



# PRESIDENCY UNIVERSITY

BENGALURU

Roll No.																			
----------	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

## Make up Examinations - December 2025

Date: 31 - 12 - 2025

Time: 09:30am - 12:30pm

School: SOCSE	Program: B.Tech	
Course Code :CSE3011	Course Name :REINFORCEMENT LEARNING	
Semester: MK	Max Marks:100	Weightage:50%

CO - Levels	CO1	CO2	CO3	CO4	CO5
Marks	16	26	34	24	

### Instructions:

- (i) Read all questions carefully and answer accordingly.
- (ii) Do not write anything on the question paper other than roll number.

### Part A

Answer ALL the Questions. Each question carries 2marks.

10Q x2M=20M

1	Define the return of a trajectory for a continuous task	2 Marks	L1	CO1
2	Illustrate the equation to find $Q(S)$ in the Monte Carlo method	2 Marks	L2	CO2
3	Interpret SARSA update rule	2 Marks	L2	CO3
4	Define the following a) Agent b) Environment	2 Marks	L1	CO1
5	Illustrate the equation to find $V(S)$ in the Monte Carlo method	2 Marks	L2	CO2
6	State optimal policy	2 Marks	L1	CO2
7	$X$ is a random variable with the outcomes of throwing a dice. Find the expectation $E[f(x)]$ where $f(x)=X^3$	2 Marks	L1	CO1
8	Distinguish between SARSA and Q-Learning algorithms	2 Marks	L2	CO3
9	Interpret upper confidence bound of an arm in MAB Problem	2 Marks	L1	CO4
10	Describe Thompson sampling	2 Marks	L2	CO4

## Part B

**Answer all the Questions**

**Total 80 Marks.**

<b>11.</b>	<b>a.</b>	Discuss stochastic environment and deterministic environment in RL with an example	<b>10 Marks</b>	<b>L2</b>	<b>CO1</b>
	<b>b.</b>	Explain Markov chain and Markov property in Markov decision process	<b>10 Marks</b>	<b>L3</b>	<b>CO2</b>

**Or**

<b>12.</b>	<b>a.</b>	Illustrate Bellman Equation for the Value function in deterministic environment	<b>10 Marks</b>	<b>L2</b>	<b>CO1</b>																																			
	<b>b.</b>	<p>Using the value iteration algorithm and the model dynamics of state A given in the table below, Compute the optimal value of state A, after the first iteration</p> <table border="1" style="margin: 10px auto; border-collapse: collapse; text-align: center;"> <thead> <tr> <th>State (<math>s</math>)</th> <th>Action (<math>a</math>)</th> <th>Next State (<math>s'</math>)</th> <th>Transition Probability <math>P(s' s,a)</math> or <math>P_{ss'}^a</math></th> <th>Reward Function <math>R(s,a,s')</math> or <math>R_{ss'}^a</math></th> </tr> </thead> <tbody> <tr><td>A</td><td>0</td><td>A</td><td>0.1</td><td>0</td></tr> <tr><td>A</td><td>0</td><td>B</td><td>0.8</td><td>-1</td></tr> <tr><td>A</td><td>0</td><td>C</td><td>0.1</td><td>1</td></tr> <tr><td>A</td><td>1</td><td>A</td><td>0.1</td><td>0</td></tr> <tr><td>A</td><td>1</td><td>B</td><td>0.0</td><td>-1</td></tr> <tr><td>A</td><td>1</td><td>C</td><td>0.9</td><td>0</td></tr> </tbody> </table>	State ( $s$ )	Action ( $a$ )	Next State ( $s'$ )	Transition Probability $P(s' s,a)$ or $P_{ss'}^a$	Reward Function $R(s,a,s')$ or $R_{ss'}^a$	A	0	A	0.1	0	A	0	B	0.8	-1	A	0	C	0.1	1	A	1	A	0.1	0	A	1	B	0.0	-1	A	1	C	0.9	0	<b>10 Marks</b>	<b>L3</b>	<b>CO2</b>
State ( $s$ )	Action ( $a$ )	Next State ( $s'$ )	Transition Probability $P(s' s,a)$ or $P_{ss'}^a$	Reward Function $R(s,a,s')$ or $R_{ss'}^a$																																				
A	0	A	0.1	0																																				
A	0	B	0.8	-1																																				
A	0	C	0.1	1																																				
A	1	A	0.1	0																																				
A	1	B	0.0	-1																																				
A	1	C	0.9	0																																				

<b>13.</b>	<b>a.</b>	Explain different types of RL environments with an example each	<b>10 Marks</b>	<b>L2</b>	<b>CO2</b>
	<b>b.</b>	Articulate TD Prediction algorithm in FZLE environment	<b>10 Marks</b>	<b>L3</b>	<b>CO3</b>

**Or**

<b>14.</b>	<b>a.</b>	Discuss Markov Decision Process in Detail	<b>10 Marks</b>	<b>L2</b>	<b>CO2</b>
	<b>b.</b>	Articulate TD Control algorithm in FZLE environment	<b>10 Marks</b>	<b>L3</b>	<b>CO3</b>

<b>15.</b>	<b>a.</b>	Determine the Algorithmic steps of the off policy TD Control-Q Learning technique .Also find the updated Q Value of state (3,2) using the following data , apply epsilon greedy technique (choose exploitation) to choose current and next actions. Assume alpha as 0.1 and gamma as 1	<b>10 Marks</b>	<b>L3</b>	<b>CO3</b>
------------	-----------	--	-----------------	-----------	------------

	1	2	3	4
1	S	F	F	F
2	F	H	F	H
3	F	F	F	H
4	H	F	F	G

State	Action	Value
(1,1)	Up	0.5
⋮	⋮	⋮
(3,2)	Up	0.1
(3,2)	Down	0.8
(3,2)	Left	0.5
(3,2)	Right	0.6
⋮	⋮	⋮
(4,4)	Right	0.5

Figure 5.13: The Frozen Lake environment with a randomly initialized Q table

State	Action	Value
(1,1)	Up	0.5
⋮	⋮	⋮
(4,2)	Up	0.3
(4,2)	Down	0.5
(4,2)	Left	0.1
(4,2)	Right	0.8
⋮	⋮	⋮
(4,4)	Right	0.5

**b** Illustrate upper confidence bound method to overcome the exploration –exploitation dilemma with the algorithm

**10 Marks**

**L3 CO4**

**Or**

**16.** **a.** Articulate MAB problem and list out the applications of MAB

**10 Marks**

**L3 CO3**

**b** Apply softmax exploration method of exploration-exploitation strategy for the following 4 arm bandit

arm	Q
arm 1	1
arm 2	0
arm 3	0
arm 4	0

**10 Marks**

**L3 CO4**

**17.** **a.** Explain SARSA in FZLE

**10 Marks**

**L2 CO3**

**b.** Implement the reinforcement learning environment namely, Frozen Lake Environment using a random policy and show the output of the following:

**10 Marks**

**L3 CO4**

	a. Create and render the environment c. Action Space	b. State Space d. State transition probabilities			
--	---	---	--	--	--

**Or**

<b>18.</b>	<b>a.</b>	Explain TD control types in brief	<b>10 Marks</b>	<b>L2</b>	<b>CO3</b>
	<b>b.</b>	<p>Write a python program to find an optimal policy using Q-learning for the frozen lake environment with alpha=0.85 and gamma=0.90 and epsilon =0.8</p> <p>Create and render the environment</p> <p>Generate policy using 30 episodes with 60 timesteps</p> <p>Print the optimal policy</p>	<b>10 Marks</b>	<b>L3</b>	<b>CO4</b>

**\*\*\*\*\* BEST WISHES \*\*\*\*\***