Roll No. 

# PRESIDENCY UNIVERSITY

## BENGALURU

| End - Term Examinations – MAY 2025 |
|---|
| **Date:** 27-05-2025        **Time:** 09:30 am – 12:30 pm |

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

| CO – Levels | CO1 | CO2 | CO3 | CO4 | CO5 |
|---|---|---|---|---|---|
| Marks | 26 | 26 | 24 | 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

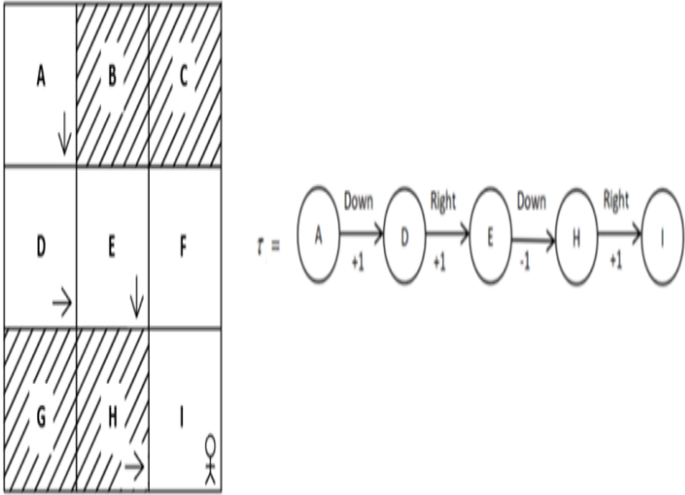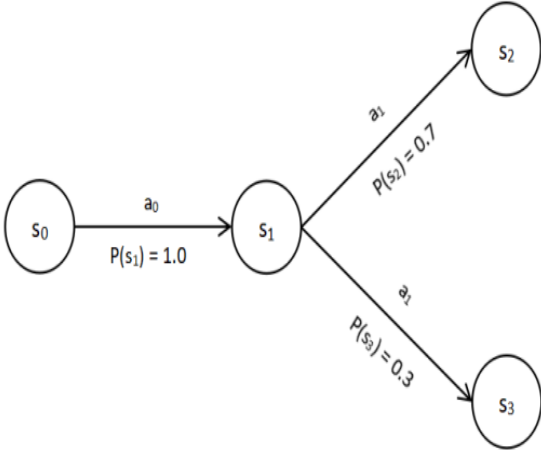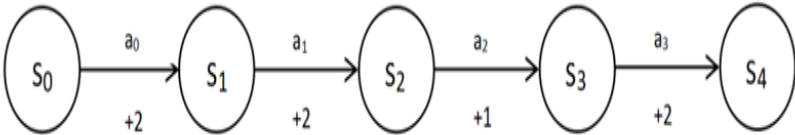**Answer ALL the Questions. Each question carries 2marks.**      **10Q x 2M=20M**

| 1. | Explain the elements of RL | 2 Marks | L2 | CO1 |
|---|---|---|---|---|
| 2. | Define 'reward' and 'return' for an episodic task with an example for each | 2 Marks | L1 | CO1 |
| 3. | Define a) Episode b) Optimal policy | 2 Marks | L1 | CO1 |
| 4. | Define Q function | 2 Marks | L1 | CO2 |
| 5. | Write any two differences between on-policy and off policy TD Control algorithms | 2 Marks | L2 | CO3 |
| 6. | Define Value function | 2 Marks | L1 | CO2 |
| 7. | What is the value of the cards J,4,Q and 'Ace' in the blackjack game? | 2 Marks | L1 | CO3 |
| 8. | What is the significance of T in softmax exploration? | 2 Marks | L1 | CO4 |
| 9. | Write the equation to find V(S) in the Monte Carlo method | 2 Marks | L2 | CO2 |
| 10. | Define Thompson sampling of an arm in MAB Problem | 2 Marks | L1 | CO4 |

# Part B

| 11. | a. | For the following grid world Environment Calculate the Value function which follows a deterministic policy | 10Marks | L3 | CO1 |
|---|---|---|---|---|---|
| | |  Figure 3.2: Transition probability of performing action $a_1$ in state $s_1$ | | | |
| | b. | Identify Bellman equation to the value function of a state in a deterministic environment and stochastic environment . Explain each term in it. Find the value of all the states in the trajectory given below using Bellman equation. Assume Y=1  | 10 Marks | L2 | CO1 |
| | | **Or** | | | |
| 12. | a. | Implement the reinforcement learning environment namely, | 10 | L3 | CO1 |

| | | | Marks | | |
|---|---|---|---|---|---|
| | | Frozen Lake Environment using a random policy and show the output of the following:<br>a. Create and render the environment<br>b. Action Space<br>c. State Space<br>d. Generate 20 Episodes and print Return of each episode | | | |
| | b. | Discuss stochastic environment and deterministic environment in RL with an example. | 10 Marks | L2 | CO1 |

| | | | | | |
|---|---|---|---|---|---|
| 13. | a. | 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<br><br>Create and render the environment<br><br>Generate policy using 20episodes with 30 timesteps<br><br>Print the optimal policy | 10 Marks | L3 | CO2 |
| | b. | Explain different types of RL environments with an example each. | 10Marks | L2 | CO2 |
| | | **Or** | | | |
| 14. | a. | Articulate TD Prediction algorithm in FZLE environment | 10 Marks | L3 | CO2 |
| | b. | Discuss the appropriate situations that are suitable to apply DP,MC or TD methods to learn optimal policy | 10 Marks | L2 | CO2 |

| | | | | | |
|---|---|---|---|---|---|
| 15. | a. | Implement SARSA algorithm to learn the optimal policy in Frozen Lake environment using Python. | 10 Marks | L3 | CO3 |
| | b. | Articulate TD Control algorithm in FZLE environment | 10Marks | L3 | CO3 |
| | | **Or** | | | |
| 16. | a. | Interpret Thompson sampling strategy to overcome the exploration – exploitation dilemma with the algorithm | 10 Marks | L3 | CO3 |
| | b. | Compare SARSA and Q Learning exploration strategies | 10 Marks | L3 | CO3 |

| | | | | | |
|---|---|---|---|---|---|
| 17. | a. | Articulate contextual Bandits and list out the applications of MAB | 10 Marks | L3 | CO4 |

| arm | Q |
|---|---|
| arm 1 | 1 |
| arm 2 | 0 |
| arm 3 | 0 |
| arm 4 | 0 |

| | | | | | |
|---|---|---|---|---|---|
| | b | Briefly Explain the applications of Reinforcement Learning | **10 Marks** | **L2** | **CO4** |
| | | **Or** | | | |
| **18.** | **a.** | Compare advantages and disadvantages of Monte carlo, Dynamic Programming and Temporal Difference in detail | **10 Marks** | **L3** | **CO4** |
| | **b** | Demonstrate Upper Confidence Bound to overcome the exploration –exploitation dilemma with the algorithm | **10Marks** | **L2** | **CO4** |