

1 Reinforcement Learning

Imagine an unknown environments with four states (A, B, C, and X), two actions (\leftarrow and \rightarrow). An agent acting in this environment has recorded the following episode:

s	a	s'	r	Q-learning iteration numbers (for part b)
A	\rightarrow	B	0	1, 10, 19, ...
B	\rightarrow	C	0	2, 11, 20, ...
C	\leftarrow	B	0	3, 12, 21, ...
B	\leftarrow	A	0	4, 13, 22, ...
A	\rightarrow	B	0	5, 14, 23, ...
B	\rightarrow	A	0	6, 15, 24, ...
A	\rightarrow	B	0	7, 16, 25, ...
B	\rightarrow	C	0	8, 17, 26, ...
C	\rightarrow	X	1	9, 18, 27, ...

- (a) Consider running model-based reinforcement learning based on the episode above. Calculate the following quantities:

$$\hat{T}(B, \rightarrow, C) = \frac{2}{3}$$

$$\hat{R}(C, \rightarrow, X) = 1$$

- (b) Now consider running Q-learning, repeating the above series of transitions in an infinite sequence. Each transition is seen at multiple iterations of Q-learning, with iteration numbers shown in the table above. After which iteration of Q-learning do the following quantities first become nonzero? (If they always remain zero, write *never*).

$$Q(A, \rightarrow)? \quad 14$$

$$Q(B, \leftarrow)? \quad 22$$

- (c) True/False: For each question, you will get positive points for correct answers, zero for blanks, and negative points for incorrect answers. Circle your answer **clearly**, or it will be considered incorrect.

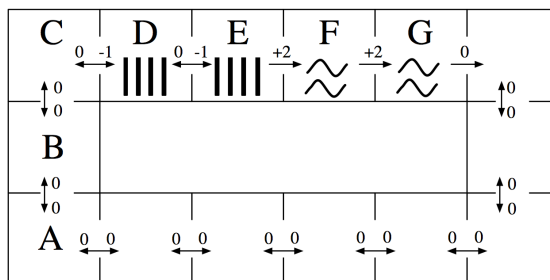
- (i) [true or false] In Q-learning, you do not learn the model.
 Q learning is model-free, you learn the optimal policy explicitly, and the model itself implicitly.
- (ii) [true or false] For TD Learning, if I multiply all the rewards in my update by some nonzero scalar p , the algorithm is still guaranteed to find the optimal policy.
 TD Learning does not necessarily find the optimal policy, it only learns the value of the states following some given policy.
- (iii) [true or false] In Direct Evaluation, you recalculate state values after each transition you experience.

In order to estimate state values, you calculate state values from episodes of training, not single transitions.

- (iv) [*true* or *false*] Q-learning requires that all samples must be from the optimal policy to find optimal q-values.
 Q-learning is off-policy, you can still learn the optimal values even if you act suboptimally sometimes.

2 MDPs: Grid-World Water Park

Consider the MDP drawn below. The state space consists of all squares in a grid-world water park. There is a single waterslide that is composed of two ladder squares and two slide squares (marked with vertical bars and squiggly lines respectively). An agent in this water park can move from any square to any neighboring square, unless the current square is a slide in which case it must move forward one square along the slide. The actions are denoted by arrows between squares on the map and all deterministically move the agent in the given direction. The agent cannot stand still: it must move on each time step. Rewards are also shown below: the agent feels great pleasure as it slides down the water slide (+2), a certain amount of discomfort as it climbs the rungs of the ladder (-1), and receives rewards of 0 otherwise. The time horizon is infinite; this MDP goes on forever.



- (a) How many (deterministic) policies π are possible for this MDP?
 2^{11}
- (b) Fill in the blank cells of this table with values that are correct for the corresponding function, discount, and state. *Hint: You should not need to do substantial calculation here.*

	γ	$s = A$	$s = E$
$V_3(s)$	1.0	0	4
$V_{10}(s)$	1.0	2	4
$V_{10}(s)$	0.1	0	2.2
$Q_1(s, \text{west})$	1.0	—	0
$Q_{10}(s, \text{west})$	1.0	—	3
$V^*(s)$	1.0	∞	∞
$V^*(s)$	0.1	0	2.2

$V_{10}^*(A), \gamma = 1$: In 10 time steps with no discounting, the rewards don't decay, so the optimal strategy is to climb the two stairs (-1 reward each), and then slide down the two slide squares (+2 rewards each). You only have time to do this once. Summing this up, we get $-1 - 1 + 2 + 2 = 2$.

$V_{10}^*(E), \gamma = 1$: No discounting, so optimal strategy is sliding down the slide. That's all you have time for. Sum of rewards = $2 + 2 = 4$.

$V_{10}^*(A), \gamma = 0.1$. The discount rate is 0.1, meaning that rewards 1 step further into the future are discounted by a factor of 0.1. Let's assume from A, we went for the slide. Then, we would have to take the actions

$A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow E, E \rightarrow F, F \rightarrow G$. We get the first -1 reward from $C \rightarrow D$, discounted by γ^2 since it is two actions in the future. $D \rightarrow E$ is discounted by γ^3 , $E \rightarrow F$ by γ^4 , and $F \rightarrow G$ by γ^5 . Since γ is low, the positive rewards you get from the slide have less of an effect as the larger negative rewards you get from climbing up. Hence, the sum of rewards of taking the slide path would be negative; the optimal value is 0.

$V_{10}^*(E), \gamma = 0.1$. Now, you don't have to do the work of climbing up the stairs, and you just take the slide down. Sum of rewards would be 2 (for $E \rightarrow F$) + 0.2 (for $F \rightarrow G$, discounted by 0.1) = 2.2.

$Q_{10}^*(E, west), \gamma = 1$. Remember that a Q-state (s,a) is when you start from state s and are committed to taking a. Hence, from E, you take the action West and land in D, using up one time step and getting an immediate reward of 0. From D, the optimal strategy is to climb back up the higher flight of stairs and then slide down the slide. Hence, the rewards would be $-1(D \rightarrow E) + 2(E \rightarrow F) + 2(F \rightarrow G) = 3$.

$V^*(s), \gamma = 1$. Infinite game with no discount? Have fun sliding down the slide to your content from anywhere.

$V^*(s), \gamma = 0.1$. Same reasoning apply to both A and E from $V_{10}^*(s)$. With discounting, the stairs are more costly to climb than the reward you get from sliding down the water slide. Hence, at A, you wouldn't want to head to the slide. From E, since you are already at the top of the slide, you should just slide down.

- (c) Fill in the blank cells of this table with the Q-values that result from applying the Q-update for the transition specified on each row. You may leave Q-values that are unaffected by the current update blank. Use discount $\gamma = 1.0$ and learning rate $\alpha = 0.5$. Assume all Q-values are initialized to 0. (Note: the specified transitions would not arise from a single episode.)

	$Q(D, west)$	$Q(D, east)$	$Q(E, west)$	$Q(E, east)$
Initial:	0	0	0	0
Transition 1: $(s = D, a = east, r = -1, s' = E)$		-0.5		
Transition 2: $(s = E, a = east, r = +2, s' = F)$				1.0
Transition 3: $(s = E, a = west, r = 0, s' = D)$				
Transition 4: $(s = D, a = east, r = -1, s' = E)$		-0.25		