

## **OEOD Quiz Preparation (Questions and Answers)**

### **MCD – Iscte, Diana Mendes – 2023**

---

1. What is Reinforcement Learning (RL)?
  - A. An unsupervised learning approach
  - B. A form of unsupervised learning
  - C. Learning from labeled data
  - D. A machine learning training method based on rewarding desired behaviors and/or punishing undesired ones
  
2. In RL, what represents the learning agent's environment?
  - A. The learner and the decision-maker
  - B. The data used for training
  - C. The model architecture
  - D. The set of actions available
  
3. What is the objective of reinforcement learning?
  - A. To minimize rewards
  - B. To maximize the loss function
  - C. To minimize the policy
  - D. To train an agent to complete a task within an uncertain environment
  
4. What is the action-value function in RL?
  - A. The probability of taking an action
  - B. The immediate reward of an action
  - C. The future reward of an action
  - D. The probability of exploring an action
  
5. What does the "discount factor" in RL determine?

- A. The learning rate
  - B. The agent's exploration rate
  - C. The value of the reward signal over time
  - D. The agent's decision-making speed
6. What is the term for the method in which an RL agent explores the environment to learn optimal actions?
- A. Exploitation
  - B. Generalization
  - C. Exploration
  - D. Policy Optimization
7. In RL, what is a policy?
- A. A set of states
  - B. A sequence of actions
  - C. A mapping of states to actions
  - D. A series of rewards
8. What is the exploration-exploitation trade-off in RL?
- A. Balancing the model complexity
  - B. Balancing the learning rate
  - C. Balancing immediate and future rewards
  - D. Balancing between exploring and exploiting
9. What sets Reinforcement Learning apart from other machine learning paradigms?
- A. Pre-trained models
  - B. Supervised labeling
  - C. Interaction with an environment
  - D. Batch processing
10. What term describes the strategy of choosing actions to maximize cumulative rewards over time?

- A. Hyperparameter tuning
- B. Feature extraction
- C. Reinforcement learning
- D. Policy optimization

11. What does the term “agent” refer to in Reinforcement Learning?

- A. A person supervising the learning process
- B. A software program making decisions
- C. A labeled data point
- D. A neural network architecture

12. What numerical values are used to evaluate the outcomes of actions taken by an agent?

- A. Observations
- B. Rewards
- C. Policies
- D. Loss functions

13. The “reward function” in RL is used for:

- A. Defining the neural network architecture
- B. Calculating the probability of actions
- C. Evaluating the quality of an agent’s actions
- D. Filtering noisy observations

14. Which algorithm is particularly well-suited for environments with continuous action spaces?

- A. Q-Learning
- B. Deep Q-Network (DQN)
- C. Policy Gradient
- D. Monte Carlo Tree Search (MCTS)

15. Which real-world applications benefit from Reinforcement Learning?

- A. Image classification
- B. Text generation
- C. Autonomous driving
- D. Data clustering

16. What is a sequential decision problem?

17. What is the Markov property?

18. What is a policy  $\pi(a|s)$ ?

19. What's on-policy?

20. What is Grid world?

21. What is  $Q(s, a)$ ?

22. Which model or function is meant when we say "model-free" or "model-based"?

23. What type of action space and environment suit value-based methods?

23 Reinforcement learning is a \_\_\_\_\_

- A. Prediction-based learning technique
- B. Feedback-based learning technique
- C. History results-based learning technique

24. Which kind of data does reinforcement learning use?

- A. Labeled data
- B. Unlabelled data
- C. None
- D. Both

25. Reinforcement learning methods learned through \_\_\_\_?
- A. Experience
  - B. Predictions
  - C. Analyzing the data
26. Which of the following is the practical example of reinforcement learning?
- A. House pricing prediction
  - B. Market basket analysis
  - C. Text classification
  - D. Driverless cars
27. What is an agent in reinforcement learning?
- A. An agent is the situation in which rewards are being exchanged
  - B. An agent is the simple value in reinforcement learning.
  - C. An agent is an entity that explores the environment.
28. What are actions in reinforcement learning?
- A. Actions are the moves that the agent takes inside the environment.
  - B. Actions are the functions that the environment takes.
  - C. Actions are the feedback that an agent provides.
29. In which of the following reinforcement learning approaches do we find the optimal value function?
- A. Value-based
  - B. Policy-based
  - C. Model-based
30. The agent's main objective is to \_\_\_\_ the total number of rewards for good actions?
- A. Minimize
  - B. Maximize
  - C. Null

31. Reinforcement learning is defined by the \_\_\_\_?
- A. Policy
  - B. Reward Signal
  - C. Value Function
  - D. Model of the environment
32. Which element in reinforcement learning defines the behavior of the agent?
- A. Policy
  - B. Reward Signal
  - C. Value Function
  - D. Model of the environment
33. Can reward signals change the policy?
- A. Yes
  - B. No
34. Which of the following elements of reinforcement learning imitates the behavior of the environment?
- A. Policy
  - B. Reward Signal
  - C. Value Function
  - D. Model of the environment
35. The approach in which reinforcement learning problems are solved with the help of models is known as \_\_\_\_?
- A. Model-based approach
  - B. Model-free approach
  - C. Model known approach
36. Gamma ( $\gamma$ ) in the bellman equation is known as?
- A. Value factor
  - B. Discount factor

C. Environment factor

37. How do you represent the agent state in reinforcement learning?

- A. Discount state
- B. Discount factor
- C. Markov state

38. Consider the following condition:  $P[S(t+1) | S(t)] = P[S(t+1) | S(1), \dots, S(t)]$ . What is the meaning of  $S(t)$ ?

- A. State factor
- B. Discount factor
- C. Markov state

39. What do you mean by MDP in reinforcement learning?

- A. Markov discount procedure
- B. Markov discount process
- C. Markov deciding procedure
- D. Markov decision process

40. Why do we use MDP in reinforcement learning?

- A. We use MDP to formalize the reinforcement learning problems.
- B. We use MDP to predict reinforcement learning problems.
- C. We use MDP to analyze the reinforcement learning problems.

41. What do you mean by SARSA in reinforcement learning?

- A. State action reward state action
- B. State achievement rewards state action
- C. State act reward achievement
- D. State act reward act

42. Which of the following policy types is a learning algorithm that evaluates and improves a policy dissimilar from the policy used for action selection?

- A. behavior policy
  - B. Target policy
  - C. On-policy
  - D. Off-policy
43. Q-learning is a model-free or model-based learning algorithm?
- A. Model-free
  - B. Model-based
44. The matrix created during the Q-learning algorithm is commonly known as \_\_\_\_?
- A. Query-table
  - B. Q-table
  - C. Quick-matrix
  - D. Table
45. Does reinforcement learning provide any previous training?
- A. Yes
  - B. No
46. Reinforcement Learning is about how computer systems can learn to \_\_\_\_\_.
47. The Markov Assumption assumes that to predict the future, you'll use \_\_\_\_\_.
48. What is the difference between state and history?
49. How is a policy evaluated?
50. What is a transition matrix for Markov Process?
51. What is a Bandit?



- A. Nickname for slot-machine
- B. Nickname for Black-Jack
- C. Losing reward

52. Which is the difference between Bellman (Value) Equation and Bellman (Value) Optimality Equation?

53. What is policy iteration?

54. What is value iteration?

55. What does  $P(s' | s, a)$  stand for?

56. How do we calculate the  $Q(s, a)$  of a given policy  $\pi_i$ ?

57. How do we calculate Total Discounted Reward?

58. What is the difference between \*Off-line\* vs \*On-line\* learning?

59. In multiarm bandits, what is regret, and why do we use it?

60. How do we calculate Reward for MABs?

61. Monte Carlo policy evaluation is a good choice when we know \_\_\_\_\_ about the dynamics and/or reward model.

62. Monte Carlo assumes \_\_\_\_\_ {one-shot, episodic, infinite} MDP.

63. An MDP given a fixed policy is a Markov chain with rewards.

64. If we know the optimal  $Q$  values, we can get the optimal  $V$  values only if we know the environment's transition function/matrix.

65. A policy that is greedy—with respect to the optimal value function—is not necessarily an optimal policy.

66. Why is the discount factor usually smaller than 1?

67. What is the idea of the Bellman equation

68. Define the Return

69. Define the state-value function

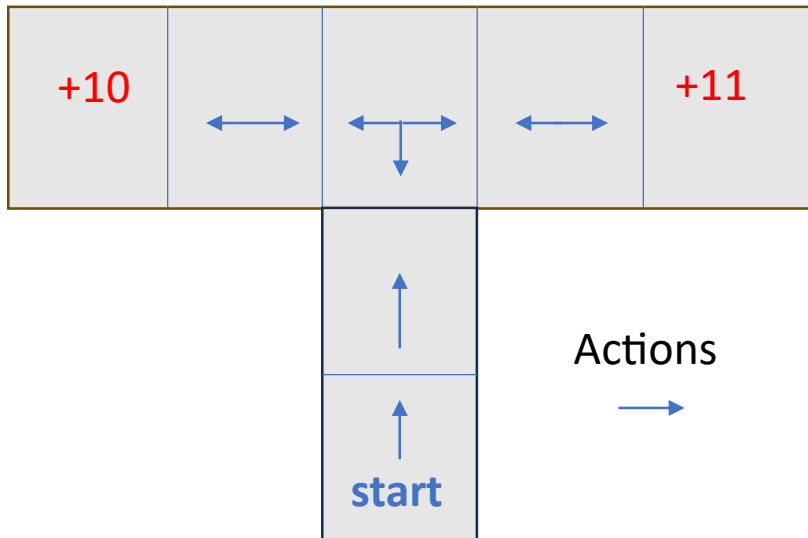
70. What is the direct solution of the Bellman equation for MDP's using matrix/vector form?

71. What happens when gamma, the discount factor, is close to 0 and 1

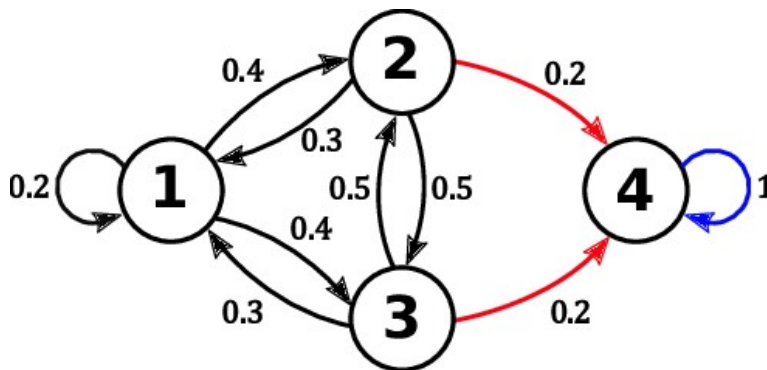
72. What is a stochastic environment?

73. Calculate the value function  $V^*$  from Bellman equation for the maze shown in the image below.

A simple T-maze with 7 states. The start location is at the bottom, the possible action of each state are indicated as arrows. When arriving at the left or right end the agent obtains a reward of 10 or 11 reward units, respectively and gets 'transferred' to the starting state. Consider each action as equally likely when writing down the Bellman equation and use  $\gamma = 0.1$  for your numerical solution.

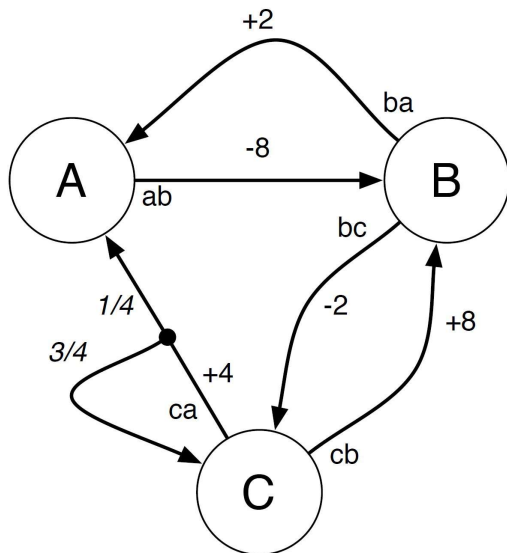


74. Consider the following transition graph associated with a Markov chain. Write down the probability transition matrix  $P$ .



75. Consider the following Markov Decision Process (MDP) with discount factor  $\gamma = 0.5$ . Upper case letters A, B, C represent states; arcs represent state transitions; lower case letters ab; ba; bc; ca; cb represent actions; signed integers fractions represent transition probabilities.

- Define the state-value function  $V(s)$  for the MDP
- Write the Bellman optimal equation for state-value functions
- Starting with an initial value function of  $V_1(A) = V_1(B) = V_1(C) = 2$ , apply one iteration of value iteration (i.e. one backup for each state) to compute a new value function  $V_2(s)$ .



76. Consider the following transition graph with Rewards:

- Compute the Return for the sequence 1-1-2-3-Exit, for  $\gamma=0.8$
- Compute the state-value function at state 2 for the previous example.

