

U.S.N.									
--------	--	--	--	--	--	--	--	--	--

# B.M.S. College of Engineering, Bengaluru-560019

Autonomous Institute Affiliated to VTU

## June 2025 Semester End Main Examinations

**Programme: B.E.**

**Semester: VI**

**Branch: Artificial Intelligence and Machine Learning**

**Duration: 3 hrs.**

**Course Code: 24AM6PCREL**

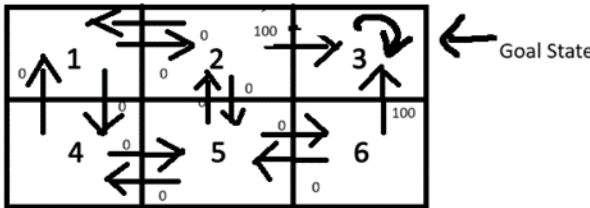
**Max Marks: 100**

**Course: Reinforcement Learning**

**Instructions:** 1. Answer any FIVE full questions, choosing one full question from each unit.  
2. Missing data, if any, may be suitably assumed.

<b>Important Note:</b> Completing your answers, compulsorily draw diagonal cross lines on the remaining blank pages. Revealing of identification, appeal to evaluator will be treated as malpractice.			<b>UNIT - I</b>	<b>CO</b>	<b>PO</b>	<b>Marks</b>
	1	a)	Explain the concept of evaluative feedback in REL (Reinforcement Learning) and highlight how it contrasts with instructive feedback. Also, mention an example of an algorithm that relies on instructive feedback.	CO1	PO1	06
		b)	Analyze how learning occurs in the Tic-Tac-Toe example using value functions.	CO2	PO1	06
		c)	i. Illustrate the working of the $\epsilon$ -greedy algorithm and its approach to balancing exploration and exploitation in reinforcement learning with an example. ii. Interpret the significance of decay in the $\epsilon$ -greedy algorithm and its impact on the learning process.	CO2	PO1	08
			<b>OR</b>			
	2	a)	Outline the main components of a reinforcement learning system and describe the role of each component with a diagram.	CO1	PO1	06
		b)	Derive the incremental implementation equation from the Action-Value Method.	CO2	PO1	06
		c)	i. Formulate the Softmax equation used in reinforcement learning and identify its key components. ii. Consider initial preferences $H(A) = 2$ , $H(B) = 1$ , and $H(C) = 0$ , compute the action probabilities using the Softmax method. iii. Based on the computed probabilities, identify the most likely action to be chosen.	CO2	PO1	08
			<b>UNIT - II</b>			
	3	a)	i. Explain the Markov Property in the context of reinforcement learning.	CO1	PO2	06

		ii. How does the state transition matrix represent the dynamics of a Markov process?																												
	b)	Identify the differences between the State Value Function and the Action Value Function, and present their mathematical definitions.	CO2	PO1	06																									
	c)	In reinforcement learning, agents may interact with their environment through tasks that are either episodic or continuing in nature. i. Describe how episodic and continuing tasks differ in terms of agent-environment interaction. ii. Explain the importance of discount factor $\gamma \in [0,1]$ in continuing tasks. iii. Provide mathematical expression for the term G. iv. An agent receives a sequence of rewards over 5-time steps as follows: $R = [3, 2, 4, 1, 5]$ . Compute the total return G using discount factor: i. $\gamma = 0.9$ ii. $\gamma = 0.1$ v. What value of $\gamma$ reflects a far-sighted agent and which reflects a near-sighted agent and compare both?	CO3	PO1	08																									
		OR																												
4	a)	Differentiate between Dynamic programming and Monte Carlo method.	CO1	PO1	06																									
	b)	The value function $v\pi(s)$ under a policy $\pi$ is defined as the expected return when starting from state $s$ and following $\pi$ thereafter. Derive the recursive form of the value function $v\pi(s)$ in terms of immediate reward, discount factor $\gamma$ , and the next state's value. Clearly show how the Bellman equation is obtained.	CO2	PO1	06																									
	c)	A Markov Decision Process has the states $S= \{S1, S2, S3\}$ and actions $A = \{A1, A2\}$ . The transition and reward probabilities are given in the table. <table><tr><td>States</td><td>Action a</td><td>Next States</td><td>Reward r</td><td>Probability <math>p(s', r s, a)</math></td></tr><tr><td>S1</td><td>A1</td><td>S2</td><td>3</td><td>0.5</td></tr><tr><td>S1</td><td>A1</td><td>S3</td><td>7</td><td>0.5</td></tr><tr><td>S2</td><td>A2</td><td>S3</td><td>4</td><td>0.6</td></tr><tr><td>S2</td><td>A2S</td><td>S1</td><td>2</td><td>0.4</td></tr></table> Compute the following: i. $p(S3,7 S1,A1)$ ii. $r(S1, A1)$ . iii. $p(S3 S1,A1)$ iv. $r(S1,A1,S3)$	States	Action a	Next States	Reward r	Probability $p(s', r s, a)$	S1	A1	S2	3	0.5	S1	A1	S3	7	0.5	S2	A2	S3	4	0.6	S2	A2S	S1	2	0.4	CO3	PO1	08
States	Action a	Next States	Reward r	Probability $p(s', r s, a)$																										
S1	A1	S2	3	0.5																										
S1	A1	S3	7	0.5																										
S2	A2	S3	4	0.6																										
S2	A2S	S1	2	0.4																										

UNIT - III																																													
5	a)	<p>Consider a 2x3 Grid world with immediate rewards given below.</p> <div></div> <p>Let the state space be <math>S = \{1,2,3,4,5,6\}</math> with grid 3 being the Goal State. The agent starts in any grid (0 to 6) and must learn to reach grid 3 using Q-learning. Moving into grid 3 (from rooms 2,6) gives a reward of +100 and all other moves give a reward of 0.</p> <p>i. Represent the environment using a reward matrix R.</p> <p>ii. Implement the Q-learning algorithm to update the Q-values in the Q Table (Consider Discount factor <math>\gamma= 0.9</math>).</p>	CO3	PO1	10																																								
	b)	<p>Imagine you are designing a reinforcement learning agent to play a complex strategy game. The environment is initially unknown to the agent, but you have two options for your learning approach:</p> <ul style="list-style-type: none"><li>Option A: The agent tries to learn a model of the game's rules, including how the game state changes based on its actions, and then uses this model to plan future moves.</li><li>Option B: The agent learns a policy directly by interacting with the game and adjusting its actions based on the rewards it receives, without explicitly modelling the game's rules.</li></ul> <p>i. Identify which option corresponds to model-based learning and which corresponds to model-free learning.</p> <p>ii. Discuss the differences between the two.</p>	CO2	PO1	06																																								
	c)	<p>Describe the <math>TD(0)</math> algorithm used for policy evaluation, including its update rule and key steps.</p>	CO2	PO1	04																																								
		OR																																											
6	a)	<p>A delivery drone is flying through different checkpoints (P, Q, R, S) to deliver packages. Each checkpoint grants the drone a reward based on delivery success or time saved. The drone's flight path and rewards are recorded over 5 episodes. The discount factor <math>\gamma=1</math> (no discounting).</p> <table><tr><th>Episode</th><th>Time Step</th><th>State</th><th>Reward</th></tr><tr><td>1</td><td>0</td><td>P</td><td>0</td></tr><tr><td></td><td>1</td><td>Q</td><td>1</td></tr><tr><td></td><td>2</td><td>R</td><td>2</td></tr><tr><td></td><td>3</td><td>S</td><td>3</td></tr><tr><td></td><td>4</td><td>Terminal</td><td>-</td></tr><tr><td>2</td><td>0</td><td>Q</td><td>1</td></tr><tr><td></td><td>1</td><td>R</td><td>0</td></tr><tr><td></td><td>2</td><td>Q</td><td>2</td></tr><tr><td></td><td>3</td><td>S</td><td>2</td></tr></table>	Episode	Time Step	State	Reward	1	0	P	0		1	Q	1		2	R	2		3	S	3		4	Terminal	-	2	0	Q	1		1	R	0		2	Q	2		3	S	2	CO3	PO1	10
Episode	Time Step	State	Reward																																										
1	0	P	0																																										
	1	Q	1																																										
	2	R	2																																										
	3	S	3																																										
	4	Terminal	-																																										
2	0	Q	1																																										
	1	R	0																																										
	2	Q	2																																										
	3	S	2																																										

			<table><tr><td></td><td>4</td><td>Terminal</td><td>-</td></tr><tr><td>3</td><td>0</td><td>R</td><td>2</td></tr><tr><td></td><td>1</td><td>Q</td><td>1</td></tr><tr><td></td><td>2</td><td>S</td><td>3</td></tr><tr><td></td><td>3</td><td>Terminal</td><td>-</td></tr><tr><td>4</td><td>0</td><td>P</td><td>1</td></tr><tr><td></td><td>1</td><td>S</td><td>3</td></tr><tr><td></td><td>2</td><td>R</td><td>2</td></tr><tr><td></td><td>3</td><td>Terminal</td><td>-</td></tr><tr><td>5</td><td>0</td><td>Q</td><td>3</td></tr><tr><td></td><td>1</td><td>R</td><td>1</td></tr><tr><td></td><td>2</td><td>P</td><td>2</td></tr><tr><td></td><td>3</td><td>S</td><td>0</td></tr><tr><td></td><td>4</td><td>Terminal</td><td>-</td></tr></table>		4	Terminal	-	3	0	R	2		1	Q	1		2	S	3		3	Terminal	-	4	0	P	1		1	S	3		2	R	2		3	Terminal	-	5	0	Q	3		1	R	1		2	P	2		3	S	0		4	Terminal	-			
	4	Terminal	-																																																											
3	0	R	2																																																											
	1	Q	1																																																											
	2	S	3																																																											
	3	Terminal	-																																																											
4	0	P	1																																																											
	1	S	3																																																											
	2	R	2																																																											
	3	Terminal	-																																																											
5	0	Q	3																																																											
	1	R	1																																																											
	2	P	2																																																											
	3	S	0																																																											
	4	Terminal	-																																																											
		<p>i. Calculate the return <math>G_t</math> (sum of future rewards) for each visit to each checkpoint in all episodes (with <math>\gamma=1</math>).</p> <p>ii. Compute the estimates for each checkpoint by averaging the returns from their first visits in each episode.</p> <p>iii. Compute the every-visit Monte Carlo estimates for each checkpoint by averaging the returns from all visits across episodes.</p>																																																												
	b)	Differentiate between SARSA and Q-Learning Algorithms	CO1	PO1	06																																																									
	c)	List the advantages of Monte Carlo over Dynamic Programming.	CO1	PO1	04																																																									
		UNIT - IV																																																												
7	a)	<p>i. Analyze how the use of the max operator in Q-learning can lead to instability during training.</p> <p>ii. Suggest any one technique that addresses the instability caused by the max operator in Deep Q-Learning.</p> <p>iii. Elaborate on the use of Experience Reply buffer in Deep Q – Learning.</p>	CO1	PO1	10																																																									
	b)	<p>An agent interacts with the environment and observes the following transition: State A <math>\rightarrow</math> Action a <math>\rightarrow</math> Reward = 2 <math>\rightarrow</math> State B using a Deep Q-Network (DQN), where the Q- function is approximated using a second-degree polynomial over state features as follows: <math>Q(s, a; w) = w_1\phi_1(S) + w_2[\phi_1(S)]^2 + w_3\phi_2(S) + w_4[(\phi_2(S))^2]</math></p> <p><math>\rightarrow</math> Feature vectors:</p> <ul style="list-style-type: none"><li><math>\phi(A) = [1,2]</math></li><li><math>\phi(B) = [2,0]</math></li></ul> <p><math>\rightarrow</math> Network Weights:</p> <ul style="list-style-type: none"><li>Online network: <math>w_{online} = [0.5,0.1,0.0,0.2]</math></li><li>Target network: <math>w_{target} = [0.2,0.1,0.1,0.1]</math></li><li>Target network is updated once in 2 iterations</li></ul> <p><math>\rightarrow</math> Parameters:</p> <ul style="list-style-type: none"><li>Discount factor: <math>\gamma = 0.9</math></li><li>Learning rate: <math>\alpha = 0.5</math></li></ul>	CO3	PO2	10																																																									

			<b>OR</b>			
	8	a)	<p>You're implementing reinforcement learning agents to solve a complex navigation task in a simulated environment. You are considering two policy-gradient methods: A2C and A3C.</p> <ol style="list-style-type: none"> <li>Outline the key architectural difference between A2C and A3C</li> <li>"A3C typically perform better in highly dynamic or large environments". Justify</li> <li>Illustrate how the update mechanism differ between A2C and A3C.</li> </ol>	CO2	PO1	<b>10</b>
		b)	Derive the mathematical formulation underlying direct policy optimization that maximizes expected return without relying on value function estimation or knowledge of environment dynamics.	CO2	PO1	<b>10</b>
			<b>UNIT - V</b>			
	9	a)	<p>A delivery robot in a warehouse starts at point A, moves to checkpoint B, and then to its destination C, completing one delivery episode. The transitions and rewards are:</p> <ul style="list-style-type: none"> <li>From A to B: reward = +1</li> <li>From B to C: reward = +2</li> <li>C is a terminal state.</li> </ul> <p>The robot wants to estimate the value of being at each location using TD(<math>\lambda</math>) with Discount factor <math>\gamma = 1</math>, Trace decay parameter <math>\lambda = 0.5</math>, Initial values: <math>V(A) = V(B) = V(C) = 0</math>, Learning rate <math>\alpha = 1</math>.</p> <ol style="list-style-type: none"> <li>Using the forward view (<math>\lambda</math>-return), compute the value update for <math>V(A)</math> after this episode.</li> <li>Using the backward view (eligibility traces), compute the value update for <math>V(A)</math> after this episode.</li> <li>Explain why the results are the same, illustrating the equivalence of the two views.</li> </ol>	CO2	PO1	<b>10</b>
		b)	Watkins' $Q(\lambda)$ is one of the earliest attempts to combine Q-learning with eligibility traces in an off-policy setting, yet it exhibits two fundamental shortcomings. Identify these two key limitations of Watkins' $Q(\lambda)$ and briefly explain why each impedes in off policy learning with eligibility traces.	CO2	PO1	<b>10</b>
			<b>OR</b>			
	10	a)	<p>Consider a robot navigating a linear environment with four states: A, B, C, and D (where D is the Terminal state). The robot has only one available action at each state: Move Forward.</p> <p>The state transitions and corresponding rewards are as follows:</p> <ul style="list-style-type: none"> <li>moving from A to B yields a reward of 1,</li> <li>from B to C yields a reward of 2, and</li> <li>from C to D yields a reward of 3.</li> </ul> <p>The robot always follows a fixed policy of moving forward at every state. Initially, the Q-values for all state-action pairs are set</p>	CO2	PO1	<b>10</b>

			<p>to 0. Use the SARSA(<math>\lambda</math>) algorithm with the following parameters: a learning rate (<math>\alpha</math>) of 0.1, a discount factor (<math>\gamma</math>) of 1.0, and a trace decay parameter (<math>\lambda</math>) of 0.8.</p> <ol style="list-style-type: none"> <li>Perform one complete episode starting from state A and ending at the terminal state D. After each step, update the Q-values using replacing eligibility traces.</li> <li>Show the eligibility trace values after each step, update the Q-values accordingly, and finally report the Q-values after the episode ends.</li> </ol>			
		b)	<p>Eligibility traces enhance learning in reinforcement learning by supporting both forward and backward perspectives. Elaborate on the following points:</p> <ol style="list-style-type: none"> <li>The role of eligibility traces in unifying TD and Monte Carlo methods using a spectrum of n-step updates.</li> <li>The function of eligibility traces as temporary memory traces for assigning credit during learning updates.</li> <li>The significance of combining forward and backward views for improved temporal credit assignment and learning performance.</li> </ol>	CO2	PO1	<b>10</b>

\*\*\*\*\*