

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:

$$\hat{T}(B, \rightarrow, C) = \underline{\hspace{2cm}}$$

$$\hat{R}(C, \rightarrow, X) = \underline{\hspace{2cm}}$$

- (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*).

$$Q(A, \rightarrow)? \underline{\hspace{2cm}}$$

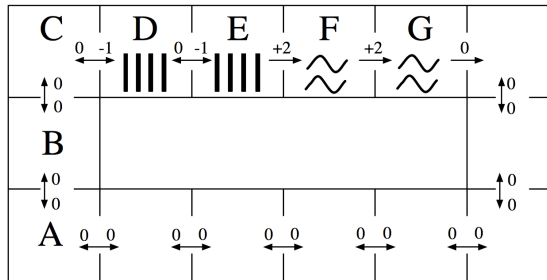
$$Q(B, \leftarrow)? \underline{\hspace{2cm}}$$

- (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.
- (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.
- (iii) [*true* or *false*] In Direct Evaluation, you recalculate state values after each transition you experience.
- (iv) [*true* or *false*] Q-learning requires that all samples must be from the optimal policy to find optimal q-values.

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?
- (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		
$V_{10}(s)$	1.0		
$V_{10}(s)$	0.1		
$Q_1(s, \text{west})$	1.0	---	
$Q_{10}(s, \text{west})$	1.0	---	
$V^*(s)$	1.0		
$V^*(s)$	0.1		

- (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.)

	$Q(D, \text{west})$	$Q(D, \text{east})$	$Q(E, \text{west})$	$Q(E, \text{east})$
Initial:	0	0	0	0
Transition 1: ($s = D, a = \text{east}, r = -1, s' = E$)				
Transition 2: ($s = E, a = \text{east}, r = +2, s' = F$)				
Transition 3: ($s = E, a = \text{west}, r = 0, s' = D$)				
Transition 4: ($s = D, a = \text{east}, r = -1, s' = E$)				