In *Proceedings of the Ninth International Conference on Machine Learning ICML'92*, pages 202–207, San Mateo, CA, 1992. Morgan Kaufmann.
17. Richard S. Sutton. Integrating architectures for learning, planning, and reacting based on approximating dynamic programming. In *Proceedings of the Seventh International Conference on Machine Learning ICML'90*, pages 216–224, San Mateo, CA, 1990. Morgan Kaufmann.
18. Richard S. Sutton. TD models: Modeling the world as a mixture of time scales. In *Proceedings of the Twelfth International Conference on Machine Learning ICML'95*, pages 531–539, San Mateo, CA, 1995. Morgan Kaufmann.
19. Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning. An Introduction*. MIT Press, Cambridge, MA, 1998.
20. Richard S. Sutton and Brian Pinette. The learning of world models by connectionist networks. In *Proceedings of the Seventh Annual Conference of the Cognitive Science Society*, pages 54–64, 1985.
21. Christopher J. C. H. Watkins. Learning with delayed rewards PhD Thesis, Cambridge University, 1989.