EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the <u>Privacy Policy</u>.



Course > Week 11 > Final E... > Q9: Rei...

Q9: Reinforcement Learning Q9: Reinforcement Learning

Imagine an unknown game which has only two states $\{A,B\}$ and in each state the agent has two actions to choose from: $\{Up,Down\}$. Suppose a game agent chooses actions according to some policy π and generates the following sequence of actions and rewards in the unknown game:

t	s_t	a_t	s_{t+1}	r_t
0	A	Down	В	2
1	В	Down	В	-4
2	В	Up	В	0
3	В	Up	A	3
4	A	Up	A	-1

Unless specified otherwise, assume a discount factor $\gamma=0.5$ and a learning rate lpha=0.5

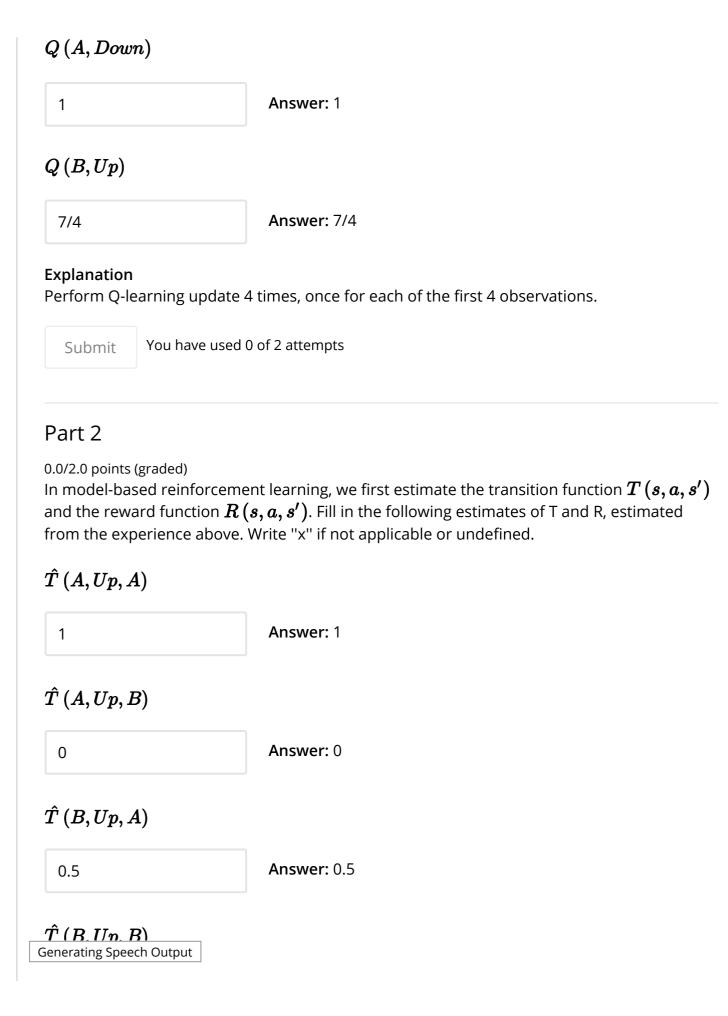
Part 1

0.0/2.0 points (graded)

Recall the update function of Q-learning is:

$$Q\left(s_{t}, a_{t}\right) \leftarrow \left(1 - \alpha\right) Q\left(s_{t}, a_{t}\right) + \alpha\left(r_{t} + \gamma \max_{a'} Q\left(s_{t+1}, a'\right)\right)$$

Assume that all Q-values initialized as 0. What are the following Q-values learned by running Generating Speech Output pove experience sequence?



0.5	Answer: 0.5
$\hat{R}\left(A,Up,A ight)$	
-1	Answer: -1
$\hat{R}\left(A,Up,B ight)$	
Х	Answer: x
$\hat{R}\left(B,Up,A ight)$	
3	Answer: 3
â (D . T. D)	

 $\hat{R}\left(B,Up,B
ight)$

0 Answer: 0

Explanation

Count transitions above and calculate frequencies. Rewards are observed rewards.

Submit

You have used 0 of 2 attempts

To decouple this question from the previous one, assume we had **a different experience** and ended up with the following estimates of the transition and reward functions:

s	a	s'	$\hat{T}(s, a, s')$	$\hat{R}(s, a, s')$
A	Up	A	1	10
A	Down	Α	0.5	2
A	Down	В	0.5	2
В	Up	A	1	-5
В	Down	В	1	8

Generating Speech Output d on the above experience.



0.0/1.0 point (graded)

Give the optimal policy $\hat{\pi}^*$ (s) for the MDP with transition function \hat{T} and reward function \hat{R}

 $\hat{\pi}^*$ (A)

Up **▼ Answer:** Up

 $\hat{\pi}^*$ (B)

Down ▼ Answer: Down

Submit

You have used 0 of 1 attempt

1 Answers are displayed within the problem

Part 4

0.0/1.0 point (graded)

Give \hat{V}^* (s) for the MDP with transition function \hat{T} and reward function \hat{R} . Hint: for any $x\in\mathbb{R}$, |x|<1, we have $1+x+x^2+x^3+x^4+\cdots=1/(1-x)$.

 ${\hat V}^*\left(A
ight)$

20 Answer: 20

 \hat{V}^* (B)

16 **Answer:** 16

Explanation

Based on the optimal policy in Part 3, calculate the value function using a Bellman equation.

Generating Speech Output

Part 5

0.0/2.0 points (graded)

If we repeatedly feed this new experience sequence through our Q-learning algorithm, what values will it converge to? Assume the learning rate α_t is properly chosen so that convergence is guaranteed.

- $^{ ext{ }}$ the values found above, \hat{V}^* 🗸
- lacksquare the optimal values, V^*
- $^{\circ}$ neither $\hat{oldsymbol{V}}^{oldsymbol{*}}$ nor $oldsymbol{V}^{oldsymbol{*}}$
- not enough information to determine

Explanation

The Q-learning algorithm will not converge to the optimal values V^* for the MDP because the experience sequence and transition frequencies replayed are not necessarily representative of the underlying MDP. (For example, the true T(A,Down,A) might be equal to 0.75, in which case, repeatedly feeding in the above experience would not provide an accurate sampling of the MDP.) However, for the MDP with transition function \hat{T} and reward function \hat{R} , replaying this experience repeatedly will result in Q-learning converging to its optimal values \hat{V}^* .

Submit

You have used 0 of 1 attempt

1 Answers are displayed within the problem