

11 Reinforcement Learning

11.1 Learning from reward and the credit assignment problem

11.1.1 The reinforcement learning problem

In supervised learning we assumed that a teacher supplies detailed information on the desired response of a learner. This was particularly suited to object recognition where we had a large number of labeled examples. However, there are different learning circumstances which seem to be much more common when agents have to learn from acting in the environment. With agent we refer here to a system that acts in the world, like a robot or a human. An example of a learning task for an agent is to learn to play tennis. Such learning comprises trying out moves and getting rewarded by points the agent score rather than a teacher who specifies every muscle movement we need to follow or an engineer who designs every sequence of motor activations. It is indeed common in such situations that the agent only get some simple feedback after long periods of actions in form of reward or punishment without even detailing which of the actions has contributed to the outcome. In this type of learning scenario is called **reinforcement learning (RL)**.

RL faces several challenges. One is called the **credit assignment** problem. This includes which action (**spatial credit assignment**) and at which time (**temporal credit assignment**) of the system should be given credit for the achievement of reward. Another important aspect that is new in contrast to supervised learning is that the agent must search for solutions by trying different actions. The agent must therefore generally play an active role in exploring options. And what if we we find a solution that give us some reward? We could then ask if this is a good solution or if the agent should search for a better solution and forgo some of the known rewards? This problem is commonly stated as **exploration versus exploitation trade-off**.

Learning with reward signals has been studied by psychologists for many years under the term **conditioning**. An example of classical conditioning in animal learning is shown in Fig. 11.1.1. In the illustrated experiment, we place a rodent in a T-maze and supply food of different sizes when the rodent goes to the end of each horizontal arm of the T-maze. The rodent might wander around. Let us assume that it found the smaller food reward at the end of the left arm of the T-maze. It is then likely that the rodent will turn left in subsequent trials to receive food reward. Thus the animal learned that the action of taking a left turn and going to the end of the arm is associated with food reward. Of course, in this case the rodent could also receive larger reward when exploring the right arm of the maze. Similar settings can be applied to robots, for example robots that must learn to find a charging station. Even our original object recognition problems of naming the correct class to which an object belongs

can be formulated in terms of reinforcement. Psychologists also distinguish between **instrumental conditioning**, in which an action and hence a decision has to be made in order to find positive reward. In contrast, in **classical conditioning** does not require an action, and such settings are ideal to illuminate the temporal order in which reward associations are made.

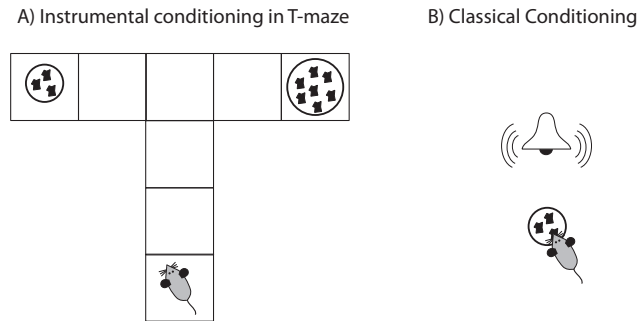


Fig. 11.1 Example of instrumental and classical conditioning. A) A rodent has to learn to transverse the maze and make a decision at the junction in which direction to go. Such as decision problem, which necessitates the action of an actor, is called instrumental conditioning in the animal learning literature. B) A slightly simpler setting is that of classical conditioning which does not require an action and thus concentrates on learning the reward associations. An typical example is when a subject is required to associate the ringing of a bell with reward.

11.1.2 Formalization of the problem setting: The Markov Decision Process

To discuss RL in more depth we need to formalize the reinforcement setting a bit more. We consider an agent that in each time interval t is in a specific state s_t . A state describes thereby the environment such as a location at which the agent could be. Furthermore, we assume that the agent can take an action a_t from each state. This action specifies a transition to a new state, and this transition is specified by the transition function τ as

$$\text{Transition function: } s_{t+1} = \tau(s_t, a_t). \quad (11.1)$$

We restrict the discussion here to the common assumption that the transition function only depend on the previous state and the intended action from the corresponding state. This is called the **Markov condition**, and the corresponding decision process is called a **Markov Decision Process (MDP)**. In contrast, a non-Markovian condition would be the case in which the next state depends on a series of previous states and actions, and our agent would then need a memory to make optimal decisions. The situation described by the Markov condition is quite natural as many decisions processes only depend on the current state, and we could even define a state as having all necessary conditions to make the decision. a MDP is therefore a good scenario to solve. We also assume that the agent knows in which state it is in. This is often a problematic assumptions with limited sensors such as human perception or in robotics. However,

the MDP problem can also be generalized to situation of a partially observable Markov decision process, usually called an POMDP.

Next, we assume that the environment or a teacher provides reward according to a reward function ρ ,

$$\text{Reward function: } r_{t+1} = \rho(s_t, a_t). \tag{11.2}$$

This reward functions returns the value of reward when the agent is in state s_t and takes intended action a_t . In the deterministic case this is of course the reward at the next state $s_{t+1} = \tau(s_t, a_t)$.

Finally, an important function in reinforcement learning is the control **policy** that specifies which action a to take from each state,

$$\text{Policy: } a_t = \pi(s_t). \tag{11.3}$$

In the context of a mobile agent, the action a_t is commonly provided by a motor command that specifies the action that the agent should take at time t , the time when the agent is in state s_t . For example, a muscle that should move the arm should be activated with a certain nerve pulse, or a robot that should turn is activated by a certain motor that should run for a certain time with a certain speed.

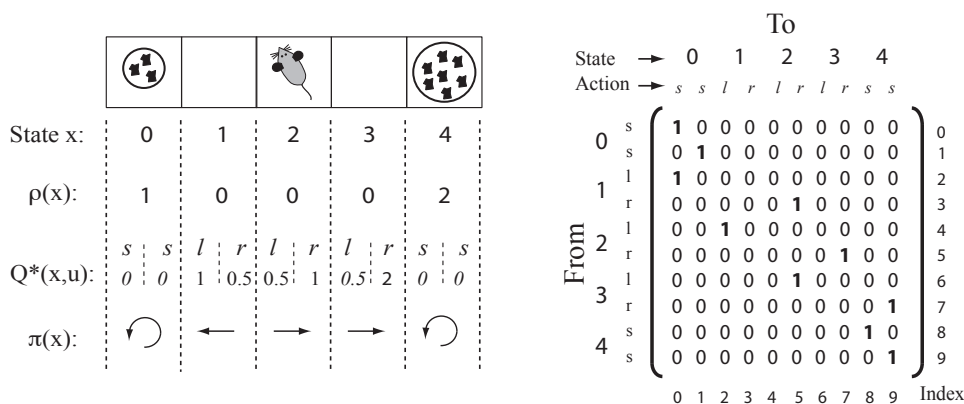


Fig. 11.2 Example experiment with the simplified T-maze where we concentrate on the more interesting horizontal portion of the maze (linear maze). The right hand side shows the corresponding transition matrix for the optimal policy.

To illustrate the different reinforcement schemes discussed in this chapter, we will apply these to the example of the T-maze outlined in the introduction. To keep the programs minimal and clean, we concentrate on the upper linear part for the maze as illustrated in Fig. 11.2. States s of the maze are labeled as 0 to 4. A reward of value 1 is provided in state 0 and a reward of 2 is provided in state 4. The discrete Q -function has 10 values corresponding to each possible action in each state. In states 1, 2, and 3 these are the actions of move *left* or move *right*. The states with the reward, states 0 and 4, are terminal states and the agent would stay in these states. We coded this with two *stay* actions to keep some consistency in the representation.

We provide here several python functions for later use. First there is the transition function $\tau(s, a)$ and the reward function $\rho(s, a)$. Next we provide a function to calculate

the policy from a value function. This will be discussed further below. Finally, we provide a helper function $\text{idx}(a)$ to transform the action representation $u \in \{-1, 1\}$ to the corresponding indices $\text{idx} \in \{0, 1\}$:

```
## Reinforcement learning in a maze
import numpy as np
import matplotlib.pyplot as plt

def tau(s, a):
    if s==0 or s==4: return(s)
    else: return(s+a)
def rho(s, a):
    return(s==1 and a== -1)+2*(s==3 and a==1)
def calc_policy(Q):
    policy=np.zeros(5)
    for s in range(0,5):
        action_idx=np.argmax(Q[s, :])
        policy[s]=2*action_idx -1
        policy[0]=policy[4]=0
    return policy.astype(int)
def idx(a):
    return(int((a+1)/2))
gamma=0.5;
```

11.1.3 Return and Value functions

The goal of the agent is to maximize the **total expected reward** in the future from every initial state. This quantity is called **return** in economics. Of course, if we assume that this goes on forever then this return should be infinite, and we have hence to be a bit more careful. One common choice is to define the return as the average reward in a finite time interval, also called the **finite horizon** case. Another common form to keep the return finite is to use a **discounted return** in which an agent values immediate reward more than reward far in the future. To capture this we define a discount factor $0 < \gamma < 1$. In the example program we will use a value of $\text{gamma} = 0.5$,

```
## discount factor
gamma=0.5;
```

although values much closer to one such as $\gamma = 0.99$ are common. For this case we now define a **state-action value function** which gives us a numerical value of the return (all future discounted reward) when the agent is in state s and takes action a and then follows the policy for the following actions,

$$\text{Value function (state-action): } Q^\pi(s, a) = \rho(s, a) + \sum_{t=1}^{\infty} \gamma^t \rho(s_t, \pi(s_t)) \quad (11.4)$$

In other words, this function tells us how good is action a in state s , and the knowledge of this value function should hence guide the actions of an agent. The aim of value-based reinforcement learning is to estimate this function.

Sometimes we are only interested in the value function when we follow the policy even for the first step in the action sequence from state s . This value function does then not depend explicitly on the action, only indirectly of course on the policy, and is defined as The total discounted return from state $s = s_0$ following policy π is,

$$\text{Value function (state): } V^\pi(s) = Q^\pi(s, \pi(s)) \quad (11.5)$$

The goal of RL is to find the policy which maximized the return. If the agents knows the

$$\text{Optimal Value function: } V^*(s) = \max_a Q^*(s, a), \quad (11.6)$$

then the optimal policy is simply given by taking the action that leads to the biggest expected return, namely

$$\text{Optimal policy: } \pi^*(s) = \arg \max_a Q^*(s, a). \quad (11.7)$$

This function is implemented above with the python code for `calc_policy`.

The optimal value function and the optimal policy are closely related. We will discuss in the following several methods to calculate or estimate the value function from which the policy can be derived. These methods can be put under the heading of **value-search**. Corresponding agents, or part of the corresponding RL algorithms, are commonly called a **critic**. However, there are also methods to learn the policy directly. Such methods are called **policy-search** is the corresponding agents are called an **actor**. At the end we will argue that combining these approaches in an **actor-critic scheme** has attractive creatures, and such schemes are increasingly used in practical applications.

11.2 Model-based Reinforcement Learning

In this section we assume that the agent has a **model** of the environment and its behaviour by knowing the reward function $\rho(s, a)$ and the transfer functions $\tau(s, a)$. The knowledge of these functions, or a model thereof, is required for **model-based RL**, which is also called **dynamic programming**.

11.2.1 The basic Bellman equation

The key to learning the value functions is the realization that the right hand side of eq.?? can be written in terms of the Q-function itself, namely

$$\begin{aligned} Q^\pi(s, a) &= \rho(s, a) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} \rho(s_t, \pi(s_t)) \\ &= \rho(s, a) + \gamma \left[\rho(\tau(s, a), \pi(\tau(s, a))) + \gamma \sum_{t=2}^{\infty} \gamma^{t-2} \rho(s_t, \pi(s_t)) \right]. \end{aligned}$$

The term in the square bracket is equal to the value function of the state that is reached after the transition $\tau(s, a)$

$$Q^\pi(\tau(s, a), \pi(\tau(s, a))) = V^\pi(\tau(s, a)). \quad (11.8)$$

The Q -function and the V -function are here equivalent since we are following the policy in these steps. Using this fact in the equation above we get the

$$\pi \text{ Bellman equation: } Q^\pi(s, a) = \rho(s, a) + \gamma Q^\pi(\tau(s, a), \pi(\tau(s, a))). \quad (11.9)$$

If we combine this with known dynamic equations in the continuous time domain, then this becomes the Hamilton-Jacobi-Bellman equation, often encountered in engineering.

As stated above, we assume here the reward function $\rho(s, a)$ and the transition functions $\tau(s, a)$ are known. At this point the agent follows a specified policy $\pi(s)$. Let us further assume that we have n_s states and n_a possible actions in each state. We have thus $n_s \times n_a$ unknown quantities $Q^\pi(s, a)$ which are governed by the Bellman equation above. More precisely, the Bellman equations 11.9 are $n_s \times n_a$ coupled linear equations of the unknowns $Q^\pi(s, a)$. It is then convenient to write this equation system with vectors

$$\mathbf{Q}^\pi = \mathbf{R} + \gamma \mathbf{T}^\pi \mathbf{Q}^\pi \quad (11.10)$$

Where T^π is an appropriate transition matrix which depends on the policy. This equation can also be written as

$$\mathbf{R} = (\mathbb{1} - \gamma \mathbf{T}^\pi) \mathbf{Q}^\pi, \quad (11.11)$$

where $\mathbb{1}$ is the identity matrix. This equation has the solution

$$\mathbf{Q}^\pi = (\mathbb{1} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{R} \quad (11.12)$$

if the inverse exists. In other words, as long as the agent knows the reward function and the transition function, it can calculate the value function for a specific policy without even taking a single step. This is an example of a deliberative system where the agent can use the models of reward and the environment to calculate optimal decisions; hence the specification of model-based RL.

To demonstrate how to solve the Bellman equation with linear algebra tools, we need to define the corresponding vectors and matrices as used in eq.11.12. We order therefore the quantities such as ρ and Q with ten indices. The first one correspond to $(s = 0, u = -1)$, the second to $(s = 0, u = 1)$, the third to $(s = 1, u = -1)$, etc. The reward vector can thus be coded as:

```
print('—> Analytic solution for optimal policy')
```

```
# Defining reward vector R
i=0; R=np.zeros(10)
for s in range(0,5):
    for a in range(-1,2,2):
        R[i]=rho(s,a)
        i += 1
```

The transition matrix depends on the policy, so we need to choose one. We chose the one specified on the left in Fig. 11.2 where the agent would move to the left in state $s = 2$ and to the right in states $s = 3$ and $s = 4$. This happens to be the optimal

solution as we will show later so that this will also give us a solution for the optimal value function. We use this policy to construct the transition matrix by hand as shown on the right in Fig. 11.2. For example, if we are in state $s = 4$ and move to the left, $a = -1$, corresponding to the from-index=7, then we end up in state $s = 3$, from which the policy say go right, $a = 1$. This correspond to the to-index=6. Thus, the transition matrix should have an entry $T(7, 6) = 1$. Going through all the cases results in

```
# Defining transition matrix
T=np.zeros([10,10]);
T[0,0]=1; T[1,1]=1; T[2,0]=1; T[3,5]=1; T[4,2]=1
T[5,7]=1; T[6,5]=1; T[7,9]=1; T[8,7]=1; T[9,9]=1
```

With this we can solve this linear matrix equations with the `inv()` function in the linear algebra package of numpy,

```
# Calculate Q-function
Q=np.linalg.inv(np.eye(10)-gamma*T) @ np.transpose(R)
Q=np.transpose(np.reshape(Q,[5,2]))
```

We reshaped the resulting Q -function so that the first row shows the values for left movements in each state and the second row shows the values for a right movement in each state. From this we can calculate which movement to take in each state, namely just the action corresponding to the maximum value in each column and print the results,

```
policy=calc_policy(Q)
print('Q values: \n', np.transpose(Q))
print('policy: \n', np.transpose(policy))
```

which gives

```
—> Analytic solution for optimal policy
Q values:
[[ 0.   1.   0.5  0.5  1. ]
 [ 0.   0.5  1.   2.   0. ]]
policy:
[ 0 -1  1  1  0]
```

That is ignoring the end states it is moving left in state 1 as this would lead to an immediate reward of 1 and moving right in the other states as this would result in a larger reward even when taking the discounting for more steps into account. Of course, at this point our argument is circular as we have started already with the assumption that we use the optimal policy as specified in the transition matrix to start with. We will soon see how to start with an arbitrary policy an improve this to find the optimal strategy. Also note that the transition matrix was perfect in the sense that the intended move always leads to the intended end state. We will later see that a probabilistic extension of this transition matrix is quite useful in describing more realistic situations.

In the code we save the optimal Q -values for the optimal policy

```
Qana=Q
```

so that we can later plot the difference of other solution methods.

The Bellman equations are a set of n coupled linear equations for n unknown Q values, and we have solved these here with linear algebra function to find an inverse of a matrix. This can be implemented with some algorithms such as Gauss elimination. Alternatively we can solve the Bellman equation with the following iterative procedure. We starts with a guess of the Q -function, let's call this Q_i^π , and improve it by calculating the right hand side of the Bellman equation,

$$\text{Dynamic Programming: } Q_{i+1}^\pi(s, a) \leftarrow \rho(s, a) + \gamma Q_i^\pi(\tau(s, a), \pi(\tau(s, a))). \quad (11.13)$$

The fixed-point of this equation, that is, the values that does not change with these iterations, are the desired values of Q^π . Another way of thinking about this algorithm is that the Bellman equality is only true for the correct Q^π values. For our guess, the difference between the left- and right-hand side is not zero, but we are minimizing this with the iterative procedure above. The corresponding code is

```
print ( '\n--> Dynamic Programming ' )

Q=np. zeros ([5 ,2])
for iter in range(3):
    for s in range(0 ,5):
        for a in range(-1 ,2 ,2):
            act = np. int ( policy [ tau ( s , a ) ] )
            Q[ s , idx ( a ) ]=rho ( s , a )+gamma*Q[ tau ( s , a ) , idx ( act ) ]

print ( 'Q values : \n' , np. transpose ( Q ) )
print ( ' policy : \n' , np. transpose ( policy ) )
```

This is a much more common implementation and it does not require the explicit coding of the transition matrix. Iterative approaches will be used in all further methods discussed below. Note that we have only used three iterations to converge on the right solution. While we set here the number of iterations by hand, in practice we iterate until the changes in the values are sufficiently small.

11.2.2 Policy Iteration

The goal of RL is of course to find the policy which maximized the return. So far we have only discussed a method to calculate the value for a given policy. However, we can start with an arbitrary policy and can use the corresponding value function to improve the policy by defining a new policy which is given by taken the actions from each state that gives us the best next return value,

$$\text{Policy iteration: } \pi(s) \leftarrow \arg \max_a Q^\pi(s, a). \quad (11.14)$$

For the new policy we can then calculate the corresponding Q -function and then use this Q -function to improve the policy again. Iterating over the policy gives us the

$$\text{Optimal policy: } \pi^*(s). \quad (11.15)$$

The corresponding value function is Q^* . In the maze example we can see that the maximum in each column of the Q-matrix is the policy we started with. This is hence the optimal policy as we stated before.

The corresponding code for our maze example is

```
print( '\n--> Policy iteration ' )

Q=np.zeros([5,2])
policy=calc_policy(Q)
for iter in range(3):
    for s in range(0,5):
        for a in range(-1,2,2):
            act = np.int(policy[tau(s,a)])
            Q[s,idx(a)]=rho(s,a)+gamma*Q[tau(s,a),idx(act)]
policy=calc_policy(Q)

print( 'Q values: \n', np.transpose(Q) )
print( 'policy: \n', np.transpose(policy) )
```

Note that in this example we iterated again only three times over the policies. In principle we could and should iterate several times for each policy in order to converge to a stable estimate for this Q^π . However, the improvements will anyhow lead very quickly to a stable state, at least in this simple example.

11.2.3 Bellman function for optimal policy and value iteration

Since we are foremost interested in the optimal policy, we could try to solve the Bellman equation right away for the optimal policy,

$$Q^*(s,a) = \rho(s,a) + \gamma Q^*(\tau(s,a), \pi^*(\tau(s,a))). \quad (11.16)$$

The problem is that this equation does now depend on the unknown π^* . However, we can check in each state all the actions and take the one which gives us the best return. This should be equivalent to the equation above in the optimal case. Hence we propose the

Optimal Bellman equation: $Q^*(s,a) = \rho(s,a) + \gamma \max_{a'} Q^*(\tau(s,a), a').$ (11.17)

We can solve this with dynamic programming when the transfer function and the reward functions are known,

Q-iteration: $Q^*(s,a) \leftarrow \rho(s,a) + \gamma \max_{a'} Q^*(\tau(s,a), a').$ (11.18)

The corresponding code for our maze example is

```
print( '\n--> Q-iteration ' )

Q_new=np.zeros([5,2])
Q=np.zeros([5,2])
```

```

policy = np.zeros(5)
for iter in range(2):
    for s in range(0,5):
        for a in range(-1,2,2):
            maxV = np.maximum(Q[tau(s,a),0],Q[tau(s,a),1])
            Q_new[s,idx(a)] = rho(s,a) + gamma*maxV
        Q = np.copy(Q_new)

print('Q values: \n', np.transpose(Q))
print('policy: \n', np.transpose(policy))

```

In this example we again used 3 iterations which are sufficient to reach the correct values. In practice we would terminate the program if the changes are sufficiently small.

11.3 Model-free Reinforcement Learning

11.3.1 Temporal Difference Method for Value Iteration

Above we assumed a model of the environment by an explicit knowledge of the functions $\tau(s, a)$ and $\rho(s, a)$. We could use modelling techniques and some sampling strategies to learn these functions explicitly and then use model-based RL as described above to find the optimal policy. Instead we will here combine here the sampling by exploring the environment directly with reinforcement learning.

We will start again with a version for a specific policy by choosing a policy and estimate the Q -function for this policy. As in dynamic programming we want to minimize the difference between the left- and right-hand side of the Bellman equation. But we can not calculate the right hand side since we do not know the transition function and the reward function. However, we can just take a step according to our policy $u = \pi(s)$ and observe a reward r_{i+1} and the next state s_{i+1} . Now, since this is only one sample we should take this only with a small learning rate α into account to update the value function

$$\text{SARSA: } Q_{i+1}(s_i, a_i) = Q_i(s_i, a_i) + \alpha [(r_{i+1} + \gamma Q_i(s_{i+1}, a_{i+1}) - Q_i(s_i, a_i))]. \quad (11.19)$$

The term in the square brackets on the right hand-side is called the **temporal difference**. Note that we are following here the policy, and the method is therefore often labeled as **on-policy**. The next step is to use the estimate of the Q -function to improve policy. However, since we are anyhow mainly interested in the optimal policy we should improve the policy by taking the steps that maximizes the payoff. However, one problem in this scheme is that we have to estimate the Q values by sampling so that we have to make sure to trade off **exploitation** with **exploration**. A common way to choose the policy in this scheme is the

$$\epsilon\text{-greedy policy: } p\left(\arg \max_a Q(s, a)\right) = 1 - \epsilon. \quad (11.20)$$

So, we are really evaluating the optimal policy that requires to make the exploration zero at the end, $\epsilon \rightarrow 0$.

The corresponding code for the SARSA is

```

print( '\n--> SARSA' )

Q=np.zeros([5,2])
error = []
alpha=1;
for trial in range(500):
    policy=calc_policy(Q)
    s=2
    for t in range(0,5):
        a=policy[s]
        if np.random.rand() < 0.1: a=-a #epsilon greedy
        a2=idx(policy[tau(s,a)])
        Q[s,idx(a)]=(1-alpha)*Q[s,idx(a)]+alpha*(rho(s,a)+gamma*Q[tau(s,a),a2])
        s=tau(s,a)
    error.append(np.sum(np.sum(np.abs(np.subtract(Q,Qana))))))

print( 'Q values: \n', np.transpose(Q))
print( 'policy: \n', np.transpose(policy))
plt.figure(); plt.plot(error); plt.show()

```

Note that we have set here the learning rate $\alpha = 1$ for this example, which makes the update rule look similar to dynamic programming

$$Q_{i+1}(s_i, a_i) = (r_{i+1} + \gamma Q_i(s_{i+1}, a_{i+1})). \quad (11.21)$$

However, there is a major difference. In dynamic programming we iterate over all possible states with the knowledge of the transition function and the reward function. Thus an agent does not really have to explore the environment and can just "sit there" and calculate what the optimal action is. This is the benefit of model-based reinforcement learning. In contrast, here we discuss the case where we do not know the transition function and the reward function and hence have to explore the environment by acting in it. As there is usually associated with a physical movement, this takes time and hence limits the exploration we can do. Hence it is common that an agent can not explore all possible states. Also, a learning rate of $\alpha = 1$ is not always advisable since a more common setting is that reward itself is probabilistic. A smaller value of α then represents a form of taking a sliding average and hence estimating the expected value of the reward.

It is illustrative to go through the SARSA algorithm by hand for our linear-maze example. An example is shown in Fig. ?? . In this example we changed the situation to a linear maze in which the state always starts at the leftmost state and a reward of $r = 1$ is received in the rightmost state. The policy is to always go right, which is also the optimal policy in this situation. At the first time step of the first episode we are in the leftmost state and evaluate the value of going right. In the corresponding state to the right there is no reward given, and the value function is also zero. So the value function of this state-action is zero. The same is true for every step until we are in the state before the reward state. At this point the value is updated to the reward of the next

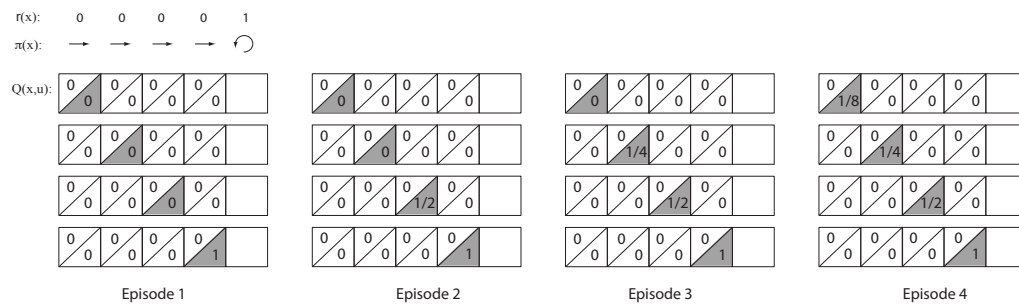


Fig. 11.3 Example of the "back-propagation" of the reward (not to be confused with the back propagation algorithm in supervised learning). In this example, an episode always starts in the leftmost state and the policy is to always go right. A reward is received in the rightmost state.

state. In the second episode the value of the first and second state remains zero, but the third state is updated to $\gamma * 1$ since the value of the next state following the policy is given by one, and we discount this by γ . Going through more episodes it can be seen that the value "back-propagates" by one step in each episode. Note that this back propagation of the value is not to be confused with the back propagation algorithm in supervised learning. Also, notice that the values for the Q -function for going left are not updated as we only followed optimal policy deterministically. Some exploration steps will eventually update these values, although it might take a long time until these values propagate through the system.

Finally, a more common version of a temporal difference learning is to use an **off-policy** approach for the estimation step from each visited state. That is we check all possible actions from the state that we evaluate state and update the value function with the maximal expected return,

$$\mathbf{Q\text{-learning:}} \quad Q_{i+1}(s_i, a_i) = Q_i(s_i, a_i) + \alpha \left[(r_{i+1} + \gamma \max_{a'} Q_i(s_{i+1}, a') - Q_i(s_i, a_i)) \right]. \quad (11.22)$$

We still have to explore the environment which is again usually following the optimal estimated policy with some allowance for exploration such as ϵ -greedy or a softmax exploration strategy.

The corresponding code for the Q learning is

```

print('\n--> Q-Learning:')

Q=np.zeros([5,2])
alpha=1
error = []
for trial in range(500):
    s=2
    for t in range(0,5):
        action_idx=np.argmax(Q[s,:])
        a=2*action_idx-1
        if np.random.rand()<0.1: a=-a #epsilon greedy
  
```

```

Q[s, idx(a)]=(1-alpha)*Q[s, idx(a)]+alpha*(rho(s,a)+gamma*np.maximum(Q[tau(s,a)]-Q[s, idx(a)], 0))
Q[0]=0;Q[4]=0
s=tau(s,a)
error.append(np.sum(np.sum(np.abs(np.subtract(Q,Qana))))))

print('Q values: \n',np.transpose(Q))
print('policy: \n',np.transpose(policy))
plt.plot(error,'r'); plt.show()

```

11.3.2 TD(λ)

The example of the linear maze in the previous section has shown that the expectation of reward propagates backwards in each episode which hence requires multiple repetition of the episodes in order to evaluate the value function. The reason for this is that we only give credit for making a step to a valuable state to the previous step and hence only update the corresponding value function. A different approach is to keep track of which states have led to the reward and assign the credit to each step that was visited. However, because we discount the reward proportional to the time it takes to get to the rewarded state, we need to also take this into account.

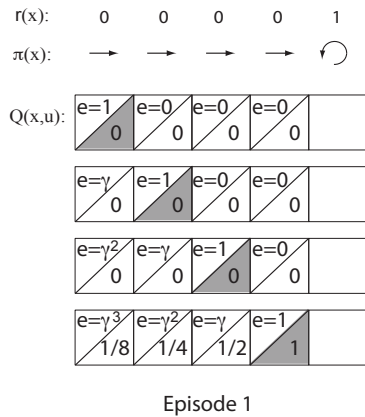


Fig. 11.4 Example of the "back-propagation" of the reward with an eligibility trace e . In this example, an episode always starts in the leftmost state and the policy is to always go right. A reward is received in the rightmost state.

To realize this we introduce an eligibility trace that we call $e(s)$. At the beginning, this eligibility trace is set to zero for all the states. The for each visited state we set this eligibility state to one for this current state, and we discount the eligibility for all the other states by γ . This is demonstrated in Fig. 11.4. In the figure we reused the place for the $Q(s, -1)$ to indicate the eligibility trace at every time step. Note how the eligibility trace is building up during one episode until reaching this the rewarded state, at which time the values of all the states are updated in the right proportion. This algorithm does therefore only need one optimal episode to find the correct value function, at least in this case with a learning rate of $\alpha = 1$. This algorithm is implemented for the

Q-learning version of temporal difference learning in the code below:

```

print ( '\n—> TD(lambda) for Q-Learning: ')

Q=np.zeros([5,2])
alpha=1
lam=0.1
error = []

for trial in range(100):
    s=2; eligibility=np.zeros(5)
    for t in range(0,5):
        if s==0 or s==4: break
        eligibility*=gamma*lam
        eligibility[s]=1
        action_idx=np.argmax(Q[s,:])
        a=2*action_idx-1
        if np.random.rand()<0.1: a=-a #epsilon greedy

        for x in range(1,4):
            Q[x,idx(a)]=Q[x,idx(a)]+alpha*(rho(x,a)+gamma*np.maximum(Q[tau(x,a),0]
            s=tau(s,a)
            error.append(np.sum(np.sum(np.abs(np.subtract(Q,Qana))))))

print ('Q-values: \n',np.transpose(Q))
print ('policy: \n',np.transpose(policy))
plt.plot(error,'r'); plt.show()

```

While this algorithms does cut down the number of episodes to learn the value function, it does so on the expense of keeping memories of visited states and their respective times. A compromise is to allow some decaying memory in the algorithm. This can be implemented simply by a factor λ ind replacing the update of the eligibility trace from

```
eligibility*=gamma
```

to

```
eligibility*=gamma*lambda
```

With a term $\lambda < 1$, this corresponds to an exponential decay of the eligibility trace and hence the necessary memory. We will see later that this can be implemented efficiently in neural networks,