

# CS181 Practice Questions: More MDPs and Reinforcement Learning

## 1. Q and V practice

Consider an MDP with two states (A and B) and two actions ( $\square$  and  $\triangle$ ). You have the following transition and reward functions with  $\gamma = 0.5$ .

$s$	$a$	$s'$	$P(s'   a, s)$		$s$	$a$	$R(s, a)$
A	$\square$	A	0.5		A	$\square$	3
A	$\square$	B	0.5		A	$\triangle$	1
A	$\triangle$	A	1		B	$\square$	0
A	$\triangle$	B	0		B	$\triangle$	1
B	$\square$	A	1				
B	$\square$	B	0				
B	$\triangle$	A	1				
B	$\triangle$	B	0				

Complete the following table. (Hint: It's pretty easy to set up a spreadsheet to propagate the information for value iteration.)

$k$ To Go	$Q(A, \square)$	$Q(A, \triangle)$	$Q(B, \square)$	$Q(B, \triangle)$	$\pi_k(A)$	$\pi_k(B)$	$V_k(A)$	$V_k(B)$
0	—	—	—	—	—	—	0	0
1	3.00	1.00	0.00	1.00	$\square$	$\triangle$	3.00	1.00
2	4.00	2.50	1.50	2.50	$\square$	$\triangle$	4.00	2.50
3	4.63	3.00	2.00	3.00	$\square$	$\triangle$	4.63	3.00
4	4.91	3.31	2.31	3.31	$\square$	$\triangle$	4.91	3.31
5	5.05	3.45	2.45	3.45	$\square$	$\triangle$	5.05	3.45
6	5.13	3.53	2.53	3.53	$\square$	$\triangle$	5.13	3.53

## 2. Stationary Policy

In the previous problem, is the policy stationary?
--

In general  $k$ -to-go policies may not be stationary, as it may make sense to take different actions closer to the horizon. However, in this case, the policy is stationary as it always makes sense to try to switch to state  $A$ .

### 3. Value Iteration for Survival (MIT 6.034)

Consider the following MDP:

Start State	Action	Probability	Next State
hungry, tired	rested	1.0	hungry, rested
	hunt	0.8	hungry, tired
		0.2	full, tired
hungry, rested	rest	1.0	hungry, rested
	hunt	0.8	full, rested
		0.2	hungry, tired
full, tired	rest	1.0	hungry, rested
	hunt	0.8	hungry, tired
		0.2	full, tired
full, rested	rest	1.0	hungry, rested
	hunt	0.8	full, tired
		0.2	hungry, tired

In this world, you care about being rested and well fed. You can either rest or hunt. Your reward function is

State	Reward
hungry, tired	0
hungry, rested	1
full, tired	1
full, rested	2

- Given a discounting constant  $\gamma$ , what is the value of the (hungry, tired) state for a policy that always rests? Give your answer in terms of  $\gamma$ . (Hint: use geometric series.)
- Assume that  $Q((\text{hungry, tired}), \text{rest}) = a$  and  $Q((\text{hungry, tired}), \text{hunt}) = b$  and that  $b > a$ . Write an expression for  $Q((\text{hungry, rested}), \text{hunt})$  in terms of these quantities,  $\gamma$ , and  $V^*(\text{full, rested})$ .

- Constantly resting causes one to go to (hungry, rested) and stick there. This has reward 1 at every time step, discounted by gamma, so the value is:

$$\sum_{t=1}^{\infty} \gamma^t = \frac{1}{1-\gamma} - 1$$

- Writing Bellman's equation:

$$Q((\text{hungry, rested}), \text{hunt}) = 1 + \gamma \times (0.8 \times V^*(\text{full, rested}) + 0.2 \times b) \quad (1)$$

#### 4. Q Learning

Imagine that you are taking actions in an MDP that are selected according to the optimal policy. Explain whether or not Q-learning will only learn optimal Q-values.

Q-learning will learn the optimal Q-values as long as the policy explores all states. The optimal policy may not necessarily do that, however.

## 5. Q Learning

Suppose you are standing on a linear board on which you can take action  $L$  or  $R$  to walk left or right, or  $S$  to sleep. If you walk, you have probability  $p_a$  that you actually walk to the next square, where  $a \in \{L, R\}$ . With probability  $1 - p_a$ , however, your cat distracts you and you stay on the same square. Staying on a square gives you reward  $r_i$ . Your learning rate is  $\alpha = 0.5$  and  $\gamma = 0.5$ .

You are on square 1, you choose  $a = L$  and receive  $r = 4$ . What is your updated value of  $Q(1, L)$ ? Assume that your initial Q-values are all zero.

All Q values are initially 0.  $Q(1, L) = 0 + 0.5(4 + .5 \times 0 - 0) = 2$

## 6. Q learning

Continuing the previous problem, in the next step, you are on square 0, you choose  $a = R$ , receive  $r = 3$  and end up in  $s = 1$ . What is your updated value of  $Q(0, R)$ ?

$$Q(2, W) = 0 + 0.5(3 + 0.5 \times 2 - 0) = 2$$

## 7. Model Based Reinforcement Learning

How do we compute the Maximum-Likelihood estimate of the expected reward from taking action  $a$  in state  $s$ ?

We take the average of the rewards received from taking action  $a$  when we are in state  $s$ . We can keep track of  $N(s, a)$ , the number of times that we visit state  $s$  and perform  $a$  and  $R^{\text{total}}(s, a)$ , the sum of the rewards received from taking action  $a$  in state  $s$ . Then the ML estimate is given by

$$\frac{R^{\text{total}}(s, a)}{N(s, a)}$$

## 8. $\epsilon$ -Greedy

Why do we generally use an  $\epsilon$ -Greedy algorithm when choosing the current action during the Q-Learning algorithm? Describe how you would change the value of  $\epsilon$  over time to encourage early exploration/learning and good decision making after the state space was explored.

First we recall that in the  $\epsilon$ -greedy algorithm we choose the best action, via

$$\max_{a' \in A} Q(s, a'),$$

with probability  $1 - \epsilon$  and a random action otherwise. We use this because it encourages exploration of the state space by introducing randomness into the choice of action at each step, which we recall is necessary for RL to work with. In the beginning, we would choose a high value of  $\epsilon$ , which would encourage a lot of random exploration of the state space, and as time went on, we could slowly decrease the value of  $\epsilon$  to encourage taking the action that maximizes the Q value in the current state, which we expect to be a good action given that we've explored the state space sufficiently to closely approximate the true Q function.



## 9. Convergence of Q-Learning

Suppose we have a deterministic world where each state-action pair  $(s, a)$  is visited infinitely often. Consider an interval during which every such pair is visited. Suppose that the largest error in the approximation of  $\hat{Q}_n$  after  $n$  iterations is  $e_n$ , i.e.

$$e_n = \max_{s,a} |\hat{Q}_n(s, a) - Q(s, a)|$$

Show that  $e_{n+1}$  is bounded above by  $\gamma e_n$ , where  $\gamma$  is the usual discount factor.

$$\begin{aligned} e_{n+1} &= \max_{s,a} |\hat{Q}_{n+1}(s, a) - Q(s, a)| \\ &= |(r + \gamma \max_{a'} \hat{Q}_n(s', a')) - (r + \gamma \max_{a'} Q(s', a'))| \\ &= \gamma |\max_{a'} \hat{Q}_n(s', a') - \max_{a'} Q(s', a')| \\ &\leq \gamma \max_{a'} |\hat{Q}_n(s', a') - Q(s', a')| \\ &\leq \gamma \max_{s'', a'} |\hat{Q}_n(s'', a) - Q(s'', a)| \\ &= \gamma e_n \end{aligned}$$