

# Notes on Reinforcement Learning

Paulo Eduardo Rauber

2014

## 1 Introduction

Reinforcement learning is the study of agents that act in an environment with the goal of maximizing cumulative reward signals. The agent is not told which actions to take, but discovers which actions yield most rewards by trying them. Its actions also may affect not only the immediate rewards but rewards for the next situations. There are several methods to solve the reinforcement learning problem.

One of the challenges in reinforcement learning is the exploration versus exploitation problem, which is the trade-off between obtaining rewards from perceived safe options against exploring other possibilities which may be advantageous.

A reinforcement learning problem can be divided into four subelements: policy, reward function, value function and a model.

A policy specifies how an agent behaves in a given scenario. It can be seen as a mapping from the perceived environmental state to actions.

A reward function defines the goal in a reinforcement learning problem. It maps each state of the environment to a numerical reward, indicating the intrinsic desirability of that state. An agent's sole objective is to maximize the total reward it receives on the long run. The reward function, therefore, should not be alterable by the agent.

A value function maps a state to an estimate of the total reward an agent can expect to accumulate over the future starting from that state. Even though a given state may have a small reward, it may be a prerequisite for achieving states with higher rewards.

A model is an optional element of a reinforcement learning system. A model describes the behavior of the environment. For example, it can be used to predict the results of actions (both states and rewards). It is possible to have agents that learn by trial and error, construct a model of the environment and use the model for planning.

## 2 Evaluative Feedback

The  $n$ -armed bandit problem is very useful to introduce a number of basic learning methods. Suppose an agent is faced with a choice among  $n$  different actions. After each action, a numerical reward is sampled based on a stationary probability distribution associated to that particular action. The agent's goal is to maximize the total reward after a number of iterations.

In this case, the value of an action is the expected reward that follows it. If an agent knew the value of every action, it would suffice to choose the action with the highest value.

The agent can maintain an estimate to the value of each action. The action with the highest estimated value is called greedy. Choosing the greedy action (exploitation) may guarantee higher rewards in the short run, while choosing non-greedy actions may lead to better outcomes later.

Let  $Q_t(a)$  denote the estimated value action  $a$  at iteration  $t$ . We also let  $Q^*(a)$  denote the value of  $a$ . A natural way to estimate the value of action  $a$  is to average the rewards already received. Concretely, if  $k_a$  denotes the number of times action  $a$  was chosen, we can compute  $Q_t(a)$  as:

$$Q_t(a) = \frac{r_1 + \dots + r_{k_a}}{k_a}.$$

By the law of large numbers, when  $k_a \rightarrow \infty$ ,  $Q_t(a) \rightarrow Q^*(a)$ .

A very simple alternative to always choosing the action with the highest estimated value is the  $\epsilon$ -greedy selection rule. This rule chooses the greedy action with probability  $1 - \epsilon$  and a random action (uniformly) with probability  $\epsilon$ . Using this rule, it is extremely important that ties be broken at random. After infinite iterations, every action will be sampled infinite times. Therefore,  $Q_t(a)$  converges to  $Q^*(a)$  and the chance of choosing the optimal action converges to  $1 - \epsilon$ .

Another possibility is to weight the probability of choosing an action with its current value estimate. Let  $\pi_t(a)$  denote the probability of choosing action  $a$  at iteration  $t$  and  $n$  denote the total number of actions. The following equation is called the softmax action selection rule:

$$\pi_t(a) = \frac{e^{\frac{Q_t(a)}{\tau}}}{\sum_{b=1}^n e^{\frac{Q_t(b)}{\tau}}}.$$

In this case, the probability of an action being selected is in proportion to its estimated value. The parameter  $\tau$  is called temperature. When  $\tau$  is large, the probabilities are more uniform. When  $\tau$  is near zero, the rule is greedier. In the case of two actions, this function becomes the sigmoidal function.

Fixing an action  $a$  and letting  $k$  denote the number of times it has been chosen, the equation given for estimating action values as an average of observed rewards can be rewritten as:

$$Q_{k+1} = Q_k + \frac{1}{k+1}(r_{k+1} - Q_k).$$

If the problem is non-stationary, however, it may be better to estimate the value of an action as:

$$Q_{k+1} = Q_k + \alpha(r_{k+1} - Q_k).$$

Where  $0 < \alpha \leq 1$  is called the step-size parameter. In this case, it can be shown that the value of  $Q_k$  is given by:

$$Q_k = (1 - \alpha)^k Q_0 + \sum_{i=1}^k \alpha(1 - \alpha)^{k-i} r_i.$$

Since it can be shown that  $(1 - \alpha)^k + \sum_{i=1}^k \alpha(1 - \alpha)^{k-i} = 1$ , this update rule is a weighted average of the initial estimate and the observed rewards. The recently observed rewards have more importance than older rewards. In fact, the weight decreases exponentially for older rewards.

If we denote by  $\alpha_k(a)$  the step-size parameter used to process the reward given after the  $k$ -th selection of action  $a$ , the choice of  $\alpha_k(a) = \frac{1}{k}$  leads to the sample average method, while  $\alpha_k(a) = \alpha$  leads to the constant step-size method.

There are two conditions that guarantee the convergence of the value estimate to the value of an action  $a$ :

$$\sum_{k=1}^{\infty} \alpha_k(a) = \infty,$$

and

$$\sum_{k=1}^{\infty} \alpha_k^2(a) < \infty.$$

Therefore, convergence is guaranteed for the sample average method, but not for the constant step size method. However, constant step sizes are advantageous if the problem is non-stationary, meaning that the value of actions may change during the iterations.

There is a simple trick that can be employed to make an agent explore in early iterations. If the initial value estimates are set optimistically, i.e., significantly higher than the value of the optimal action, even a greedy agent will explore considerably in the initial iterations until the value estimates are lowered.

Suppose that an observed variable (the state of the system) contained information about the next reward for a given action. It would be possible to achieve better results than using the strategies described in this chapter. This motivates the full reinforcement learning problem.

### 3 The Reinforcement Learning Problem

A reinforcement learning agent interacts with the environment at each of a sequence of discrete time steps  $t = 0, 1, 3, \dots$ . At each time step  $t$ , the agent receives some representation of the environment's state  $s_t \in \mathcal{S}$ , where  $\mathcal{S}$  is the set of possible states. The agent then selects an action  $a_t \in \mathcal{A}(s_t)$ , where  $\mathcal{A}(s_t)$  is the set of available actions in state  $s_t$ . One time step later, the agent receives a reward  $r_{t+1}$  and a new state  $s_{t+1}$ . The agent implements a policy  $\pi_t(s, a)$ : the probability that  $a_t = a$  given that  $s_t = s$ .

Modeling a problem as a reinforcement learning problem is challenging. Particularly, the boundary between agent and environment is sometimes not clear. A good strategy may be to abstract as much of the problem as

possible and let the agent focus on learning that which is hard to explicitly program. For instance, the action primitives for a robot may be pre-programmed tasks instead of mechanical control. Anything that cannot be controlled directly by the agent should be considered part of the environment, including rewards.

At each time step, the agent receives a reward  $r_t \in \mathbb{R}$ . The objective of the agent is to maximize the total amount of reward it receives on the long run. The choice of a reward function is critical in a reinforcement learning problem. It is particularly important not to impart previous knowledge into the reward function about how to solve a problem, but only what the agent should achieve. Previous knowledge should be imparted into state value estimates, which will be discussed later.

Let  $r_{t+1}, r_{t+2}, r_{t+3}, \dots$  denote the sequence of rewards received after time step  $t$ . A function  $R_t$  from a sequence of rewards to a real number is called a return function, and its value for a specific sequence is called return. A reinforcement learning agent maximizes the expected return.

In a simple formulation, the return function could be defined as  $R_t = r_{t+1} + r_{t+2} + \dots + r_T$ , where  $T$  is the final time step. When  $T$  is finite, we call the problem episodic. An episode ends when the agent reaches a terminal state. After that, the time is reset and the agent samples from a distribution of initial states.

However, when  $T = \infty$ , the simple formulation of  $R_t$  may be problematic, since the return could be infinite. This is called a continuous problem. In this case, it is better to consider the discounted return:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^T \gamma^k r_{t+k+1}, \quad (1)$$

where  $0 \leq \gamma < 1$  is called the discount rate. The discount rate determines the present value of future rewards. A reward received  $k$  time steps into the future is only worth  $\gamma^{k-1}$  times what it would be worth if it were received immediately. In this formulation, as long as the sequence  $\{r_k\}$  is bounded, the (discounted) return is finite. The choice of  $\gamma$  is important to determine how future-oriented is the agent.

In both continuous and episodic problems, the expected return for each state can be estimated by averaging the returns observed in several instances of the problem (which may be approximations in the continuous case).

The choice of a reward function is dependent on the definition of return. For instance, consider an agent which tries to escape a maze. In the episodic case, the reward could be  $-1$  for every move the agent makes until it reaches the exit. In this case, after  $T$  time steps, the reward for the agent would be  $-T$ , and maximizing the expected return would imply in finding the exit early. In the discounted case, the reward could be 1 whenever the agent found the final step and 0 otherwise. In this case, the expected return would be  $\gamma^k$ , and minimizing  $k$  would be implied. However, if the same formulation were used in the episodic, undiscounted case, the agent would have no incentive to find the exit early.

An episodic problem can be modeled as a continuous problem by simply creating a state which yields no rewards and transitions only to itself. In fact, if we let  $\gamma = 1$ , we can use the same formulation given in Eq. 1 for continuous tasks, as long as  $T = \infty$  and  $\gamma = 1$  are never both true.

In a reinforcement learning problem the state represents whatever information is available to the agent. States may be immediate sensations or highly processed versions of data and are often represented by vectors. States don't need to represent the environment perfectly or even all that could be useful in making marginally better decisions.

However, a state should respect the Markov property as well as possible. A state that summarizes all that is relevant for decision making that was observed in past states is said to have the Markov property. More precisely, given the sequence of states, action and rewards  $s_t, a_t, r_t, \dots, r_1, s_0, a_0$  (called history) observed by an agent, the state signal is said to have the Markov property if:

$$P(s_{t+1} = s', \mathbf{r}_{t+1} = r' | s_t, a_t, r_t, \dots, r_1, s_0, a_0) = P(s_{t+1} = s', \mathbf{r}_{t+1} = r' | s_t, a_t),$$

for all  $s', r'$  and histories. We call the probability distribution function on the right side the one-step dynamics of the reinforcement learning problem. It can be shown that the one-step dynamics can be used to derive the probability of all future states and expected rewards only from the initial state.

A reinforcement learning problem that satisfies the Markov property is called a Markov decision process, often called *MDP*. In particular, if the state and action sets are finite, the problem is called a finite MDP.

A Markov decision process is defined by its state and action sets and by the one-step dynamics of the environment. Given any state  $s$  and action  $a$ , the probability  $\mathcal{P}_{ss'}^a$  of each next state  $s'$  is given by:

$$\mathcal{P}_{ss'}^a = P(s_{t+1} = s' | s_t = s, \mathbf{a}_t = a) = \sum_{r'} P(s_{t+1} = s', \mathbf{r}_{t+1} = r' | s_t = s, \mathbf{a}_t = a).$$

Given any state  $s$ , action  $a$  and next state  $s'$ , the expected value of the next reward is given by:

$$\mathcal{R}_{ss'}^a = \mathbb{E}[\mathbf{r}_{t+1} | \mathbf{s}_t = s, \mathbf{a}_t = a, \mathbf{s}_{t+1} = s'] = \frac{1}{\mathcal{P}_{ss'}^a} \sum_{r'} r' P(\mathbf{s}_{t+1} = s', \mathbf{r}_{t+1} = r' | \mathbf{a}_t = a, \mathbf{s}_t = s).$$

The value  $V^\pi(s)$  of a state  $s$  is defined as the expected return of starting in state  $s$  and following the policy  $\pi$ . The policy assigns a probability  $\pi(s, a)$  to taking an action  $a$  when in a state  $s$ . More precisely:

$$V^\pi(s) = \mathbb{E}_\pi[\mathbf{R}_t | \mathbf{s}_t = s] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k \mathbf{r}_{t+k+1} | \mathbf{s}_t = s\right].$$

Similarly, the value of taking an action  $a$  when in state  $s$  and afterwards following the policy  $\pi$  is denoted by  $Q^\pi(s, a)$  and defined as:

$$Q^\pi(s, a) = \mathbb{E}_\pi[\mathbf{R}_t | \mathbf{s}_t = s, \mathbf{a}_t = a] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k \mathbf{r}_{t+k+1} | \mathbf{s}_t = s, \mathbf{a}_t = a\right].$$

The function  $V^\pi$  is called state value function for policy  $\pi$  and the function  $Q^\pi$  the action value function for policy  $\pi$ . These two functions can be estimated, for a given policy, by observing several instances of the reinforcement learning problem.

Given a policy, another way of stating the state and action value functions is:

$$V^\pi(s) = \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')], \quad (2)$$

and

$$Q^\pi(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a')].$$

These two equations are called the Bellman equations and can be interpreted intuitively. Note that, given a policy and the one-step dynamics of the problem, each of the equations defines a system of  $|\mathcal{S}|$  linear equations and variables. For any given policy, it is possible to find the *unique* expected return for any given state and action at a state.

It is possible to prove, in the case of a continuous discounted return problem, that adding a constant to all the rewards only adds another constant to the value of all states. However, in the case of an episodic task, this is not generally true, due to the variable number of time steps that may occur in each episode.

There are two alternative ways of writing the state and action value functions in terms of each other:

$$V^\pi(s) = \sum_a \pi(s, a) Q^\pi(s, a),$$

and

$$Q^\pi(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')].$$

Solving a reinforcement learning problem consists in finding a policy that has high expected return. Let  $\pi > \pi'$  if and only if  $V^\pi(s) \geq V^{\pi'}(s)$  for all  $s \in \mathcal{S}$ . The optimal policy  $\pi^*$  is such that  $\pi^* \geq \pi$  for any policy  $\pi$ . It is possible to show that this policy exists but is not necessarily unique. We define the optimal state value function  $V^*$  and the optimal action value function  $Q^*$  as:

$$V^*(s) = \max_{\pi} V^\pi(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^*(s')],$$

and

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a) = \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma \max_{a'} Q^*(s', a')],$$

for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ . These are called the Bellman optimality equations. Each equation gives a system of nonlinear equations that can be used to find  $V^*(s)$  or  $Q^*(s, a)$  for any state  $s$  and action  $a$ . In practice, however, solving nonlinear systems of equations in problems with many states may be unfeasible.

An optimal policy  $\pi^*$  can be found given  $V^*$  or  $Q^*$ . In the case of  $V^*$ , it suffices to choose one of the actions  $a$  that maximizes the right hand side of its Bellman optimality equation with uniform probability. In the case of  $Q^*$ , it suffices to choose uniformly the  $a$  such that  $Q^*(s, a)$  is maximal.

A large amount of memory may be required to build up approximations of value functions, policies and models. Therefore, in many practical problems, these functions must be approximated by a parameterized function.

## 4 Dynamic Programming

The term dynamic programming refers to a collection of reinforcement learning algorithms that can be used to compute optimal policies given a perfect model of the environment (one step dynamics). It is assumed that the set of states and actions are finite.

Policy evaluation is the name given to computing the state value function  $V^\pi$  given an arbitrary policy  $\pi$ . As already stated, it is possible to solve the system of linear equations with variables  $V^\pi(s)$  (Eq. 2). However, it is also possible to create a sequence of approximations to  $V^\pi$  by the following rule:

$$V_{k+1}(s) = \sum_a \pi(s, a) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s)],$$

for all  $s \in \mathcal{S}$ . This sequence is guaranteed to converge to  $V^\pi(s)$  for all  $s \in \mathcal{S}$ , since  $V_k = V_{k+1}$  for some  $k$  implies on the solution to the Bellman equations. The initial value for each state can be arbitrary. This method is called iterative policy evaluation.

In a given iteration of policy evaluation, instead of computing the new estimate for each state using the estimate from the previous iteration, it is also valid to consider the latest estimate available for each state. This variation is classified as an in-place algorithm.

A deterministic policy  $\pi$  is one such that, for all  $s \in \mathcal{S}$ ,  $\pi(s, a) = 1$  for some  $a$  and  $\pi(s, b) = 0$  for all  $b \neq a$ . Thus, in the case of deterministic policies, it is simpler to consider  $\pi$  as a function from states to actions.

Let  $\pi$  and  $\pi'$  be any pair of deterministic policies such that, for all  $s \in \mathcal{S}$ ,  $Q^\pi(s, \pi'(s)) \geq V^\pi(s)$ . The policy improvement theorem guarantees that  $V^{\pi'}(s) \geq V^\pi(s)$  for all  $s \in \mathcal{S}$ . This gives a natural way of improving a given policy  $\pi$  to a better policy  $\pi'$  by the following rule:

$$\pi'(s) = \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s)].$$

A sequence of evaluations and improvements can be applied to an arbitrary initial policy. Once a policy can no longer be improved, it is guaranteed to be optimal by the Bellman optimality equations. This technique for finding an optimal policy is called policy iteration.

A very simple and efficient alternative to policy iteration is called value iteration. This technique successively improves the estimates for the value of each state  $s \in \mathcal{S}$ , beginning with an arbitrary value  $V_0(s)$ , by the following rule:

$$V_{k+1} = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')].$$

It can be shown that the sequence  $\{V_k\}$  converges to  $V^*$ . Algorithm 4 describes how an optimal policy can be obtained by value iteration.

---

### Algorithm 1 Value iteration

---

**Input:** set of states  $\mathcal{S}$ , one step dynamics functions  $\mathcal{P}$  and  $\mathcal{R}$ , discount factor  $\gamma$ .

**Output:** optimal deterministic policy  $\pi$ .

```

1: for each  $s \in \mathcal{S}$  do
2:    $V(s) \leftarrow 0$ 
3: end for
4: repeat
5:    $\Delta \leftarrow 0$ 
6:   for each  $s \in \mathcal{S}$  do
7:      $v \leftarrow V(s)$ 
8:      $V(s) \leftarrow \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$ 
9:      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
10:  end for
11: until  $\Delta < \theta$ 
12: for each  $s \in \mathcal{S}$  do
13:    $\pi(s) = \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s)]$ 
14: end for
```

---

Asynchronous algorithms can also be used to update the estimates for  $V$ .

There are two main disadvantages of both policy iteration and value iteration: they required expensive sweeps over the state space, which may be prohibitively large, and a complete model of the environment.

## 5 Monte Carlo Methods

Monte Carlo methods can be used to solve the reinforcement learning problem without a model (one-step dynamics) for the environment. The agent learns by averaging observed returns. Even if it is possible to construct a model of the environment, it is often easier to let the agent learn in this way.

These methods also work by successive policy evaluation and improvement. The most simple way to estimate the value of a state is to average the observed returns after a number of episodes in which the state appears. If only the first occurrence of a state  $s$  in an episode is considered to improve the estimate for  $V(s)$ , then the method is called first-visit. Otherwise, when every occurrence of a state  $s$  is considered, the method is called every-visit. In both cases, when the number of visits to  $s$  tends to infinity, this kind of estimate tends to  $V(s)$ .

Since a model is not available for Monte Carlo methods, it is useful to estimate action values instead of state values. Evaluating a policy  $\pi$  by a Monte Carlo method consists on experiencing several episode and averaging the returns that follow every possible state action pair  $(s, a)$  to obtain an estimate for  $Q(s, a)$ .

As in the previous section, we are interested in improving a policy  $\pi$ . In this case, we do this episode-by-episode. However, unlike in the previous section, the policies must be  $\epsilon$ -soft, since we need to guarantee that the method improves its estimate for the action values. Algorithm 5 illustrates a Monte Carlo based algorithm to solve the reinforcement learning problem. It is assumed that there is no discounting, since the length of each episode *should* be finite.

---

**Algorithm 2** Monte Carlo control algorithm

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , probability of choosing random action  $\epsilon$ .

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```
1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:      $n(s, a) \leftarrow 0$ 
5:   end for
6: end for
7: for each  $i$  in  $\{1, \dots, N\}$  do
8:   Experience a new episode  $e$  following an  $\epsilon$ -greedy policy based on  $Q$ .
9:   for each state-action pair  $(s, a)$  in the episode  $e$  do
10:     $R \leftarrow$  return following  $(s, a)$  in the episode  $e$ .
11:     $n(s, a) \leftarrow n(s, a) + 1$ 
12:     $Q(s, a) \leftarrow Q(s, a) + \frac{1}{n(s, a)}[R - Q(s, a)]$ 
13:   end for
14: end for
15: for each state  $s \in \mathcal{S}$  do
16:    $\pi(s) \leftarrow \arg \max_a Q(s, a)$ 
17: end for
```

---

A method to solve the reinforcement learning problem that uses the value of other states to compute the value of a state  $s$  is said to bootstrap. Since Monte Carlo methods do not bootstrap, in contrast with dynamic programming, it is possible to evaluate the value of each state separately, which can be necessary in some situations. It is also possible to evaluate a policy while following another, which is described in detail in [1].



## 6 Temporal-Difference Learning

Temporal-difference learning methods are also capable of learning from experience without a model of the environment. However, differently from Monte Carlo methods, they bootstrap. Intuitively, this allows them to extract more information from their experiences.

Instead of waiting for the end of an episode to update their state or action values, a TD agent updates at every time step. The simplest TD method is called  $TD(0)$ , and updates the state value estimates using the following rule:

$$V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)], \quad (3)$$

where  $\alpha$  is the learning rate and  $\gamma$  is a discount parameter. This rule updates the value of a state based on the reward after one time-step and the estimated value of the next state.

Sarsa is an on-policy method to solve the control problem: finding an optimal policy by interacting with the environment. The idea is to adapt Equation 3 to estimate action values, which involves using the tuple  $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$ , which gives the method's name. As with other control algorithms, the starting point is an arbitrary  $\epsilon$ -soft policy, which is used to estimate action values and improve the policy, which is done after each time step. Sarsa is exemplified by algorithm 6, which uses  $\epsilon$ -greedy policies.

---

### Algorithm 3 Sarsa control algorithm

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ .

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```

1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:   end for
5: end for
6: for each  $i$  in  $\{1, \dots, N\}$  do
7:    $s \leftarrow$  initial state for episode  $i$ 
8:   Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
9:   while state  $s$  is not terminal do
10:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
11:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
12:    Select action  $a'$  for state  $s'$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
13:     $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$ 
14:     $s \leftarrow s'$ 
15:     $a \leftarrow a'$ 
16:   end while
17: end for
18: for each state  $s \in \mathcal{S}$  do
19:    $\pi(s) \leftarrow \arg \max_a Q(s, a)$ 
20: end for
```

---

Another advantage of TD methods over Monte Carlo methods is that episodes don't need to be finite. Also, while Monte Carlo methods minimized the mean squared-error on the training set, TD methods find the maximum likelihood model for the underlying Markov process, which achieves better generalization. Finding the maximum likelihood model requires that a set of episodes be experienced repeatedly by the algorithms until convergence of the action value function  $Q$ , in a process that is called batch-updating.

An alternative to Sarsa is  $Q$ -Learning, an off-policy control method. This means that  $Q$ -learning learns the action values for a policy that it does not use to act on the environment. The main difference to Sarsa is the update rule, which takes the form:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)].$$

Instead of choosing  $a'$  using the current policy, which may be  $\epsilon$ -soft for learning purposes, this update rule considers the best possible action according to the current estimate. This avoids the problem of misattributing an inherent low value for a state in which a bad action choice was made. For the same reason,  $Q$ -learning may

converge faster but have lower efficacy during the learning phase (called online efficacy). Algorithm 6 exemplifies a particular implementation of  $Q$ -learning.

---

**Algorithm 4**  $Q$ -learning control algorithm

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ .

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```

1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:   end for
5: end for
6: for each  $i$  in  $\{1, \dots, N\}$  do
7:    $s \leftarrow$  initial state for episode  $i$ 
8:   while state  $s$  is not terminal do
9:     Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
10:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
11:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
12:     $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}(s')} Q(s', a') - Q(s, a)]$ 
13:     $s \leftarrow s'$ 
14:   end while
15: end for
16: for each state  $s \in \mathcal{S}$  do
17:    $\pi(s) \leftarrow \arg \max_a Q(s, a)$ 
18: end for

```

---

Alternatively, the update rule for  $Q$ -learning can be taken as:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \sum_a \pi(s_t, a) Q(s_{t+1}, a) - Q(s_t, a_t)].$$

In this case, the probability of choosing each action is taken into account, which may improve its online efficacy.

It is often possible to design a more effective learning algorithm when some part of the environment dynamics is known. For instance, if the state following a state action pair is always known, it is possible to base policies on the value of possible successor states, by learning what is called an afterstate value function. In that case, the control algorithm can learn the state value function  $V$  instead of  $Q$ .

## 7 Eligibility Traces

The  $n$ -step return at time  $t$  is defined as:

$$R_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n V_t(s_{t+n}).$$

This return can be used to update state (action) values:

$$V(s_t) \leftarrow V(s_t) + \alpha [R_t^{(n)} - V(s_t)].$$

When  $n = 1$ , the update rule above is equivalent to the temporal difference learning rule. When  $n \rightarrow \infty$ , this rule is equivalent to the rule used in Monte Carlo methods.

It is also possible to average  $R_t^{(i)}$  for different  $i$ . One particularly important weighted average of variable step returns is  $R_t^\lambda$ , called  $\lambda$ -return, which is defined by:

$$R_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} R_t^{(n)}.$$

In this way, the weights given to each  $n$ -step return sum to 1. The 1-step return has the highest weight, which decreases by a factor of  $\lambda$  for each following return. When  $\lambda = 0$ , the definition is once again equivalent to temporal difference learning with one step return. When  $\lambda \rightarrow 1$ , the method is equivalent to Monte Carlo. This gives a somewhat simple way to control the balance between state (action) value estimates and looking ahead for rewards. This idea has been shown to be very effective in practice for good choices of  $\lambda$ .

In practice, the update rule based on  $R_t^\lambda$  is efficiently implemented by using eligibility traces. The eligibility trace  $e_t$  for a state  $s$  at time  $t$  is defined as:

$$e_t(s) = \begin{cases} 0, & t = 0 \\ \gamma \lambda e_{t-1}(s) + 1, & s = s_t \text{ and } t \neq 0 \\ \gamma \lambda e_{t-1}(s), & s \neq s_t \text{ and } t \neq 0 \end{cases}$$

The following update rule should be applied to *every* state (action) value:

$$V(s) \leftarrow V(s) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] e_t(s). \quad (4)$$

Intuitively, the observed error for the estimate of the value of the current state  $s_t$  changes the estimate of  $s$  by a factor of  $\alpha e_t(s)$ , which is proportional to how soon state  $s$  was last visited. In this context,  $\lambda$  is called the trace decay factor.

The control algorithms presented in the previous section can be adapted to use eligibility traces as described by 7 and 7. Note that the algorithm  $Q(\lambda)$  needs to erase its eligibility traces once a non-greedy action is taken because of its off-policy evaluation.

The eligibility traces presented in this section are also called accumulating traces. Replacing traces are an alternative definition that may be better suited to some applications, particularly when the agent is likely to take the same action at a given state several times, causing a detrimental accumulation of its eligibility trace. The replacing trace  $e_t$  for a state  $s$  at time  $t$  is defined as:

$$e_t(s) = \begin{cases} 0, & t = 0 \\ 1, & s = s_t \text{ and } t \neq 0 \\ \gamma \lambda e_{t-1}(s), & s \neq s_t \text{ and } t \neq 0 \end{cases}$$

---

**Algorithm 5** Sarsa( $\lambda$ ) control algorithm

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ , trace decay  $\lambda$

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```
1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:      $e(s, a) \leftarrow 0$ 
5:   end for
6: end for
7: for each  $i$  in  $\{1, \dots, N\}$  do
8:   Set the eligibility trace  $e$  to zero for every valid state action pair.
9:    $s \leftarrow$  initial state for episode  $i$ 
10:  Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
11:  while state  $s$  is not terminal do
12:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
13:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
14:    Select action  $a'$  for state  $s'$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
15:     $e(s, a) \leftarrow e(s, a) + 1$ 
16:    for each  $s'' \in \mathcal{S}$  do
17:      for each action  $a'' \in \mathcal{A}(s'')$  do
18:         $Q(s'', a'') \leftarrow Q(s'', a'') + \alpha[r + \gamma Q(s', a') - Q(s, a)]e(s'', a'')$ 
19:         $e(s'', a'') \leftarrow \gamma\lambda e(s'', a'')$ 
20:      end for
21:    end for
22:     $s \leftarrow s'$ 
23:     $a \leftarrow a'$ 
24:  end while
25: end for
26: for each state  $s \in \mathcal{S}$  do
27:    $\pi(s) \leftarrow \arg \max_a Q(s, a)$ 
28: end for
```

---

---

**Algorithm 6**  $Q(\lambda)$  control algorithm

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ , trace decay  $\lambda$

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```
1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:      $e(s, a) \leftarrow 0$ 
5:   end for
6: end for
7: for each  $i$  in  $\{1, \dots, N\}$  do
8:   Set the eligibility trace  $e$  to zero for every valid state action pair.
9:    $s \leftarrow$  initial state for episode  $i$ 
10:  Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
11:  while state  $s$  is not terminal do
12:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
13:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
14:    Select action  $a'$  for state  $s'$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
15:    Select the greedy action  $a^*$  for state  $s'$  according to  $Q$ , preferring  $a'$  when tied.
16:     $e(s, a) \leftarrow e(s, a) + 1$ 
17:    for each  $s'' \in \mathcal{S}$  do
18:      for each action  $a'' \in \mathcal{A}(s'')$  do
19:         $Q(s'', a'') \leftarrow Q(s'', a'') + \alpha[r + \gamma Q(s', a^*) - Q(s, a)]e(s'', a'')$ 
20:        if  $a^* = a'$  then
21:           $e(s'', a'') \leftarrow \gamma \lambda e(s'', a'')$ 
22:        else
23:           $e(s'', a'') \leftarrow 0$ 
24:        end if
25:      end for
26:    end for
27:     $s \leftarrow s'$ 
28:     $a \leftarrow a'$ 
29:  end while
30: end for
31: for each state  $s \in \mathcal{S}$  do
32:    $\pi(s) \leftarrow \arg \max_a Q(s, a)$ 
33: end for
```

---

## 8 Generalization and Function Approximation

There are situations in which it is inadvisable to store state or action values in arrays. In large state spaces, some spaces may be seen rarely. In these cases, the state (action) values must be generalized across states. The kind of generalization discussed in this section is based on function approximation, which is studied extensively in machine learning.

The state value function (at time  $t$ )  $V_t$  can be approximated by a function  $f$  such that  $f(\vec{\theta}_t, s_t) = V_t(s_t)$ , where  $\vec{\theta}_t$  is a vector in  $\mathbb{R}^m$  and  $s_t$  is a state. In this case, the value  $V(s_t)$  should depend only on the parameters  $\vec{\theta}_t$ , and the problem is then formulated as finding the best  $\vec{\theta}_t$  to fit the real value function. In this case, changing the value of  $\vec{\theta}_t$  typically changes the value estimates for several states.

Several methods already presented shifted the value  $V_t(s_t)$  of a state  $s_t$  towards a new value  $V_{t+1}(s_t)$ . In the case of function approximation, the new value may be considered as new sample for supervised learning. Concretely, the approximate value function  $V_t$  can be learned from pairs  $(s_i, v_i)$ , for  $i \leq t$ , where  $s_i$  is a state and  $v_i$  is its observed value. In the case of  $TD(\lambda)$ , for instance,  $v_i = R_i^\lambda$ .

Using this formulation, any regression method may be used to approximate  $V_t$ . This section will focus on linear regression by gradient descent, mainly due to its simplicity.

A common performance measure for evaluating function approximation is the mean squared error:

$$\text{MSE}(\vec{\theta}_t) = \sum_{s \in S} P(s) [V^\pi(s) - V_t(s)]^2.$$

In the equation above,  $S$  is the set of states,  $P$  is a probability distribution over states,  $V^\pi(s)$  is the value of state  $s$  under the policy  $\pi$  and  $V_t(s) = f(\vec{\theta}_t, s)$  is the current value estimate for state  $s$ . In most cases, the distribution  $P$  should describe the frequency with which each state is encountered while the agent interacts with the environment under  $\pi$ .

A vector  $\vec{\theta}^*$  such that  $\text{MSE}(\vec{\theta}^*) \leq \text{MSE}(\vec{\theta})$  for every vector  $\vec{\theta}$  is called a global optimum. Finding a global optimum is rarely practical for complex functions. Several methods used in practice are guaranteed only to find local optima, which are global optima in a restricted domain of the function.

Gradient descent is a very widespread technique for finding local optima in smooth functions. Stochastic gradient descent is a variation that is particularly useful in reinforcement learning. Given an estimate for the parameter vector  $\vec{\theta}_t$  and a current state  $s_t$ , the estimate may be updated in the following way:

$$\vec{\theta}_{t+1} = \vec{\theta}_t - \frac{1}{2} \alpha \nabla_{\vec{\theta}_t} [V^\pi(s_t) - V_t(s_t)]^2.$$

In the equation above,  $\alpha$  is the learning rate and  $\nabla_{\vec{\theta}_t} [V^\pi(s_t) - V_t(s_t)]^2$  is the gradient of  $[V^\pi(s_t) - V_t(s_t)]^2$  with respect to  $\vec{\theta}_t$ . This gradient is the vector of partial derivatives of the error term (the squared term) with respect to each component of  $\vec{\theta}_t$ . Intuitively, at each time step, the estimate is changed in the direction that minimizes the error (squared term) observed for state  $s_t$ . The direction in which the error is lowered for each component of  $\vec{\theta}_t$ , for a particular state  $s_t$ , is given by this gradient.

The equation above can be rewritten in the following way:

$$\vec{\theta}_{t+1} = \vec{\theta}_t + \alpha [V^\pi(s_t) - V_t(s_t)] \nabla_{\vec{\theta}_t} V_t(s_t).$$

If  $\alpha$  decays with time in the same way as described in section 2, this method is guaranteed to converge to a local optimum in a finite number of steps.

Naturally, if  $V^\pi(s)$  were available for all observed states  $s$ , there would be no need for function approximation. Therefore, in practice, an estimate of  $V^\pi(s_t)$ , denoted  $v_t$ , should be used to update the current parameter vector:

$$\vec{\theta}_{t+1} = \vec{\theta}_t + \alpha [v_t - V_t(s_t)] \nabla_{\vec{\theta}_t} V_t(s_t).$$

Suppose the states are sampled according to policy  $\pi$ . Let  $R_t$  denote the return following each state  $s_t$ . Since the expected value of  $R_t$  is the true value of  $s_t$ ,  $R_t$  is called an unbiased estimate and gradient descent is guaranteed to converge to a locally optimum approximation of  $V^\pi(s_t)$  when  $v_t = R_t$ . Therefore, Monte Carlo policy evaluation has the same guarantee. Unfortunately, the same guarantee does not apply to methods based on  $TD(\lambda)$  for  $\lambda < 1$ . However, in practice, temporal difference methods are often more successful in policy evaluation.

The updating of  $\vec{\theta}_t$  using  $TD(\lambda)$  can be formulated in the following way:

$$\vec{\theta}_{t+1} = \vec{\theta}_t + \alpha [R_t^\lambda - V_t(s_t)] \nabla_{\vec{\theta}_t} V_t(s_t).$$

It can be shown that this is equivalent to the following update rule using eligibility traces:

$$\vec{\theta}_{t+1} = \vec{\theta}_t + \alpha \delta_t \vec{e}_t,$$

with

$$\delta_t = r_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t),$$

and

$$\vec{e}_t = \gamma \lambda \vec{e}_{t-1} + \nabla_{\vec{\theta}_t} V_t(s_t),$$

and  $\vec{e}_0 = \vec{0}$ .

In this case, there is a vector of eligibility traces with the same number of components as  $\vec{\theta}_t$ .

Algorithm 7 is a complete algorithm for approximating  $V^\pi$  using gradient descent. The function  $V$  is, implicitly, a function of  $\vec{\theta}$ .

---

**Algorithm 7** Gradient descent TD( $\lambda$ ) value estimation algorithm

---

**Input:** policy  $\pi$ , number of episodes  $N$ , learning rate  $\alpha$ , discount factor  $\gamma$ , trace decay  $\lambda$

**Output:** parameter vector  $\vec{\theta}$

```

1:  $\vec{\theta} \leftarrow \vec{0}$ 
2: for each  $i$  in  $\{1, \dots, N\}$  do
3:    $\vec{e} \leftarrow \vec{0}$ 
4:    $s \leftarrow$  initial state for episode  $i$ 
5:   while state  $s$  is not terminal do
6:      $a \leftarrow \pi(s)$ 
7:      $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
8:      $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
9:      $\delta \leftarrow r + \gamma V(s') - V(s)$ 
10:     $\vec{e} \leftarrow \gamma \lambda \vec{e} + \nabla_{\vec{\theta}} V(s)$ 
11:     $\vec{\theta} \leftarrow \vec{\theta} + \alpha \delta \vec{e}$ 
12:     $s \leftarrow s'$ 
13:   end while
14: end for
```

---

Linear functions on the parameter vector  $\vec{\theta}_t$  are widely used for function approximation mainly due to their simplicity. In this case, a state is often represented by a *feature vector*  $\vec{\phi}_s$ , which should be useful for generalization across states. Given a parameter vector  $\vec{\theta}_t$ , the value  $V_t(s)$  can be computed as:

$$V_t(s) = f(\vec{\theta}_t, s) = \sum_i \theta_t(i) \phi_s(i).$$

In this case, it is clear that  $\nabla_{\vec{\theta}_t} V_t(s) = \vec{\phi}(s)$ .

In the case of linear function approximation, there is only one local optimum and, therefore, only one global optimum

The choice of feature vectors is crucial to the success of any function approximation method. Intuitively, the features should correspond to natural features of the state, which should facilitate generalization. For instance, a mobile robot might represent its state as a combination of its position, velocity, remaining power etc. In the case of linear models, it might be necessary to create features that are the combinations of other features, since linear models are, in general, incapable of modeling interactions such as feature  $i$  being good only on the absence of feature  $j$ .

High dimensional state spaces often need to be simplified for use with linear models. There are several techniques to accomplish this. The states can be represented in feature vectors by discretizing the state space into overlapping hyperspheres (coarse coding) or by tiling into hypercubes (tile coding, which may combine several tilings with different offsets). States can also be represented by feature vectors of distance calculations to random points sampled in the state space (similar to coding by radial basis functions).

The extension of function approximation from value prediction to action value prediction is straightforward. In this case,  $Q$  should be approximated by a function  $f$  of  $\vec{\theta}_t$ ,  $s$  and  $a$ . Algorithms 8 and 9 illustrate Sarsa( $\lambda$ ) and Q( $\lambda$ ) based on linear function approximation.

---

**Algorithm 8** Sarsa( $\lambda$ ) control algorithm for linear function approximation

---

**Input:** feature vector  $\phi_{s,a}$  for all state-action pairs  $(s, a)$ , number of episodes  $N$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ , trace decay  $\lambda$

**Output:** parameter vector  $\vec{\theta}$

```
1:  $\vec{\theta} \leftarrow \vec{0}$ 
2: for each  $i$  in  $\{1, \dots, N\}$  do
3:    $\vec{e} \leftarrow \vec{0}$ 
4:    $s \leftarrow$  initial state for episode  $i$ 
5:   for each  $a \in A(s)$ : do
6:      $Q(a) \leftarrow \sum_i \theta(i) \phi_{s,a}(i)$ 
7:   end for
8:    $a \leftarrow \arg \max_{a \in A(s)} Q(a)$ 
9:   With probability  $\epsilon$ :  $a \leftarrow$  random action in  $A(s)$ 
10:  while state  $s$  is not terminal do
11:     $\vec{e} \leftarrow \gamma \lambda \vec{e} + \vec{\phi}_{(s,a)}$ 
12:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
13:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
14:     $\delta \leftarrow r - Q(a)$ 
15:    for each  $a \in A(s)$ : do
16:       $Q(a) \leftarrow \sum_i \theta(i) \phi_{s',a}(i)$ 
17:    end for
18:     $a' \leftarrow \arg \max_{a \in A(s')} Q(a)$ 
19:    With probability  $\epsilon$ :  $a' \leftarrow$  random action in  $A(s')$ 
20:     $\delta \leftarrow \delta + \gamma Q(a')$ 
21:     $\vec{\theta} \leftarrow \vec{\theta} + \alpha \delta \vec{e}$ 
22:     $s \leftarrow s'$ 
23:     $a \leftarrow a'$ 
24:  end while
25: end for
```

---



---

**Algorithm 9**  $Q(\lambda)$  control algorithm for linear function approximation

---

**Input:** feature vector  $\phi_{s,a}$  for all state-action pairs  $(s, a)$ , number of episodes  $N$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ , trace decay  $\lambda$

**Output:** parameter vector  $\vec{\theta}$

```
1:  $\vec{\theta} \leftarrow \vec{0}$ 
2: for each  $i$  in  $\{1, \dots, N\}$  do
3:    $\vec{e} \leftarrow \vec{0}$ 
4:    $s \leftarrow$  initial state for episode  $i$ 
5:   for each  $a \in A(s)$ : do
6:      $Q(a) \leftarrow \sum_i \theta(i) \phi_{s,a}(i)$ 
7:   end for
8:   while state  $s$  is not terminal do
9:     if with probability  $1 - \epsilon$ : then
10:       $a \leftarrow \arg \max_a Q(a)$ 
11:       $\vec{e} \leftarrow \gamma \lambda \vec{e} + \vec{\phi}_{(s,a)}$ 
12:    else
13:       $a \leftarrow$  random action in  $A(s)$ 
14:       $\vec{e} \leftarrow \vec{\phi}_{(s,a)}$ 
15:    end if
16:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
17:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
18:     $\delta \leftarrow r - Q(a)$ 
19:    for each  $a \in A(s')$ : do
20:       $Q(a) \leftarrow \sum_i \theta(i) \phi_{s',a}(i)$ 
21:    end for
22:     $a' \leftarrow \arg \max_{a \in A(s')} Q(a)$ 
23:     $\delta \leftarrow \delta + \gamma Q(a')$ 
24:     $\vec{\theta} \leftarrow \vec{\theta} + \alpha \delta \vec{e}$ 
25:     $s \leftarrow s'$ 
26:  end while
27: end for
```

---

## 9 Planning and Learning

Model free reinforcement learning methods, such as temporal difference learning, can be used together with model based methods, such as dynamic programming. In section 3, the one step dynamics of the environment was presented as a complete model, i.e., a mechanism that can be used by the agent to predict the consequence of its actions.

In the case of deterministic environments, where the outcome of an action at a given state is always the same, it is sufficient to store these outcomes in a table as the agent learns. In stochastic environments, it is necessary to learn distribution models that represent the probabilities of each possible outcome and expected rewards.

In this context, planning is the process where the agent simulates, using its internal model, the environment and updates its value estimates.

Algorithm 10 illustrates the combination of direct reinforcement learning, using one-step  $Q$ -learning, and model learning in a deterministic environment. This method could be enhanced by eligibility traces, but that was omitted for simplicity.

---

**Algorithm 10** Dyna-Q control algorithm for deterministic environments

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , number of planning steps per episode  $K$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ .

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```

1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:   end for
5: end for
6: for each  $i$  in  $\{1, \dots, N\}$  do
7:    $s \leftarrow$  initial state for episode  $i$ 
8:   while state  $s$  is not terminal do
9:     Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
10:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
11:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
12:     $M(s, a) \leftarrow s', r$ 
13:     $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}(s')} Q(s', a') - Q(s, a)]$ 
14:     $s_{\text{next}} \leftarrow s'$ 
15:    for each  $k$  in  $\{1, \dots, K\}$  do
16:       $s \leftarrow$  random previously observed state
17:       $a \leftarrow$  random action previously taken in state  $s$ 
18:       $s', r \leftarrow M(s, a)$ 
19:       $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}(s')} Q(s', a') - Q(s, a)]$ 
20:    end for
21:     $s \leftarrow s_{\text{next}}$ 
22:  end while
23: end for
24: for each state  $s \in \mathcal{S}$  do
25:    $\pi(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} Q(s, a)$ 
26: end for
```

---

The main advantage of using model learning is the possibility of updating the value of any state at a given time step. This includes states that are not even visited in a given episode. In practice, this may considerably reduce the necessity of interaction with the environment.

Algorithm 11 presents a variation of the previous algorithm for stochastic environments. The algorithm performs full backup, as in dynamic programming methods, using estimates of  $\mathcal{P}_{ss'}^a$  and  $\mathcal{R}_{ss'}^a$  learned during interaction with the environment.

However, in the case of non-stationary environments, meaning that the environment changes its behavior with time, a model may become incorrect. In this case, the agent may become stuck behaving suboptimally. Even using  $\epsilon$ -greedy action selection, some action sequences may become too unlikely. This is a classical instance of the exploitation versus exploration problem. There is a simple heuristic that can overcome this problem in planning agents. If the agent stores the time  $n$  since the last execution of action  $a$  at state  $s$ , the planning step may consider the reward for the state-action combination as  $r + \kappa\sqrt{n}$  for a small parameter  $\kappa$  and the reward estimate (given

---

**Algorithm 11** Dyna-Q control algorithm for stochastic environments

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , number of planning steps per episode  $K$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ .

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```
1: for each  $s \in \mathcal{S}$  do
2:   for each action  $a \in \mathcal{A}(s)$  do
3:      $Q(s, a) \leftarrow 0$ 
4:      $R(s, s', a) \leftarrow 0$ 
5:      $\mathcal{N}(s, a) \leftarrow 0$ 
6:      $\mathcal{N}_n(s, a, s') \leftarrow 0$ 
7:   end for
8: end for
9: for each  $i$  in  $\{1, \dots, N\}$  do
10:   $s \leftarrow$  initial state for episode  $i$ 
11:  while state  $s$  is not terminal do
12:    Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
13:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
14:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
15:     $\mathcal{N}(s, a) \leftarrow \mathcal{N}(s, a) + 1$ 
16:     $\mathcal{N}_n(s, a, s') \leftarrow \mathcal{N}_n(s, a, s') + 1$ 
17:     $R(s, s', a) \leftarrow R(s, s', a) + \frac{r - R(s, s', a)}{\mathcal{N}_n(s, a, s')}$ 
18:     $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}(s')} Q(s', a') - Q(s, a)]$ 
19:     $s_{\text{next}} \leftarrow s'$ 
20:    for each  $k$  in  $\{1, \dots, K\}$  do
21:       $s \leftarrow$  random previously observed state
22:       $a \leftarrow$  random action previously taken in state  $s$ 
23:       $Q(s, a) \leftarrow \sum_{s'} [\frac{\mathcal{N}_n(s, a, s')}{\mathcal{N}(s, a)}][R(s, a, s') + \gamma \max_{a'} Q(s', a')]$ 
24:    end for
25:     $s \leftarrow s_{\text{next}}$ 
26:  end while
27: end for
28: for each state  $s \in \mathcal{S}$  do
29:   $\pi(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} Q(s, a)$ 
30: end for
```

---

by the model)  $r$ . In this case, a whole sequence of exploratory actions may be undertaken. In practice, however, controlling  $\kappa$  may be challenging.

In algorithm 10, simulated transitions are sampled uniformly. However, this is usually not the best choice, even for the dynamic programming methods presented in section 4. A simple heuristic called prioritized sweeping may be employed to choose transitions for planning. This heuristic prefers to update values for transitions that go into states that had their values considerably changed in the current episode. A method based on this heuristic, for deterministic environments, is presented in algorithm 12.

---

**Algorithm 12** Prioritized sweeping Dyna-Q control algorithm for deterministic environments

---

**Input:** set of states  $\mathcal{S}$ , number of episodes  $N$ , number of planning steps per episode  $K$ , learning rate  $\alpha$ , probability of choosing random action  $\epsilon$ , discount factor  $\gamma$ , heap threshold  $h$ .

**Output:** deterministic policy  $\pi$ , optimal when  $N \rightarrow \infty$ .

```

1:  $\mathcal{H} \leftarrow$  empty maximum priority queue (max-heap)
2: for each  $s \in \mathcal{S}$  do
3:   for each action  $a \in \mathcal{A}(s)$  do
4:      $Q(s, a) \leftarrow 0$ 
5:   end for
6: end for
7: for each  $i$  in  $\{1, \dots, N\}$  do
8:    $s \leftarrow$  initial state for episode  $i$ 
9:   while state  $s$  is not terminal do
10:    Select action  $a$  for state  $s$  according to an  $\epsilon$ -greedy policy based on  $Q$ .
11:     $r \leftarrow$  observed reward for action  $a$  at state  $s$ 
12:     $s' \leftarrow$  observed next state for action  $a$  at state  $s$ 
13:     $M(s, a) = s', r$ 
14:     $p \leftarrow |r + \gamma \max_{a' \in \mathcal{A}(s')} Q(s', a') - Q(s, a)|$ 
15:    if  $p > h$  then
16:      Insert  $(s, a)$  into  $\mathcal{H}$  with priority  $p$ 
17:    end if
18:     $s_{\text{next}} \leftarrow s'$ 
19:    for each  $k$  in  $\{1, \dots, K\}$  or  $\text{empty}(\mathcal{H})$  do
20:       $s, a \leftarrow \text{first}(\mathcal{H})$ 
21:       $s', r \leftarrow M(s, a)$ 
22:       $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}(s')} Q(s', a') - Q(s, a)]$ 
23:      for each  $\bar{s}, \bar{a}$  predicted to lead to  $s$  do
24:         $\bar{r} \leftarrow$  predicted reward for transition  $\bar{s}, \bar{a}$ 
25:         $p \leftarrow |\bar{r} + \gamma \max_{a \in \mathcal{A}(s)} Q(s, a) - Q(\bar{s}, \bar{a})|$ 
26:        if  $p > h$  then
27:          Insert  $(\bar{s}, \bar{a})$  into  $\mathcal{H}$  with priority  $p$ 
28:        end if
29:      end for
30:    end for
31:     $s \leftarrow s_{\text{next}}$ 
32:  end while
33: end for
34: for each state  $s \in \mathcal{S}$  do
35:    $\pi(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} Q(s, a)$ 
36: end for
```

---

Algorithm 12 maintains a priority queue (max-heap) with every state-action pair whose estimate would change considerably when updated, ordered by the size of the change. When the pair is updated, the change is propagated to transitions that are predicted to lead to it. This idea dramatically increases the performance of agents in some environments.

Another important heuristic is called trajectory sampling. In this case, transitions are selected for planning following the on-policy distribution. In other words, state-action pairs that are frequently visited are more likely to be sampled for updates by planning.

## 10 Dimensions of Reinforcement Learning

This section presents a summary of the previous sections.

All the reinforcement learning methods have three ideas in common: estimation of value functions, backing up values and generalized policy iteration.

Value backups may vary in the following dimensions:

Table 1: Variation between backup techniques

	Shallow backups	...	Deep backups
Full backups	Dynamic programming	...	Exhaustive search
...	...	...	...
Sample backups	Temporal difference learning	...	Monte Carlo

Another important dimension to be explored in reinforcement learning is function approximation.

In generalized policy iteration, states may be evaluated on-policy (as in Sarsa) or off-policy (as in Q-learning).

There are several important elements to consider while modeling a reinforcement learning problem:

- Return: discounted (continuing) or undiscounted (episodic)
- Value estimation: action value, state value or afterstate value
- Exploration:  $\epsilon$ -greedy, softmax
- Eligibility traces
- Planning: none, deterministic or stochastic

Some less developed research areas in reinforcement learning include the study of non-Markov problems. An important formulation of the reinforcement learning problem in this scenario is called partially observable Markov decision problem (POMDP). The methods used to solve this problem are computationally expensive and harder to apply for model free agents.

Another important area of research is hierarchical learning, where the agent learns important sequences of action-states as primitive actions.

Teaching reinforcement learning agents is also an important area of study. This involves, for example, building problems of increasing difficulty that an agent must learn in order to achieve its most important goals.

## 11 Case Studies

Reinforcement learning techniques have been applied extensively to problems in the real world. Framing a real problem in terms of reinforcement learning is often very challenging. Case studies are useful in training the skills necessary to model problems.

In TD-Gammon, for instance, temporal difference learning with eligibility traces has been applied to the game of backgammon with great success. The agent uses multi-layer artificial neural networks as a function approximator and learns by self-training. In competitive games, self-training is the process whereby an agent learns by playing against itself.

Other case studies presented in [1] include the game of checkers, acrobatic robots, dynamic channel allocation and elevator dispatching,

Interestingly, reinforcement learning has also been applied to job shop scheduling using abstract actions on the state space, composed of schedules. This example shows that some useful agents might have to deal with very abstract states and actions.

Challenging decisions in modeling reinforcement learning problems include the quantization of the state space and formulation of the goal in terms of rewards.

## 12 Policy gradient methods \*

This additional section employs a notation that deviates significantly from the notation employed so far. For details on the (heavily abused) probability notation, see the machine learning notes by the same author.

Consider an agent that interacts with its environment in a sequence of episodes, each of which lasts for exactly  $T$  time steps (without loss of generality for episodic problems).

At each time step  $t < T$ , suppose that the agent observes a state  $s_t \in \text{Val}(S_t)$ , chooses an action  $a_t \in \text{Val}(A_t)$ , and receives a reward  $r_t \in \text{Val}(R_t)$ . Furthermore, suppose that  $\text{Val}(S_t)$ ,  $\text{Val}(A_t)$ , and  $\text{Val}(R_t)$  are finite for all  $t$ .

A policy gradient method represents a policy by a probability distribution over actions given states, and its goal is finding parameters for such a policy that achieve maximum expected return.

Let  $\tau = s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T$  denote a trajectory in a particular episode, such that  $\tau \in \text{Val}(\mathbf{T})$ . Assuming the Markov property, the probability  $p(\tau \mid \theta)$  of trajectory  $\tau$  given the parameters  $\theta$  is given by

$$p(\tau \mid \theta) = p(s_0) \prod_{t=0}^{T-1} p(s_{t+1}, r_t \mid s_t, a_t) p(a_t \mid s_t, \theta),$$

where  $p(a_t \mid s_t, \theta)$  is the probability of action  $a_t$  given state  $s_t$  and parameters  $\theta$ , which is given by the policy.

Given a policy parameterized by  $\theta$ , consider the expected return  $J(\theta)$  given by

$$J(\theta) = \mathbb{E} \left[ \sum_{t=0}^{T-1} R_t \mid \theta \right] = \sum_{t=0}^{T-1} \mathbb{E} [R_t \mid \theta].$$

Using the law of the unconscious statistician,

$$J(\theta) = \sum_{\tau} p(\tau \mid \theta) \sum_{t=0}^{T-1} r_t = \sum_{t=0}^{T-1} \sum_{\tau} r_t p(\tau \mid \theta).$$

Consider the task of finding a  $\theta^*$  such that  $J(\theta^*) = \max_{\theta} J(\theta)$ . Assuming  $J(\theta)$  is differentiable with respect to  $\theta$ , it is possible to attempt to solve this optimization task by gradient-based methods, such as gradient ascent.

The partial derivative  $\frac{\partial}{\partial \theta_j} J(\theta)$  of  $J$  with respect to  $\theta_j$  at  $\theta$  is given by

$$\frac{\partial}{\partial \theta_j} J(\theta) = \sum_{t=0}^{T-1} \sum_{\tau} r_t \frac{\partial}{\partial \theta_j} p(\tau \mid \theta).$$

Suppose that  $p(\tau \mid \theta)$  is positive for any  $\tau$  and  $\theta$ . The so-called likelihood ratio trick uses the fact that

$$\frac{\partial}{\partial \theta_j} p(\tau \mid \theta) = p(\tau \mid \theta) \frac{1}{p(\tau \mid \theta)} \frac{\partial}{\partial \theta_j} p(\tau \mid \theta) = p(\tau \mid \theta) \frac{\partial}{\partial \theta_j} \log p(\tau \mid \theta).$$

By using the expression above,

$$\frac{\partial}{\partial \theta_j} J(\theta) = \sum_{t=0}^{T-1} \sum_{\tau} p(\tau \mid \theta) r_t \frac{\partial}{\partial \theta_j} \log p(\tau \mid \theta).$$

Because we have already assumed that  $p(\tau \mid \theta)$  is positive for all  $\tau$  and  $\theta$ ,

$$\log p(\tau \mid \theta) = \log p(s_0) + \sum_{t=0}^{T-1} \log p(s_{t+1}, r_t \mid s_t, a_t) + \sum_{t=0}^{T-1} \log p(a_t \mid s_t, \theta).$$

Therefore,

$$\frac{\partial}{\partial \theta_j} \log p(\tau \mid \theta) = \sum_{t'=0}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} \mid s_{t'}, \theta).$$

By using the expression above,

$$\frac{\partial}{\partial \theta_j} J(\theta) = \sum_{t=0}^{T-1} \sum_{\tau} p(\tau \mid \theta) r_t \left[ \sum_{t'=0}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} \mid s_{t'}, \theta) \right].$$

It will be useful to split the innermost summation in the expression above into before and after  $t$ , leading to

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=0}^{T-1} \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | \boldsymbol{\theta}) \left[ r_t \sum_{t'=0}^t \frac{\partial}{\partial \theta_j} \log p(a_{t'} | s_{t'}, \boldsymbol{\theta}) + r_t \sum_{t'=t+1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} | s_{t'}, \boldsymbol{\theta}) \right].$$

Alternatively, the expression above can be written as

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=0}^{T-1} \sum_{t'=0}^t \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] + \sum_{t=0}^{T-1} \sum_{t'=t+1}^{T-1} \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right].$$

We will now show that the rightmost nested summations in the expression above can be dismissed.

By representing the random variables involved in a trajectory using a Bayesian network, it can be seen that  $A_{t'} \perp\!\!\!\perp R_t | S_{t'}, \boldsymbol{\theta}$  for  $t' > t$ . In other words, the action at time step  $t' > t$  does not depend on the reward at time step  $t$  given the state at time step  $t'$  and the policy parameters. The analogous statement is not generally true for  $t' \leq t$ . For  $t' > t$ , this independence leads to

$$\mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] = \sum_{r_t} \sum_{a_{t'}} \sum_{s_{t'}} p(a_{t'} | s_{t'}, \boldsymbol{\theta}) p(r_t, s_{t'} | \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} \log p(a_{t'} | s_{t'}, \boldsymbol{\theta}).$$

By reversing the likelihood-ratio trick,

$$\mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] = \sum_{r_t} \sum_{a_{t'}} \sum_{s_{t'}} p(r_t, s_{t'} | \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} p(a_{t'} | s_{t'}, \boldsymbol{\theta}).$$

By changing the order of summations and pushing constants outside the innermost summation,

$$\mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] = \sum_{r_t} \sum_{s_{t'}} p(r_t, s_{t'} | \boldsymbol{\theta}) r_t \sum_{a_{t'}} \frac{\partial}{\partial \theta_j} p(a_{t'} | s_{t'}, \boldsymbol{\theta}).$$

Finally, using the fact that  $\frac{\partial}{\partial \theta_j} 1 = 0$ ,

$$\mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] = \sum_{r_t} \sum_{s_{t'}} p(r_t, s_{t'} | \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} \sum_{a_{t'}} p(a_{t'} | s_{t'}, \boldsymbol{\theta}) = 0.$$

We may now remove the rightmost nested summations in the previous expression for  $\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$ , which gives

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=0}^{T-1} \sum_{t'=0}^t \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] = \mathbb{E} \left[ \sum_{t=0}^{T-1} R_t \sum_{t'=0}^t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | S_{t'}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right].$$

By reordering the summations, the expression above can be conveniently rewritten as

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t | S_t, \boldsymbol{\theta}) \sum_{t'=t}^{T-1} R_{t'} | \boldsymbol{\theta} \right].$$

Finally, the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  of the expected return  $J$  at  $\boldsymbol{\theta}$  is given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(A_t | S_t, \boldsymbol{\theta}) \sum_{t'=t}^{T-1} R_{t'} | \boldsymbol{\theta} \right].$$

This result has a subtle intuitive interpretation. In simple terms, the gradient of the expected return is a sum of expected values of random vectors that correspond to each time step. The expected value for time step  $t$  weights a direction that locally increases the probability of each possible decision by its expected (positive or negative) outcome. Therefore, in gradient ascent, positive expected outcomes contribute towards making the probability of a decision higher, while negative expected outcomes contribute towards making the probability of a decision lower.

Consider a sequence  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_N$  of  $N$  trajectories obtained by following the policy parameterized by  $\boldsymbol{\theta}$ , and let

$$\boldsymbol{\tau}_i = s_{i,0}, a_{i,0}, r_{i,0}, s_{i,1}, a_{i,1}, r_{i,1}, \dots, s_{i,T-1}, a_{i,T-1}, r_{i,T-1}, s_{i,T}.$$

A Monte Carlo estimate  $\hat{\mathbf{g}}(\boldsymbol{\theta})$  to  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  such that  $\hat{\mathbf{g}}(\boldsymbol{\theta}) \rightarrow \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  when  $N \rightarrow \infty$  is given by

$$\hat{\mathbf{g}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(a_{i,t} \mid s_{i,t}, \boldsymbol{\theta}) \sum_{t'=t}^{T-1} r_{i,t'},$$

and may be used for gradient ascent. However, in practice, this estimate may require too many trajectories for accuracy. The remainder of this section presents some alternatives.

Firstly, we will show that introducing a so-called baseline function results in an alternative expression for the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ . Concretely, consider the expected value

$$\mathbb{E} \left[ \sum_{t=0}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, \boldsymbol{\theta}) \left[ \left[ \sum_{t'=t}^{T-1} R_{t'} \right] - B^{\boldsymbol{\theta}}(S_t) \right] \mid \boldsymbol{\theta} \right] = \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) - \sum_{t=0}^{T-1} \mathbb{E} \left[ \frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(S_t) \right],$$

where  $B^{\boldsymbol{\theta}}$  is an arbitrary function that may depend on the policy parameters  $\boldsymbol{\theta}$ . By definition,

$$\mathbb{E} \left[ \frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(S_t) \right] = \sum_{s_t} \sum_{a_t} p(a_t \mid s_t, \boldsymbol{\theta}) p(s_t \mid \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid s_t, \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(s_t).$$

Reversing the likelihood ratio trick,

$$\mathbb{E} \left[ \frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(S_t) \right] = \sum_{s_t} p(s_t \mid \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(s_t) \sum_{a_t} \frac{\partial}{\partial \theta_j} p(a_t \mid s_t, \boldsymbol{\theta}) = \sum_{s_t} p(s_t \mid \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(s_t) \frac{\partial}{\partial \theta_j} 1 = 0.$$

Therefore, an alternative expression for the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  is given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(A_t \mid S_t, \boldsymbol{\theta}) \left[ \left[ \sum_{t'=t}^{T-1} R_{t'} \right] - B^{\boldsymbol{\theta}}(S_t) \right] \mid \boldsymbol{\theta} \right].$$

A common choice of baseline function  $B^{\boldsymbol{\theta}}$  is an approximation to the (typically unknown) value function  $V^{\boldsymbol{\theta}}$  given by

$$V^{\boldsymbol{\theta}}(s) = \mathbb{E} \left[ \sum_{t'=t}^{T-1} R_{t'} \mid s_t = s, \boldsymbol{\theta} \right].$$

Because  $V^{\boldsymbol{\theta}}(s)$  is the expected return for starting at state  $s$  and following a policy parameterized by  $\boldsymbol{\theta}$ , the corresponding estimate to  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  intuitively tends to increase the probability of actions that are followed by surprisingly large returns, and is generally more efficient in practice than the alternative presented so far. For details on value function approximation methods, see the previous sections.

Consider once again the partial derivative  $\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$  of  $J$  with respect to  $\theta_j$  at  $\boldsymbol{\theta}$  given by

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, \boldsymbol{\theta}) \sum_{t'=t}^{T-1} R_{t'} \mid \boldsymbol{\theta} \right].$$

By definition,

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) &= \sum_{t=0}^{T-1} \sum_{s_t} \sum_{a_t} \sum_{r_t} \cdots \sum_{r_{T-1}} p(s_t, a_t, r_t, \dots, r_{T-1} \mid \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid s_t, \boldsymbol{\theta}) \sum_{t'=t}^{T-1} r_{t'} \\ &= \sum_{t=0}^{T-1} \sum_{s_t} \sum_{a_t} p(s_t, a_t \mid \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid s_t, \boldsymbol{\theta}) \sum_{r_t} \cdots \sum_{r_{T-1}} p(r_t, \dots, r_{T-1} \mid s_t, a_t, \boldsymbol{\theta}) \sum_{t'=t}^{T-1} r_{t'} \\ &= \mathbb{E} \left[ \sum_{t=0}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, \boldsymbol{\theta}) Q^{\boldsymbol{\theta}}(S_t, A_t) \mid \boldsymbol{\theta} \right], \end{aligned}$$



where the so-called action value function  $Q^\theta$  is given by

$$Q^\theta(s, a) = \mathbb{E} \left[ \sum_{t'=t}^{T-1} R_{t'} \mid s_t = s, a_t = a, \theta \right].$$

Intuitively,  $Q^\theta(s, a)$  is the expected return for starting at state  $s$ , choosing action  $a$ , and then following the policy parameterized by  $\theta$ .

Therefore, another expression for the gradient  $\nabla_\theta J(\theta)$  of the expected return  $J$  at  $\theta$  is given by

$$\nabla_\theta J(\theta) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log p(A_t \mid S_t, \theta) Q^\theta(S_t, A_t) \mid \theta \right].$$

However, using this expression to obtain an estimate to the gradient  $\nabla_\theta J(\theta)$  typically requires approximating the unknown action value function  $Q^\theta$ . Such approximation will generally introduce bias (convergence to an undesired expected value) in the corresponding estimate to  $\nabla_\theta J(\theta)$ . For details on action value function approximation methods, see the previous sections.

The advantage function  $A^\theta$  is defined by  $A^\theta(s, a) = Q^\theta(s, a) - V^\theta(s)$ , for all  $s$  and  $a$ . By combining the previous results, yet another expression for the gradient  $\nabla_\theta J(\theta)$  of the expected return  $J$  at  $\theta$  is given by

$$\nabla_\theta J(\theta) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log p(A_t \mid S_t, \theta) A^\theta(S_t, A_t) \mid \theta \right].$$

Using this expression to obtain an estimate to the gradient  $\nabla_\theta J(\theta)$  typically requires approximating the unknown advantage function  $A^\theta$ . Such approximation also generally introduces a bias in the corresponding estimate to  $\nabla_\theta J(\theta)$ . However, in practice, such estimate is generally more efficient than the alternatives presented so far.

## 13 Recurrent policy gradient methods \*

In contrast to the reinforcement learning methods presented in the previous sections, recurrent policy gradient methods admit partially observable environments [5]. This additional section employs a notation that deviates significantly from the notation employed so far, except for the previous section. For details on the (heavily abused) probability notation, see the machine learning notes by the same author.

Consider an agent that interacts with its environment in a sequence of episodes, each of which lasts for exactly  $T > 1$  time steps (without loss of generality for episodic problems).

At each time step  $t$ , suppose that the environment is in hidden state  $s_t \in \text{Val}(S_t)$ , the agent receives an observation  $x_t \in \text{Val}(X_t)$ , a reward  $r_t \in \text{Val}(R_t)$ , and chooses an action  $a_t \in \text{Val}(A_t)$ . Furthermore, suppose that  $\text{Val}(S_t)$ ,  $\text{Val}(X_t)$ ,  $\text{Val}(R_t)$ , and  $\text{Val}(A_t)$  are finite, for all  $t$ .

A recurrent policy gradient method represents a policy by a probability distribution over actions given the history of observations and actions, and its goal is finding parameters for such a policy that achieve maximum expected return.

Let  $\tau = s_1, x_1, r_1, a_1, \dots, s_{T-1}, x_{T-1}, r_{T-1}, a_{T-1}, s_T, x_T, r_T$  denote a trajectory in a particular episode, such that  $\tau \in \text{Val}(\mathbf{T})$ . Furthermore, let  $v_{i:j} = v_i, v_{i+1}, \dots, v_j$ .

We will assume that the probability  $p(\tau \mid \theta)$  of trajectory  $\tau$  given the parameters  $\theta$  is given by

$$p(\tau \mid \theta) = p(s_1)p(x_1 \mid s_1)p(r_1 \mid s_1) \prod_{t=1}^{T-1} p(a_t \mid x_{1:t}, a_{1:t-1}, \theta)p(s_{t+1} \mid s_t, a_t)p(x_{t+1} \mid s_{t+1})p(r_{t+1} \mid s_{t+1}),$$

where  $p(a_t \mid x_{1:t}, a_{1:t-1}, \theta)$  is the probability of action  $a_t$  given the history of observations and actions up to time  $t$ , which is given by the policy. Intuitively, the hidden state  $s_t$  encodes all the information required to generate the observation  $x_t$  and the reward  $r_t$ . Furthermore, the action  $a_t$  may depend on the observed history up to time  $t$  (excluding the rewards, which may be encoded into the observations without loss of generality).

Given a policy parameterized by  $\theta$ , consider the expected return  $J(\theta)$  given by

$$J(\theta) = \mathbb{E} \left[ \sum_{t=1}^T R_t \mid \theta \right] = \sum_{t=1}^T \mathbb{E} [R_t \mid \theta].$$

Using the law of the unconscious statistician,

$$J(\boldsymbol{\theta}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | \boldsymbol{\theta}) \sum_{t=1}^T r_t = \sum_{t=1}^T \sum_{\boldsymbol{\tau}} r_t p(\boldsymbol{\tau} | \boldsymbol{\theta}).$$

Consider the task of finding a  $\boldsymbol{\theta}^*$  such that  $J(\boldsymbol{\theta}^*) = \max_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ . Assuming  $J(\boldsymbol{\theta})$  is differentiable with respect to  $\boldsymbol{\theta}$ , it is possible to attempt to solve this optimization task by gradient-based methods, such as gradient ascent.

The partial derivative  $\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$  of  $J$  with respect to  $\theta_j$  at  $\boldsymbol{\theta}$  is given by

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{\boldsymbol{\tau}} r_t \frac{\partial}{\partial \theta_j} p(\boldsymbol{\tau} | \boldsymbol{\theta}).$$

Suppose that  $p(\boldsymbol{\tau} | \boldsymbol{\theta})$  is positive for any  $\boldsymbol{\tau}$  and  $\boldsymbol{\theta}$ . The so-called likelihood ratio trick uses the fact that

$$\frac{\partial}{\partial \theta_j} p(\boldsymbol{\tau} | \boldsymbol{\theta}) = p(\boldsymbol{\tau} | \boldsymbol{\theta}) \frac{1}{p(\boldsymbol{\tau} | \boldsymbol{\theta})} \frac{\partial}{\partial \theta_j} p(\boldsymbol{\tau} | \boldsymbol{\theta}) = p(\boldsymbol{\tau} | \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(\boldsymbol{\tau} | \boldsymbol{\theta}).$$

By using the expression above,

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} \log p(\boldsymbol{\tau} | \boldsymbol{\theta}).$$

Because we have already assumed that  $p(\boldsymbol{\tau} | \boldsymbol{\theta})$  is positive for all  $\boldsymbol{\tau}$  and  $\boldsymbol{\theta}$ ,

$$\log p(\boldsymbol{\tau} | \boldsymbol{\theta}) = \log [p(s_1)p(x_1 | s_1)p(r_1 | s_1)] + \sum_{t=1}^{T-1} \log [p(a_t | x_{1:t}, a_{1:t-1}, \boldsymbol{\theta})p(s_{t+1} | s_t, a_t)p(x_{t+1} | s_{t+1})p(r_{t+1} | s_{t+1})].$$

Therefore,

$$\frac{\partial}{\partial \theta_j} \log p(\boldsymbol{\tau} | \boldsymbol{\theta}) = \sum_{t'=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}).$$

By using the expression above,

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | \boldsymbol{\theta}) r_t \left[ \sum_{t'=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}) \right].$$

It will be useful to split the innermost summation in the expression above into before and after  $t$ , leading to

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | \boldsymbol{\theta}) \left[ r_t \sum_{t'=1}^{t-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}) + r_t \sum_{t'=t}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}) \right].$$

Alternatively, the expression above can be written as

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{t'=1}^{t-1} \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | X_{1:t'}, A_{1:t'-1}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] + \sum_{t=1}^T \sum_{t'=t}^{T-1} \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | X_{1:t'}, A_{1:t'-1}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right].$$

We will now show that the rightmost nested summations in the expression above can be dismissed, because  $E = \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} | X_{1:t'}, A_{1:t'-1}, \boldsymbol{\theta}) | \boldsymbol{\theta} \right] = 0$  for  $t' \geq t$ .

By representing the random variables involved in a trajectory using a Bayesian network, it can be seen that  $A_{t'} \perp\!\!\!\perp R_t | X_{1:t'}, A_{1:t'-1}, \boldsymbol{\theta}$  for  $t' \geq t$ . In other words, the action at time step  $t' \geq t$  does not depend on the reward at time step  $t$  given the history of observations and actions up to time  $t'$  and the policy parameters. The analogous statement is not generally true for  $t' < t$ . For  $t' \geq t$ , this independence leads to

$$\begin{aligned} E &= \sum_{r_t} \sum_{x_{1:t'}} \sum_{a_{1:t'-1}} \sum_{a_{t'}} p(r_t, x_{1:t'}, a_{1:t'-1}, a_{t'} | \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} \log p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}) \\ &= \sum_{r_t} \sum_{x_{1:t'}} \sum_{a_{1:t'-1}} \sum_{a_{t'}} p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}) p(r_t, x_{1:t'}, a_{1:t'-1} | \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} \log p(a_{t'} | x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}). \end{aligned}$$

By reversing the likelihood-ratio trick,

$$E = \sum_{r_t} \sum_{x_{1:t'}} \sum_{a_{1:t'-1}} \sum_{a_{t'}} p(r_t, x_{1:t'}, a_{1:t'-1} \mid \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} p(a_{t'} \mid x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}).$$

By pushing constants outside the innermost summation and using the fact that  $\frac{\partial}{\partial \theta_j} 1 = 0$ ,

$$E = \sum_{r_t} \sum_{x_{1:t'}} \sum_{a_{1:t'-1}} p(r_t, x_{1:t'}, a_{1:t'-1} \mid \boldsymbol{\theta}) r_t \frac{\partial}{\partial \theta_j} \sum_{a_{t'}} p(a_{t'} \mid x_{1:t'}, a_{1:t'-1}, \boldsymbol{\theta}) = 0.$$

We may now remove the rightmost nested summations in the previous expression for  $\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$ , which gives

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{t'=1}^{t-1} \mathbb{E} \left[ R_t \frac{\partial}{\partial \theta_j} \log p(A_{t'} \mid X_{1:t'}, A_{1:t'-1}, \boldsymbol{\theta}) \mid \boldsymbol{\theta} \right] = \mathbb{E} \left[ \sum_{t=1}^T R_t \sum_{t'=1}^{t-1} \frac{\partial}{\partial \theta_j} \log p(A_{t'} \mid X_{1:t'}, A_{1:t'-1}, \boldsymbol{\theta}) \mid \boldsymbol{\theta} \right].$$

By reordering the summations, the expression above can be conveniently rewritten as

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) \sum_{t'=t+1}^T R_{t'} \mid \boldsymbol{\theta} \right].$$

Finally, the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  of the expected return  $J$  at  $\boldsymbol{\theta}$  is given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=1}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) \sum_{t'=t+1}^T R_{t'} \mid \boldsymbol{\theta} \right].$$

This result has a subtle intuitive interpretation. In simple terms, the gradient of the expected return is a sum of expected values of random vectors that correspond to each time step. The expected value for time step  $t$  weights a direction that locally increases the probability of each possible decision by its expected (positive or negative) outcome. Therefore, in gradient ascent, positive expected outcomes contribute towards making the probability of a decision higher, while negative expected outcomes contribute towards making the probability of a decision lower.

Consider a sequence  $\boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_N$  of  $N$  trajectories obtained by following the policy parameterized by  $\boldsymbol{\theta}$ , and let

$$\boldsymbol{\tau}_i = s_{i,1}, x_{i,1}, r_{i,1}, a_{i,1}, \dots, s_{i,T-1}, x_{i,T-1}, r_{i,T-1}, a_{i,T-1}, s_{i,T}, x_{i,T}, r_{i,T}.$$

A Monte Carlo estimate  $\hat{\mathbf{g}}(\boldsymbol{\theta})$  to  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  such that  $\hat{\mathbf{g}}(\boldsymbol{\theta}) \rightarrow \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  when  $N \rightarrow \infty$  is given by

$$\hat{\mathbf{g}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(a_{i,t} \mid x_{i,1:t}, a_{i,1:t-1}, \boldsymbol{\theta}) \sum_{t'=t+1}^T r_{i,t'},$$

and may be used for gradient ascent. Notice that this estimate does not require the hidden states of the environment. However, in practice, it may require too many trajectories for accuracy. The remainder of this section presents some alternatives.

Firstly, we will show that introducing a so-called baseline function results in an alternative expression for the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ . Concretely, consider the expected value

$$\mathbb{E} \left[ \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) \left[ \sum_{t'=t+1}^T R_{t'} - B^{\boldsymbol{\theta}}(X_{1:t}, A_{1:t-1}) \right] \mid \boldsymbol{\theta} \right] = \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) - C,$$

where  $B^{\boldsymbol{\theta}}$  is an arbitrary function that may depend on the policy parameters  $\boldsymbol{\theta}$ , and

$$C = \mathbb{E} \left[ \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) B^{\boldsymbol{\theta}}(X_{1:t}, A_{1:t-1}) \mid \boldsymbol{\theta} \right].$$

By definition,

$$\begin{aligned} C &= \sum_{t=1}^{T-1} \sum_{x_{1:t}} \sum_{a_{1:t-1}} \sum_{a_t} p(x_{1:t}, a_{1:t-1}, a_t \mid \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid x_{1:t}, a_{1:t-1}, \boldsymbol{\theta}) B^\theta(x_{1:t}, a_{1:t-1}) \\ &= \sum_{t=1}^{T-1} \sum_{x_{1:t}} \sum_{a_{1:t-1}} p(x_{1:t}, a_{1:t-1} \mid \boldsymbol{\theta}) B^\theta(x_{1:t}, a_{1:t-1}) \sum_{a_t} p(a_t \mid x_{1:t}, a_{1:t-1}, \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid x_{1:t}, a_{1:t-1}, \boldsymbol{\theta}). \end{aligned}$$

Reversing the likelihood ratio trick and using the fact that  $\frac{\partial}{\partial \theta_j} 1 = 0$ ,

$$C = \sum_{t=1}^{T-1} \sum_{x_{1:t}} \sum_{a_{1:t-1}} p(x_{1:t}, a_{1:t-1} \mid \boldsymbol{\theta}) B^\theta(x_{1:t}, a_{1:t-1}) \frac{\partial}{\partial \theta_j} \sum_{a_t} p(a_t \mid x_{1:t}, a_{1:t-1}, \boldsymbol{\theta}) = 0.$$

Therefore, an alternative expression for the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  is given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=1}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) \left[ \left[ \sum_{t'=t+1}^T R_{t'} \right] - B^\theta(X_{1:t}, A_{1:t-1}) \right] \mid \boldsymbol{\theta} \right].$$

A possible choice for baseline function  $B^\theta$  is an approximation to the (typically unknown) value function  $V^\theta$  given by

$$V^\theta(x'_{1:t}, a'_{1:t-1}) = \sum_{t'=t+1}^T \mathbb{E} [R_{t'} \mid X_{1:t} = x'_{1:t}, A_{1:t-1} = a'_{1:t-1}, \boldsymbol{\theta}].$$

Because  $V^\theta$  gives the expected return for following a policy parameterized by  $\boldsymbol{\theta}$  after a particular history of observations and actions, the corresponding estimate to  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  intuitively tends to increase the probability of actions that are followed by surprisingly large returns, and may be more efficient in practice than the alternative presented so far. Value function approximation for partially observable environments requires extensions to the methods presented in the previous sections.

Consider once again the partial derivative  $\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$  of  $J$  with respect to  $\theta_j$  at  $\boldsymbol{\theta}$  given by

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) \sum_{t'=t+1}^T R_{t'} \mid \boldsymbol{\theta} \right].$$

By definition,

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) &= \sum_{t=1}^{T-1} \sum_{x_{1:t}} \sum_{a_{1:t}} \sum_{r_{t+1:T}} p(x_{1:t}, a_{1:t}, r_{t+1:T} \mid \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid x_{1:t}, a_{1:t-1}, \boldsymbol{\theta}) \sum_{t'=t+1}^T r_{t'} \\ &= \sum_{t=1}^{T-1} \sum_{x_{1:t}} \sum_{a_{1:t}} p(x_{1:t}, a_{1:t} \mid \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t \mid x_{1:t}, a_{1:t-1}, \boldsymbol{\theta}) \sum_{r_{t+1:T}} p(r_{t+1:T} \mid x_{1:t}, a_{1:t}, \boldsymbol{\theta}) \sum_{t'=t+1}^T r_{t'} \\ &= \mathbb{E} \left[ \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) Q^\theta(X_{1:t}, A_{1:t}) \mid \boldsymbol{\theta} \right], \end{aligned}$$

where the so-called action value function  $Q^\theta$  is given by

$$Q^\theta(x'_{1:t}, a'_{1:t}) = \sum_{t'=t+1}^T \mathbb{E} [R_{t'} \mid X_{1:t} = x'_{1:t}, A_{1:t} = a'_{1:t}, \boldsymbol{\theta}].$$

Intuitively,  $Q^\theta$  also gives the expected return for following the policy parameterized by  $\boldsymbol{\theta}$  after a particular history of observations and actions.

Therefore, another expression for the gradient  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  of the expected return  $J$  at  $\boldsymbol{\theta}$  is given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E} \left[ \sum_{t=1}^{T-1} \nabla_{\boldsymbol{\theta}} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \boldsymbol{\theta}) Q^\theta(X_{1:t}, A_{1:t}) \mid \boldsymbol{\theta} \right].$$


However, using this expression to obtain an estimate to the gradient  $\nabla_{\theta}J(\theta)$  typically requires approximating the unknown action value function  $Q^{\theta}$ . Such approximation will generally introduce bias (convergence to an undesired expected value) in the corresponding estimate to  $\nabla_{\theta}J(\theta)$ . Action value function approximation for partially observable environments requires extensions to the methods presented in the previous sections.

The advantage function  $A^{\theta}$  is defined by  $A^{\theta}(x_{1:t}, a_{1:t}) = Q^{\theta}(x_{1:t}, a_{1:t}) - V^{\theta}(x_{1:t}, a_{1:t-1})$ , for all  $x_{1:t}$  and  $a_{1:t}$ . By combining the previous results, yet another expression for the gradient  $\nabla_{\theta}J(\theta)$  of the expected return  $J$  at  $\theta$  is given by

$$\nabla_{\theta}J(\theta) = \mathbb{E} \left[ \sum_{t=1}^{T-1} \nabla_{\theta} \log p(A_t \mid X_{1:t}, A_{1:t-1}, \theta) A^{\theta}(X_{1:t}, A_{1:t}) \mid \theta \right].$$

Using this expression to obtain an estimate to the gradient  $\nabla_{\theta}J(\theta)$  typically requires approximating the unknown advantage function  $A^{\theta}$ . Such approximation also generally introduces a bias in the corresponding estimate to  $\nabla_{\theta}J(\theta)$ . However, in practice, such estimate may be more efficient than the alternatives presented so far.

## License

This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License .

## References

- [1] Sutton, R.S. and Barto, A.G., *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [2] Peters, J. and Schaal, S. *Policy gradient methods for robotics*. International Conference on Intelligent Robots and Systems, 2006.
- [3] Schulman, J. and Abbeel, P. *Deep Reinforcement Learning*, Available at <http://rll.berkeley.edu/deeprlcourse>, 2016.
- [4] Silver, D. *Reinforcement learning*, Available at <http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>, 2016.
- [5] Wierstra, D., Förster, A., Peters, J., Schmidhuber, J. *Recurrent policy gradients*. Logic Journal of IGPL, 2010.