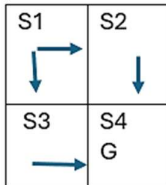


- The question paper contains 5 questions each of 5 marks and total 25 marks.
- Attempt all questions.
- The missing data, if any, may be assumed suitably.
- Tables/Data handbook/Graph paper etc., if applicable, will be supplied to the candidates

- Q.1(a) Explain with neat block diagram the framework of reinforcement learning [2] CO 1 BL 2
- Q.1(b) [3] CO 1 BL 3



Design the tabular representation of a policy for the above grid .the probabilities are 0.5 for a1(R) and a2(D) respectively, the other actions as policies are a3(U),a4(L),a5(still).

- Q.2(a) [2] CO 1 BL 3

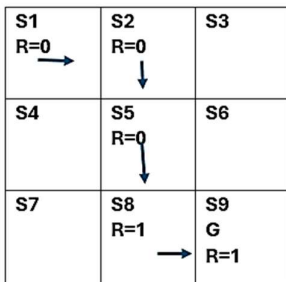


Fig1(a)

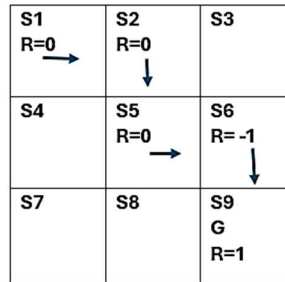


Fig1(b)

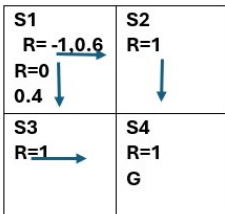
Determine the best policy by considering the finite length and infinite length trajectory. Assume γ as discount factor.

- Q.2(b) Explain mathematically why Markov property refers to memoryless property of a stochastic process. What is the difference between Markov Process and Markov Decision Process. [3] CO 1 BL 3

- Q.3(a) Is the reward a function of next state ? Justify with a mathematical equation [2] CO 1 BL 3
- Q.3(b) Derive Bellman's equation from $G_t=R_{t+1}+\gamma R_{t+2}+\gamma^2 R_{t+3}+\dots$ and reduce it to [3] CO 2 BL 3

$$v_{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} p(s'|s, a) [r(s') + \gamma v_{\pi}(s')]$$

- Q.4(a) What is the relationship between state values and action values? [2] CO 2 BL 2
- Q.4(b) [3] CO 2 BL 3



Determine the matrix form of Bellman's equations for the above grid. For actions refer question 1(b)

- Q.5(a) Derive action values in terms of state values [2] CO 2 BL 1
- Q.5(b) Why do we care about the values of the actions that a given policy cannot select? Referring to the figure given in question 4(b). Determine action values of state s1 in terms of state values .For actions refer question 1(b) [3] CO 2 BL 3