

1. The learner and decision maker is the _____.

1 point

- ☐ Reward
- ☐ State
- ☐ Environment
- ☒ Agent

2. At each time step the agent takes an _____.

1 point

- ☐ State
- ☐ Reward
- ☒ Action
- ☐ Environment

3. Imagine the agent is learning in an episodic problem. Which of the following is true?

1 point

- ☐ The number of steps in an episode is always the same.
- ☐ The agent takes the same action at each step during an episode.
- ☒ The number of steps in an episode is stochastic: each episode can have a different number of steps.

4. If the reward is always +1 what is the sum of the discounted infinite return when $\gamma < 1$

1 point

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- ☐ Infinity.
- ☐ $G_t = \frac{\gamma}{1-\gamma}$
- ☒ $G_t = \frac{1}{1-\gamma}$
- ☐ $G_t = 1 * \gamma^k$

5. How does the magnitude of the discount factor (γ) affect learning?

1 point

- ☒ With a smaller discount factor the agent is more far-sighted and considers rewards farther into the future.
- ☐ The magnitude of the discount factor has no effect on the agent.
- ☐ With a larger discount factor the agent is more far-sighted and considers rewards farther into the future.

6. Suppose $\gamma = 0.8$ and we observe the following sequence of rewards: $R_1 = -3$, $R_2 = 5$, $R_3 = 2$, $R_4 = 7$, and $R_5 = 1$, with $T = 5$. What is G_0 ? Hint: Work Backwards and recall that $G_t = R_{t+1} + \gamma G_{t+1}$.

1 point

- ☐ 11.592
- ☐ 12
- ☐ -3
- ☐ 8.24
- ☒ 6.2736

7. What does MDP stand for?

1 point

- ☐ Markov Deterministic Policy
- ☐ Meaningful Decision Process
- ☐ Markov Decision Protocol
- ☒ Markov Decision Process

8. Consider using reinforcement learning to control the motion of a robot arm to pick up objects and place them into new positions. The actions in this case might be the voltages applied to each motor at each joint, and the states might be the latest readings of joint angles and velocities. The reward might be +1 for each object successfully picked up and placed. To encourage smooth movements, on each time step a small, negative reward can be given as a function of the moment-to-moment “jerkiness” of the motion. Is this a valid MDP?

1 point

- ☒ Yes

☐ No

9. **Case 1:** Imagine that you are a vision system. When you are first turned on for the day, an image floods into your camera. You can see lots of things, but not all things. You can't see objects that are occluded, and of course you can't see objects that are behind you. After seeing that first scene, do you have access to the Markov state of the environment?

1 point

Case 2: Imagine that the vision system never worked properly: it always returned the same static image, forever. Would you have access to the Markov state then? (Hint: Reason about $P(S_{t+1}|S_t, \dots, S_0)$, where $S_t = \text{AllWhitePixels}$)

- ☒ You have access to the Markov state in both Case 1 and 2.
- ☐ You have access to the Markov state in Case 1, but you don't have access to the Markov state in Case 2.
- ☐ You don't have access to the Markov state in Case 1, but you do have access to the Markov state in Case 2.
- ☐ You don't have access to the Markov state in both Case 1 and 2.

10. What is the reward hypothesis?

1 point

- ☐ That all of what we mean by goals and purposes can be well thought of as the minimization of the expected value of the cumulative sum of a received scalar signal (called reward)
- ☐ Always take the action that gives you the best reward at that point.
- ☐ Ignore rewards and find other signals.
- ☒ That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward)

11. Imagine, an agent is in a maze-like gridworld. You would like the agent to find the goal, as quickly as possible. You give the agent a reward of +1 when it reaches the goal and the discount rate is 1.0, because this is an episodic task. When you run the agent it finds the goal, but does not seem to care how long it takes to complete each episode. How could you fix this? (**Select all that apply**)

1 point

- ☐ Give the agent a reward of +1 at every time step.
- ☐ Give the agent a reward of 0 at every time step so it wants to leave.

- ☒ Give the agent -1 at each time step.
- ☒ Set a discount rate less than 1 and greater than 0, like 0.9.

12. When may you want to formulate a problem as episodic?

- ☐ When the agent-environment interaction does not naturally break into sequences. Each new episode begins independently of how the previous episode ended.
- ☒ When the agent-environment interaction naturally breaks into sequences. Each sequence begins independently of how the episode ended.