

# 18 *Reinforcement Learning*

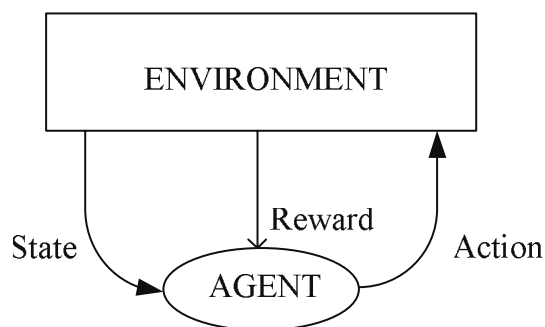
*In reinforcement learning, the learner is a decision-making agent that takes actions in an environment and receives reward (or penalty) for its actions in trying to solve a problem. After a set of trial-and-error runs, it should learn the best policy, which is the sequence of actions that maximize the total reward.*

## 18.1 Introduction

LET US SAY we want to build a machine that learns to play chess. In this case we cannot use a supervised learner for two reasons. First, it is very costly to have a teacher that will take us through many games and indicate us the best move for each position. Second, in many cases, there is no such thing as the best move; the goodness of a move depends on the moves that follow. A single move does not count; a sequence of moves is good if after playing them we win the game. The only feedback is at the end of the game when we win or lose the game.

Another example is a robot that is placed in a maze. The robot can move in one of the four compass directions and should make a sequence of movements to reach the exit. As long as the robot is in the maze, there is no feedback and the robot tries many moves until it reaches the exit and only then does it get a reward. In this case there is no opponent, but we can have a preference for shorter trajectories, implying that in this case we play against time.

These two applications have a number of points in common: There is a decision maker, called the *agent*, that is placed in an *environment* (see figure 18.1). In chess, the game-player is the decision maker and the environment is the board; in the second case, the maze is the environment



**Figure 18.1** The agent interacts with an environment. At any state of the environment, the agent takes an action that changes the state and returns a reward.

of the robot. At any time, the environment is in a certain *state* that is one of a set of possible states—for example, the state of the board, the position of the robot in the maze. The decision maker has a set of *actions* possible: legal movement of pieces on the chess board, movement of the robot in possible directions without hitting the walls, and so forth. Once an action is chosen and taken, the state changes. The solution to the task requires a sequence of actions, and we get feedback, in the form of a *reward* rarely, generally only when the complete sequence is carried out. The reward defines the problem and is necessary if we want a *learning* agent. The learning agent learns the best sequence of actions to solve a problem where “best” is quantified as the sequence of actions that has the maximum cumulative reward. Such is the setting of *reinforcement learning*.

Reinforcement learning is different from the learning methods we discussed before in a number of respects. It is called “learning with a critic,” as opposed to learning with a teacher which we have in supervised learning. A *critic* differs from a teacher in that it does not tell us what to do but only how well we have been doing in the past; the critic never informs in advance. The feedback from the critic is scarce and when it comes, it comes late. This leads to the *credit assignment* problem. After taking several actions and getting the reward, we would like to assess the individual actions we did in the past and find the moves that led us to win the reward so that we can record and recall them later on. As we see shortly, what a reinforcement learning program does is that it learns to generate

CREDIT ASSIGNMENT

an *internal value* for the intermediate states or actions in terms of how good they are in leading us to the goal and getting us to the real reward. Once such an internal reward mechanism is learned, the agent can just take the local actions to maximize it.

The solution to the task requires a *sequence* of actions, and from this perspective, we remember the Markov models we discussed in chapter 15. Indeed, we use a Markov decision process to model the agent. The difference is that in the case of Markov models, there is an external process that generates a sequence of signals, for example, speech, which we observe and model. In the current case, however, it is the agent that generates the sequence of actions. Previously, we also made a distinction between observable and hidden Markov models where the states are observed or hidden (and should be inferred) respectively. Similarly here, sometimes we have a partially observable Markov decision process in cases where the agent does not know its state exactly but should infer it with some uncertainty through observations using sensors. For example, in the case of a robot moving in a room, the robot may not know its exact position in the room, nor the exact location of obstacles nor the goal, and should make decisions through a limited image provided by a camera.

## 18.2 Single State Case: *K*-Armed Bandit

### *K*-ARMED BANDIT

We start with a simple example. The *K-armed bandit* is a hypothetical slot machine with  $K$  levers. The action is to choose and pull one of the levers, and we win a certain amount of money that is the reward associated with the lever (action). The task is to decide which lever to pull to maximize the reward. This is a classification problem where we choose one of  $K$ . If this were supervised learning, then the teacher would tell us the correct class, namely, the lever leading to maximum earning. In this case of reinforcement learning, we can only try different levers and keep track of the best. This is a simplified reinforcement learning problem because there is only one state, or one slot machine, and we need only decide on the action. Another reason why this is simplified is that we immediately get a reward after a single action; the reward is not delayed, so we immediately see the value of our action.

Let us say  $Q(a)$  is the value of action  $a$ . Initially,  $Q(a) = 0$  for all  $a$ . When we try action  $a$ , we get reward  $r_a \geq 0$ . If rewards are deterministic, we always get the same  $r_a$  for any pull of  $a$  and in such a case, we can

just set  $Q(a) = r_a$ . If we want to exploit, once we find an action  $a$  such that  $Q(a) > 0$ , we can keep choosing it and get  $r_a$  at each pull. However, it is quite possible that there is another lever with a higher reward, so we need to explore.

We can choose different actions and store  $Q(a)$  for all  $a$ . Whenever we want to exploit, we can choose the action with the maximum value, that is,

$$(18.1) \quad \text{choose } a^* \text{ if } Q(a^*) = \max_a Q(a)$$

If rewards are not deterministic but stochastic, we get a different reward each time we choose the same action. The amount of the reward is defined by the probability distribution  $p(r|a)$ . In such a case, we define  $Q_t(a)$  as the estimate of the value of action  $a$  at time  $t$ . It is an average of all rewards received when action  $a$  was chosen before time  $t$ . An online update can be defined as

$$(18.2) \quad Q_{t+1}(a) \leftarrow Q_t(a) + \eta[r_{t+1}(a) - Q_t(a)]$$

where  $r_{t+1}(a)$  is the reward received after taking action  $a$  at time  $(t + 1)$ st time.

Note that equation 18.2 is the *delta rule* that we have used on many occasions in the previous chapters:  $\eta$  is the learning factor (gradually decreased in time for convergence),  $r_{t+1}$  is the desired output, and  $Q_t(a)$  is the current prediction.  $Q_{t+1}(a)$  is the *expected* value of action  $a$  at time  $t + 1$  and converges to the mean of  $p(r|a)$  as  $t$  increases.

The full reinforcement learning problem generalizes this simple case in a number of ways. First, we have several states. This corresponds to having several slot machines with different reward probabilities,  $p(r|s_i, a_j)$ , and we need to learn  $Q(s_i, a_j)$ , which is the value of taking action  $a_j$  when in state  $s_i$ . Second, the actions affect not only the reward but also the next state, and we move from one state to another. Third, the rewards are delayed and we need to be able to estimate immediate values from delayed rewards.

### 18.3 Elements of Reinforcement Learning

The learning decision maker is called the *agent*. The agent interacts with the *environment* that includes everything outside the agent. The agent has sensors to decide on its *state* in the environment and takes an *action*

that modifies its state. When the agent takes an action, the environment provides a *reward*. Time is discrete as  $t = 0, 1, 2, \dots$ , and  $s_t \in S$  denotes the state of the agent at time  $t$  where  $S$  is the set of all possible states.  $a_t \in \mathcal{A}(s_t)$  denotes the action that the agent takes at time  $t$  where  $\mathcal{A}(s_t)$  is the set of possible actions in state  $s_t$ . When the agent in state  $s_t$  takes the action  $a_t$ , the clock ticks, reward  $r_{t+1} \in \mathfrak{R}$  is received, and the agent moves to the next state,  $s_{t+1}$ . The problem is modeled using a *Markov decision process* (MDP). The reward and next state are sampled from their respective probability distributions,  $p(r_{t+1}|s_t, a_t)$  and  $P(s_{t+1}|s_t, a_t)$ . Note that what we have is a *Markov system* where the state and reward in the next time step depend only on the current state and action. In some applications, reward and next state are deterministic, and for a certain state and action taken, there is one possible reward value and next state.

MARKOV DECISION  
PROCESS

Depending on the application, a certain state may be designated as the initial state and in some applications, there is also an absorbing terminal (goal) state where the search ends; all actions in this terminal state transition to itself with probability 1 and without any reward. The sequence of actions from the start to the terminal state is an *episode*, or a *trial*.

EPISODE  
POLICY

The *policy*,  $\pi$ , defines the agent's behavior and is a mapping from the states of the environment to actions:  $\pi : S \rightarrow \mathcal{A}$ . The policy defines the action to be taken in any state  $s_t$ :  $a_t = \pi(s_t)$ . The *value* of a policy  $\pi$ ,  $V^\pi(s_t)$ , is the expected cumulative reward that will be received while the agent follows the policy, starting from state  $s_t$ .

FINITE-HORIZON

In the *finite-horizon* or *episodic* model, the agent tries to maximize the expected reward for the next  $T$  steps:

$$(18.3) \quad V^\pi(s_t) = E[r_{t+1} + r_{t+2} + \dots + r_{t+T}] = E \left[ \sum_{i=1}^T r_{t+i} \right]$$

INFINITE-HORIZON

Certain tasks are continuing, and there is no prior fixed limit to the episode. In the *infinite-horizon* model, there is no sequence limit, but future rewards are discounted:

$$(18.4) \quad V^\pi(s_t) = E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] = E \left[ \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i} \right]$$

DISCOUNT RATE

where  $0 \leq \gamma < 1$  is the *discount rate* to keep the return finite. If  $\gamma = 0$ , then only the immediate reward counts. As  $\gamma$  approaches 1, rewards further in the future count more, and we say that the agent becomes more farsighted.  $\gamma$  is less than 1 because there generally is a time limit

to the sequence of actions needed to solve the task. The agent may be a robot that runs on a battery. We prefer rewards sooner rather than later because we are not certain how long we will survive.

OPTIMAL POLICY For each policy  $\pi$ , there is a  $V^\pi(s_t)$ , and we want to find the *optimal policy*  $\pi^*$  such that

$$(18.5) \quad V^*(s_t) = \max_{\pi} V^\pi(s_t), \forall s_t$$

In some applications, for example, in control, instead of working with the values of states,  $V(s_t)$ , we prefer to work with the values of state-action pairs,  $Q(s_t, a_t)$ .  $V(s_t)$  denotes how good it is for the agent to be in state  $s_t$ , whereas  $Q(s_t, a_t)$  denotes how good it is to perform action  $a_t$  when in state  $s_t$ . We define  $Q^*(s_t, a_t)$  as the value, that is, the expected cumulative reward, of action  $a_t$  taken in state  $s_t$  and then obeying the optimal policy afterward. The value of a state is equal to the value of the best possible action:

$$\begin{aligned}
 V^*(s_t) &= \max_{a_t} Q^*(s_t, a_t) \\
 &= \max_{a_t} E \left[ \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i} \right] \\
 &= \max_{a_t} E \left[ r_{t+1} + \gamma \sum_{i=1}^{\infty} \gamma^{i-1} r_{t+i+1} \right] \\
 &= \max_{a_t} E [r_{t+1} + \gamma V^*(s_{t+1})] \\
 (18.6) \quad V^*(s_t) &= \max_{a_t} \left( E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, a_t) V^*(s_{t+1}) \right)
 \end{aligned}$$

To each possible next state  $s_{t+1}$ , we move with probability  $P(s_{t+1}|s_t, a_t)$ , and continuing from there using the optimal policy, the expected cumulative reward is  $V^*(s_{t+1})$ . We sum over all such possible next states, and we discount it because it is one time step later. Adding our immediate expected reward, we get the total expected cumulative reward for action  $a_t$ . We then choose the best of possible actions. Equation 18.6 is known as *Bellman's equation* (Bellman 1957). Similarly, we can also write

BELLMAN'S EQUATION

$$(18.7) \quad Q^*(s_t, a_t) = E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1}|s_t, a_t) \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$$

```

Initialize  $V(s)$  to arbitrary values
Repeat
  For all  $s \in \mathcal{S}$ 
    For all  $a \in \mathcal{A}$ 
       $Q(s, a) \leftarrow E[r|s, a] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)V(s')$ 
       $V(s) \leftarrow \max_a Q(s, a)$ 
Until  $V(s)$  converge

```

**Figure 18.2** Value iteration algorithm for model-based learning.

Once we have  $Q^*(s_t, a_t)$  values, we can then define our policy  $\pi$  as taking the action  $a_t^*$ , which has the highest value among all  $Q^*(s_t, a_t)$ :

$$(18.8) \quad \pi^*(s_t) : \text{Choose } a_t^* \text{ where } Q^*(s_t, a_t^*) = \max_{a_t} Q^*(s_t, a_t)$$

This means that if we have the  $Q^*(s_t, a_t)$  values, then by using a greedy search at each *local* step we get the optimal sequence of steps that maximizes the *cumulative* reward.

## 18.4 Model-Based Learning

We start with model-based learning where we completely know the environment model parameters,  $p(r_{t+1}|s_t, a_t)$  and  $P(s_{t+1}|s_t, a_t)$ . In such a case, we do not need any exploration and can directly solve for the optimal value function and policy using dynamic programming. The optimal value function is unique and is the solution to the simultaneous equations given in equation 18.6. Once we have the optimal value function, the optimal policy is to choose the action that maximizes the value in the next state:

$$(18.9) \quad \pi^*(s_t) = \arg \max_{a_t} \left( E[r_{t+1}|s_t, a_t] + \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1}|s_t, a_t)V^*(s_{t+1}) \right)$$

### 18.4.1 Value Iteration

To find the optimal policy, we can use the optimal value function, and there is an iterative algorithm called *value iteration* that has been shown to converge to the correct  $V^*$  values. Its pseudocode is given in figure 18.2.

```

Initialize a policy  $\pi'$  arbitrarily
Repeat
   $\pi \leftarrow \pi'$ 
  Compute the values using  $\pi$  by
    solving the linear equations
       $V^\pi(s) = E[r|s, \pi(s)] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, \pi(s)) V^\pi(s')$ 
  Improve the policy at each state
       $\pi'(s) \leftarrow \arg \max_a (E[r|s, a] + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V^\pi(s'))$ 
Until  $\pi = \pi'$ 

```

**Figure 18.3** Policy iteration algorithm for model-based learning.

We say that the values converged if the maximum value difference between two iterations is less than a certain threshold  $\delta$ :

$$\max_{s \in \mathcal{S}} |V^{(l+1)}(s) - V^{(l)}(s)| < \delta$$

where  $l$  is the iteration counter. Because we care only about the actions with the maximum value, it is possible that the policy converges to the optimal one even before the values converge to their optimal values. Each iteration is  $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|)$ , but frequently there is only a small number  $k < |S|$  of next possible states, so complexity decreases to  $\mathcal{O}(k|S||\mathcal{A}|)$ .

### 18.4.2 Policy Iteration

In policy iteration, we store and update the policy rather than doing this indirectly over the values. The pseudocode is given in figure 18.3. The idea is to start with a policy and improve it repeatedly until there is no change. The value function can be calculated by solving for the linear equations. We then check whether we can improve the policy by taking these into account. This step is guaranteed to improve the policy, and when no improvement is possible, the policy is guaranteed to be optimal. Each iteration of this algorithm takes  $\mathcal{O}(|\mathcal{A}||S|^2 + |S|^3)$  time that is more than that of value iteration, but policy iteration needs fewer iterations than value iteration.



## 18.5 Temporal Difference Learning

Model is defined by the reward and next state probability distributions, and as we saw in section 18.4, when we know these, we can solve for the optimal policy using dynamic programming. However, these methods are costly, and we seldom have such perfect knowledge of the environment. The more interesting and realistic application of reinforcement learning is when we do not have the model. This requires exploration of the environment to query the model. We first discuss how this exploration is done and later see model-free learning algorithms for deterministic and nondeterministic cases. Though we are not going to assume a full knowledge of the environment model, we will however require that it be stationary.

TEMPORAL  
DIFFERENCE

As we will see shortly, when we explore and get to see the value of the next state and reward, we use this information to update the value of the current state. These algorithms are called *temporal difference* algorithms because what we do is look at the difference between our current estimate of the value of a state (or a state-action pair) and the discounted value of the next state and the reward received.

### 18.5.1 Exploration Strategies

To explore, one possibility is to use  $\epsilon$ -greedy search where with probability  $\epsilon$ , we choose one action uniformly randomly among all possible actions, namely, explore, and with probability  $1 - \epsilon$ , we choose the best action, namely, exploit. We do not want to continue exploring indefinitely but start exploiting once we do enough exploration; for this, we start with a high  $\epsilon$  value and gradually decrease it. We need to make sure that our policy is *soft*, that is, the probability of choosing any action  $a \in \mathcal{A}$  in state  $s \in S$  is greater than 0.

We can choose probabilistically, using the softmax function to convert values to probabilities

$$(18.10) \quad P(a|s) = \frac{\exp Q(s, a)}{\sum_{b \in \mathcal{A}} \exp Q(s, b)}$$

and then sample according to these probabilities. To gradually move from exploration to exploitation, we can use a “temperature” variable  $T$  and define the probability of choosing action  $a$  as

$$(18.11) \quad P(a|s) = \frac{\exp[Q(s, a)/T]}{\sum_{b \in \mathcal{A}} \exp[Q(s, b)/T]}$$

When  $T$  is large, all probabilities are equal and we have exploration. When  $T$  is small, better actions are favored. So the strategy is to start with a large  $T$  and decrease it gradually, a procedure named *annealing*, which in this case moves from exploration to exploitation smoothly in time.

### 18.5.2 Deterministic Rewards and Actions

In model-free learning, we first discuss the simpler deterministic case, where at any state-action pair, there is a single reward and next state possible. In this case, equation 18.7 reduces to

$$(18.12) \quad Q(s_t, a_t) = r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

and we simply use this as an assignment to update  $Q(s_t, a_t)$ . When in state  $s_t$ , we choose action  $a_t$  by one of the stochastic strategies we saw earlier, which returns a reward  $r_{t+1}$  and takes us to state  $s_{t+1}$ . We then update the value of *previous* action as

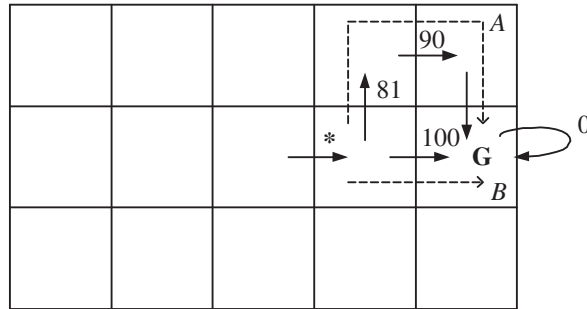
$$(18.13) \quad \hat{Q}(s_t, a_t) \leftarrow r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1})$$

where the hat denotes that the value is an estimate.  $\hat{Q}(s_{t+1}, a_{t+1})$  is a later value and has a higher chance of being correct. We discount this by  $\gamma$  and add the immediate reward (if any) and take this as the new estimate for the previous  $\hat{Q}(s_t, a_t)$ . This is called a *backup* because it can be viewed as taking the estimated value of an action in the next time step and “backing it up” to revise the estimate for the value of a current action.

BACKUP

For now we assume that all  $\hat{Q}(s, a)$  values are stored in a table; we will see later on how we can store this information more succinctly when  $|S|$  and  $|\mathcal{A}|$  are large.

Initially all  $\hat{Q}(s_t, a_t)$  are 0, and they are updated in time as a result of trial episodes. Let us say we have a sequence of moves and at each move, we use equation 18.13 to update the estimate of the  $Q$  value of the previous state-action pair using the  $Q$  value of the current state-action pair. In the intermediate states, all rewards and therefore values are 0, so no update is done. When we get to the goal state, we get the reward  $r$  and then we can update the  $Q$  value of the previous state-action pair as  $\gamma r$ . As for the preceding state-action pair, its immediate reward is 0 and the contribution from the next state-action pair is discounted by  $\gamma$  because it is one step later. Then in another episode, if we reach this



**Figure 18.4** Example to show that  $Q$  values increase but never decrease. This is a deterministic grid-world where  $G$  is the goal state with reward 100, all other immediate rewards are 0, and  $\gamma = 0.9$ . Let us consider the  $Q$  value of the transition marked by asterisk, and let us just consider only the two paths  $A$  and  $B$ . Let us say that path  $A$  is seen before path  $B$ , then we have  $\gamma \max(0, 81) = 72.9$ ; if afterward  $B$  is seen, a shorter path is found and the  $Q$  value becomes  $\gamma \max(100, 81) = 90$ . If  $B$  is seen before  $A$ , the  $Q$  value is  $\gamma \max(100, 0) = 90$ ; then when  $A$  is seen, it does not change because  $\gamma \max(100, 81) = 90$ .

state, we can update the one preceding that as  $\gamma^2 r$ , and so on. This way, after many episodes, this information is backed up to earlier state-action pairs.  $Q$  values increase until they reach their optimal values as we find paths with higher cumulative reward, for example, shorter paths, but they never decrease (see figure 18.4).

Note that we do not know the reward or next state functions here. They are part of the environment, and it is as if we query them when we explore. We are not modeling them either, though that is another possibility. We just accept them as given and learn directly the optimal policy through the estimated value function.

### 18.5.3 Nondeterministic Rewards and Actions

If the rewards and the result of actions are not deterministic, then we have a probability distribution for the reward  $p(r_{t+1}|s_t, a_t)$  from which rewards are sampled, and there is a probability distribution for the next state  $P(s_{t+1}|s_t, a_t)$ . These help us model the uncertainty in the system that may be due to forces we cannot control in the environment: for instance, our opponent in chess, the dice in backgammon, or our lack of

```

Initialize all  $Q(s, a)$  arbitrarily
For all episodes
  Initialize  $s$ 
  Repeat
    Choose  $a$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
    Take action  $a$ , observe  $r$  and  $s'$ 
    Update  $Q(s, a)$ :
       $Q(s, a) \leftarrow Q(s, a) + \eta(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$ 
     $s \leftarrow s'$ 
  Until  $s$  is terminal state

```

**Figure 18.5**  $Q$  learning, which is an off-policy temporal difference algorithm.

knowledge of the system. For example, we may have an imperfect robot which sometimes fails to go in the intended direction and deviates, or advances shorter or longer than expected.

In such a case, we have

$$(18.14) \quad Q(s_t, a_t) = E[r_{t+1}] + \gamma \sum_{s_{t+1}} P(s_{t+1} | s_t, a_t) \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

We cannot do a direct assignment in this case because for the same state and action, we may receive different rewards or move to different next states. What we do is keep a running average. This is known as the  $Q$  learning algorithm:

$Q$  LEARNING

$$(18.15) \quad \hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \eta(r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t))$$

We think of  $r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}(s_{t+1}, a_{t+1})$  values as a sample of instances for each  $(s_t, a_t)$  pair and we would like  $\hat{Q}(s_t, a_t)$  to converge to its mean. As usual  $\eta$  is gradually decreased in time for convergence, and it has been shown that this algorithm converges to the optimal  $Q^*$  values (Watkins and Dayan 1992). The pseudocode of the  $Q$  learning algorithm is given in figure 18.5.

We can also think of equation 18.15 as reducing the difference between the current  $Q$  value and the backed-up estimate, from one time step later. Such algorithms are called *temporal difference* (TD) algorithms (Sutton 1988).

TEMPORAL  
DIFFERENCE

OFF-POLICY  
ON-POLICY

This is an *off-policy* method as the value of the best next action is used without using the policy. In an *on-policy* method, the policy is used to

```

Initialize all  $Q(s, a)$  arbitrarily
For all episodes
  Initialize  $s$ 
  Choose  $a$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
  Repeat
    Take action  $a$ , observe  $r$  and  $s'$ 
    Choose  $a'$  using policy derived from  $Q$ , e.g.,  $\epsilon$ -greedy
    Update  $Q(s, a)$ :
       $Q(s, a) \leftarrow Q(s, a) + \eta(r + \gamma Q(s', a') - Q(s, a))$ 
     $s \leftarrow s', a \leftarrow a'$ 
  Until  $s$  is terminal state

```

**Figure 18.6** Sarsa algorithm, which is an on-policy version of  $Q$  learning.

SARSA determine also the next action. The on-policy version of  $Q$  learning is the *Sarsa* algorithm whose pseudocode is given in figure 18.6. We see that instead of looking for all possible next actions  $a'$  and choosing the best, the on-policy Sarsa uses the policy derived from  $Q$  values to choose one next action  $a'$  and uses its  $Q$  value to calculate the temporal difference. On-policy methods estimate the value of a policy while using it to take actions. In off-policy methods, these are separated, and the policy used to generate behavior, called the *behavior* policy, may in fact be different from the policy that is evaluated and improved, called the *estimation* policy.

Sarsa converges with probability 1 to the optimal policy and state-action values if a *GLIE policy* is employed to choose actions. A GLIE (greedy in the limit with infinite exploration) policy is where (1) all state-action pairs are visited an infinite number of times, and (2) the policy converges in the limit to the greedy policy (which can be arranged, e.g., with  $\epsilon$ -greedy policies by setting  $\epsilon = 1/t$ ).

TD LEARNING The same idea of temporal difference can also be used to learn  $V(s)$  values, instead of  $Q(s, a)$ . *TD learning* (Sutton 1988) uses the following update rule to update a state value:

$$(18.16) \quad V(s_t) \leftarrow V(s_t) + \eta[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

This again is the delta rule where  $r_{t+1} + \gamma V(s_{t+1})$  is the better, later prediction and  $V(s_t)$  is the current estimate. Their difference is the temporal difference, and the update is done to decrease this difference. The

update factor  $\eta$  is gradually decreased, and TD is guaranteed to converge to the optimal value function  $V^*(s)$ .

#### 18.5.4 Eligibility Traces

ELIGIBILITY TRACE

The previous algorithms are one-step—that is, the temporal difference is used to update only the previous value (of the state or state-action pair). An *eligibility trace* is a record of the occurrence of past visits that enables us to implement temporal credit assignment, allowing us to update the values of previously occurring visits as well. We discuss how this is done with Sarsa to learn  $Q$  values; adapting this to learn  $V$  values is straightforward.

To store the eligibility trace, we require an additional memory variable associated with each state-action pair,  $e(s, a)$ , initialized to 0. When the state-action pair  $(s, a)$  is visited, namely, when we take action  $a$  in state  $s$ , its eligibility is set to 1; the eligibilities of all other state-action pairs are multiplied by  $\gamma\lambda$ .  $0 \leq \lambda \leq 1$  is the trace decay parameter.

$$(18.17) \quad e_t(s, a) = \begin{cases} 1 & \text{if } s = s_t \text{ and } a = a_t, \\ \gamma\lambda e_{t-1}(s, a) & \text{otherwise} \end{cases}$$

If a state-action pair has never been visited, its eligibility remains 0; if it has been, as time passes and other state-actions are visited, its eligibility decays depending on the value of  $\gamma$  and  $\lambda$  (see figure 18.7).

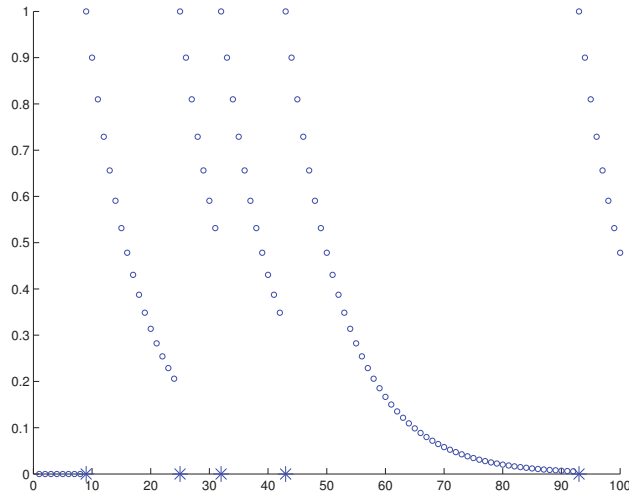
We remember that in Sarsa, the temporal error at time  $t$  is

$$(18.18) \quad \delta_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

In Sarsa with an eligibility trace, named Sarsa( $\lambda$ ), *all* state-action pairs are updated as

$$(18.19) \quad Q(s, a) \leftarrow Q(s, a) + \eta \delta_t e_t(s, a), \quad \forall s, a$$

This updates all eligible state-action pairs, where the update depends on how far they have occurred in the past. The value of  $\lambda$  defines the temporal credit: If  $\lambda = 0$ , only a one-step update is done. The algorithms we discussed in section 18.5.3 are such, and for this reason they are named  $Q(0)$ , Sarsa(0), or TD(0). As  $\lambda$  gets closer to 1, more of the previous steps are considered. When  $\lambda = 1$ , all previous steps are updated and the credit given to them falls only by  $\gamma$  per step. In online updating, all eligible values are updated immediately after each step; in offline updating, the updates are accumulated and a single update is done at



**Figure 18.7** Example of an eligibility trace for a value. Visits are marked by an asterisk.

the end of the episode. Online updating takes more time but converges faster. The pseudocode for *Sarsa*( $\lambda$ ) is given in figure 18.8.  $Q(\lambda)$  and TD( $\lambda$ ) algorithms can similarly be derived (Sutton and Barto 1998).

## 18.6 Generalization

Until now, we assumed that the  $Q(s, a)$  values (or  $V(s)$ , if we are estimating values of states) are stored in a lookup table, and the algorithms we considered earlier are called *tabular* algorithms. There are a number of problems with this approach: (1) when the number of states and the number of actions is large, the size of the table may become quite large; (2) states and actions may be continuous, for example, turning the steering wheel by a certain angle, and to use a table, they should be discretized which may cause error; and (3) when the search space is large, too many episodes may be needed to fill in all the entries of the table with acceptable accuracy.

Instead of storing the  $Q$  values as they are, we can consider this a regression problem. This is a supervised learning problem where we define a regressor  $Q(s, a|\theta)$ , taking  $s$  and  $a$  as inputs and parameterized by a