# HW3

• This is a preview of the published version of the quiz

Started: Aug 15 at 7:05pm

# **Quiz Instructions**

Question 1 0 pts

**Honor Pledge:** Please type your name as a signature in the form.

I avow that I will not give or receive any unauthorized help on this exam, and that all work will be my own.

Edit View Insert Format Tools Table (i) 0 words </> / iii р

## Question 2 2 pts

Answer the following question?

- 1. Explain how the  $\varepsilon$ -greedy approach balances exploration and exploitation.
- 2. Explain how the incremental mean method in the utility mean update reduces the memory usage.
- 3. Explain the difference between the Monte-Carlo (MC) method and the Temporal-Difference (TD) method.

<ol> <li>Explain how the temporal-difference (TD) method in the utility update reduces the memory usage.</li> <li>Explain in what occasion the TD Q-value update might result different results under SARSA learning and Q-Learning, respectively.</li> </ol>								
Edit View Insert Format Tools Table								
12pt $\vee$ Paragraph $\vee$ $\bigcirc$								
p								

Question 3 2 pts

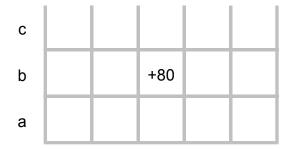
We extend the windy maze defined in HW1 with probabilistic outcome after an action. The maze map is shown here for your convenience.

d	←	<b>←</b>	<b>←</b>	<b>←</b>	←
С	<b>1</b>			↓ ↓	<b>1</b>
b	$\rightarrow$	$\rightarrow$	+80		1
а	$\rightarrow$	$\rightarrow$	1	<b>←</b>	←
	1	2	3	4	5

The agent can drift to the left or the right with probability 0.2 or go straight with probability 0.6. If the drifting direction is a wall, it will be bounced back to the original position. We consider the policy defined in the map. Calculate which open squares can be reached from 4d after three sequential actions under this policy and with what probabilities.

Edit View Insert Format Tools Table **≣** ∨ **!≡** ∨ **!≡** ∨ √x ↔ **★ (†)** 0 words </> **(\*) !!** р **Question 4** 

# We consider Value Iteration for the above MDP problem with the same drift probability assumption but without the given policy as shown here.



We still assume the wind comes from the north and the cost of one step is defined as follows (the reward will be the negation of the cost): 1 for moving southward; 2 for moving eastward or westward; 3 for moving northward. The reward function R(s,a) will be the negative of the cost. We also assume that the reward reaching at the goal is 80. Since the reward function R(s;a) here depends on both the state and the action taken at this state, all utility equations should be revised as

$$U(s) \leftarrow \max_{a} (R(s,a) + \gamma \Sigma_{s'} P(s'|s,a) U(s'))$$

We choose  $\gamma=1$ . We assume the initial utility at any state is 0 except for 80 at the goal state. We perform one-iteration update of the utility of all states in this order: a1, a2, ..., a5, b1, b2, ..., b5, c1, ..., c5, d1, d2, ..., d5. Please give the updated utility at each state.

Edit View Insert Format Tools Table (i) 0 words </> // !! р

Question 5

We consider SARSA with ε-Greedy algorithm for the above MDP problem. However, the model including transition probabilities and reward functions is unknown to the agent. The following tables show the latest values of N(s,a) (top) and Q(s,a) (bottom):

			1					1
	10	5	5	12	0	24		
	11		3		1			
	5					11		
b	18	9			55	2		
	13				4			
		6	7	•		23	4	
а	8	19	3	15	1	8	7	8
		6	4		4		74	
	1		2		3		4	
	-2.4		-1.8		41.0			
С	-2.4	16.4	-1.8	60.8	0.0	87.8	+10	00
	-2.4		-1.8		3.4			
	-2.9				44.8			
b	-2.9	-2.9			12.2	-101.0	-10	0
	-2.9				-51.5			
	-2.6		-1.8		12.0		-101.0	
а	-2.5	-2.3	-2.0	0.4	-2.0	-2.1	-16.0	-14.9
	-2.4		-2.1		-1.9		-3.5	
	1		2		3		4	

Then we run an additional trial in sequence (drifting might happen or random action might be chosen during the trial): a1,E,-2; a2,N,-1; a2,E,-2; a3,N,-1; a4,S,-3; a4,W,-2; a3,N,-1; b3,N,-1; c3,E,-2; c4. The agent performs TD updates sequentially regarding this trial with the following equations:

$$N(s,a) \leftarrow N(s,a) + 1$$
  
 $Q(s,a) \leftarrow Q(s,a) + 1/N(s,a) (R(s,a) + yQ(s',a') - Q(s,a))$ 

We choose  $\gamma=1$ . Please show the updated Q-values Q(s,a) during this trial (Keep two decimal values for each Q-value):

- 1. Q(a1,E):
- 2. Q(a2,N):
- 3. Q(a2,E):
- 4. Q(a3,N):
- 5. Q(a4,S):
- 6. Q(a4,W):
- 7. Q(a3,N):
- 8. Q(b3,N):
- 9. Q(c3,E):

р

### Question 6

3 pts

We replace SARSA learning with Q-Learning in the above question and redo the question, where you should use the following TD Q-value updates:

$$N(s,a) \leftarrow N(s,a) + 1$$



Not saved

Submit Quiz