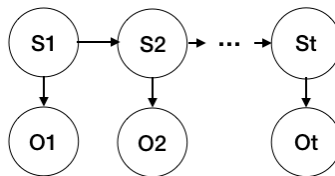


This assignment contains 3 questions. Please read and follow the following instructions.

- **DUE DATE: Oct 30th, 11:45 PM**
  - **TOTAL NUMBER OF POINTS: 100**
  - Clearly list your team ID, each team member's **names and Unity IDs** at the top of your submission.
  - Submit only a **single pdf file** of your answers.
- 

1. (30 points) [Farzaneh Khoshnevisan] [HMM]

Infection is a common condition among patients in ICU settings and can have various roots, which makes it challenging to be determined. Assume that we want to model infection using an HMM, while *infection* is the hidden state and the only available observation is *blood pressure* (0 for normal and 1 for abnormal). When patients entering ICU, the probability of being infected is 0.75. At any given time, infected patients have 40% chance of improving to be uninfected and uninfected patients have 20% chance of becoming infected. There is 80% chance of observing an abnormal blood pressure for infected patients while only 10% chance of observing abnormal observation for uninfected patients.

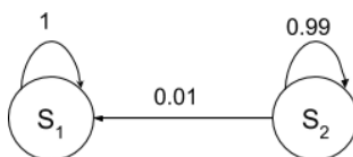


- (5 points) How many parameters are required to fully define this HMM. Justify your answer.
- (3 points) Create initial, transition, and emission probability tables based on the problem statement given above.
- (6 points) Using the described HMM and the generated probability tables, apply the forward algorithm to compute the probability that we observe the sequence  $\{0, 1, 1\}$  blood pressure. Show your work (i.e., show each of your  $\alpha$ s).
- (6 points) Using the backward algorithm, compute the probability that we observe the aforementioned sequence ( $\{0, 1, 1\}$ ). Again, show your work (i.e., show each of your  $\beta$ s).

- (e) (4 points) Using the forward-backward algorithm, compute the most likely setting for each state. Show your work.
- (f) (6 points) Use the Viterbi algorithm to compute the most likely sequence of states. Show your work.

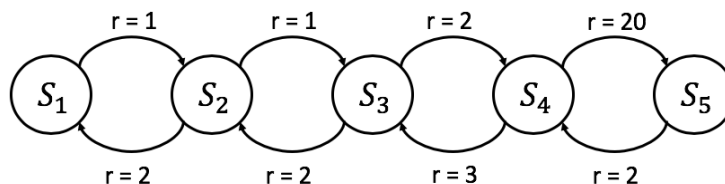
2. (20 points) [Farzaneh Khoshnevisan] [**HMM**]

We design a two-state HMM for a dice toss game between two players where the output of their toss is our only observation.  $S_1$  and  $S_2$  indicate the hidden state where the dice is tossed by player 1 and player 2, respectively. The transition probabilities between these two states are given in diagram below. Both players are equally likely to start the game and they each play with a biased dice. The output distribution corresponding to each player is defined over  $\{1, 2, 3, 4, 5, 6\}$  and is given in the table below the diagram.



	$S_1$	$S_2$
$P(x=1)$	0.01	0.32
$P(x=2)$	0.21	0.14
$P(x=3)$	0.3	0.14
$P(x=4)$	0	0.08
$P(x=5)$	0.18	0
$P(x=6)$	0.3	0.32

- (a) (5 points) Give an example of one output sequence of length 2, which cannot be generated by the given HMM. Justify your answer.
- (b) (5 points) We generated a sequence of 20,601<sup>2020</sup> observations from this HMM, and found that the last observation in the sequence was 1. What is the most likely hidden state corresponding to the last observation? Justify your answer.
- (c) (5 points) Consider an output sequence  $\{1, 1\}$ . What is the most likely sequence of hidden states corresponding to this output observation sequence? Show your work.
- (d) (5 points) Now, consider an output sequence  $\{1, 1, 3\}$ . What are the first two states of the most likely hidden state sequence? Show your work.
3. (50 points) [Song Ju] [**Reinforcement Learning**]
- Consider the following Markov Decision Process:



Our state space  $S$ :  $\{S_1, S_2, S_3, S_4, S_5\}$  and our action space  $A$ :  $\{\text{"Left"}, \text{"Right"}\}$ . For all parts of this problem, assume that  $\gamma = 0.8$ .

For subquestions (a)-(c) below, we assume that all actions are deterministic:

- (5 points) What is the optimal policy for this MDP?
- (15 points) Calculate  $V_{S_5}^*$ . Show your work. (The acceptable answer should be a numeric number and you should show all the key steps).
- (10 points) Consider executing Q-learning on this MDP. Assume that 1) all of the initial Q values are 0, 2)  $\alpha = 0.5$ , and 3) it uses a greedy exploration policy by always choosing the action with maximum Q value at any given state. The algorithm breaks ties by choosing "Left". What are the first 10 (state, action) pairs if our robot learns using Q-learning and starts in state  $S_3$ ? (A candidate answer can be expressed in the form of:  $(S_3, \text{Left}), (S_2, \text{Right}), (S_3, \text{Right}), \dots$ )

For the subquestions [d]-[e], assume the actions are not deterministic in that after taking an action, it's possible to stay in the same state:

- (10 points) Consider executing Value Iteration on this MDP. The transition matrix shown below indicates the probability of transitioning from state  $s$  to state  $s'$  by taking action  $a$ . In the matrix, the first column is the start state and the first row is the ending state. For example, when take action *left* in state  $S_2$ , the probability of transit to  $S_1$  is 0.6 while the probability of stay in  $S_2$  is 0.4. All the empty cell in the matrix means there's no transition.

Action: Left					
$S_t \backslash S_{t+1}$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$
$S_1$					
$S_2$	0.6	0.4			
$S_3$		0.7	0.3		
$S_4$			0.6	0.4	
$S_5$				1	

Action: Right					
$S_t \backslash S_{t+1}$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$
$S_1$		1			
$S_2$		0.3	0.7		
$S_3$			0.8	0.2	
$S_4$				0.5	0.5
$S_5$					

For a given iteration  $t$ , the value functions of each state are:  $V_{S_1}^t = 20$ ,  $V_{S_2}^t = 30$ ,  $V_{S_3}^t = 20$ ,  $V_{S_4}^t = 30$ ,  $V_{S_5}^t = 10$ , compute new value function of all states in the next iteration,  $t + 1$ , using Value Iteration.

- (e) (10 points) Based on the Value function and transition matrix in previous sub-question (d), answer the following questions: 1) what's the optimal policy at time  $t$ ? 2) what's the optimal policy at  $t + 1$ , after running Value Iteration? 3) Are the two policies different?