

1. Which of the following are true? (Select all that apply)

1 / 1 point

☒ In the tabular case, updating the value of one state does not affect the values of other states.

✓ **Correct**

Correct! In the tabular case, there is one estimate of the value function for each state, so updating the value function of one state does not affect the value of other states.

☐ When using state aggregation or coarse coding, updating the value of one state does not affect the values of other states.

☐ In the tabular case, there is generalization across states, i.e., updates to the value function of one state influences the value function of other states.

☒ When using state aggregation or coarse coding, there is generalization across states.

✓ **Correct**

Correct: Since states share weights, updating the value of one state also affects the value of other states; hence, learning generalizes across states.

2. To turn the update of Expected Sarsa algorithm to the update of Q-learning, one must:

1 / 1 point

- ☒ Use the maximum over all the actions instead of the expectation in the update function.
- ☐ Behave greedily with respect to the action-value function.
- ☐ Expected Sarsa cannot be adapted to represent Q-learning.
- ☐ Use a neural network to approximate the action-value function.

✓ **Correct**

Correct: The difference between the update of Expected and the update of Q-learning is that the expectation over all the actions is replaced with the maximum over all the actions.

3. Which of the following are true:

1 / 1 point

- ☐ When using function approximation and with discrete actions, there is no straightforward way to turn Sarsa into Expected Sarsa.
- ☐ Exploration is not a problem when using function approximation because learning generalizes across states and actions.
- ☐ When using function approximation, Q-learning will always result in better performance than Sarsa and Expected Sarsa.
- ☒ When learning value functions in the mountain car domain and using tile coding, where all the rewards are -1 except when the terminal state is reached, initializing the weights of a linear function approximation to zero is an example of optimistic initialization.

✓ **Correct**

Correct: In the mountain car domain, all the true action-values functions are negative. Initializing the weights to zero in the linear function approximation case is equivalent to initializing all the action-values to zero, which is greater than the true values of the action-value function. Hence, this is an example of optimistic initialization.

4.

1 / 1 point

Which of the following statements about discounted return and differential return algorithms are true:

☒ The performance of discounted return algorithms can suffer because the optimal value of gamma is problem dependent.

 **Correct**

Correct: Remember the two-choice MDP in Lesson 3 of this module. In such MDPs, different values of the discount factor result in different policies.

☒ Algorithms that maximize the discounted return can suffer from an exponentially large variance when using a large discount factor.

 **Correct**

Correct: The longer the time horizon of the return in the discounted setting — the larger the value of discount factor — the wider the range of values of the return. Hence, the variance is larger for larger values of discount factor.

☐ Discounted return algorithms can diverge if the average reward term is slightly larger or slightly lower than the true value.

☐ There is a set of Bellman equations for the value function corresponding to the discounted return. However, there is no set of Bellman equations for differential value functions.

5. Imagine we have 7 state features and 6 actions and we want to use linear function approximation to compute the action-value function. If we use feature stacking (as explained in Week 3), how many features would we need in total to approximate the action-value function of all the different actions?

1 / 1 point

☐ 13☐ 67☒ 42☐ 76 **Correct**

Correct! When stacking features in linear function approximation, we must stack them as many times there are actions.

6.

1 / 1 point

Imagine we are approximating the action-value function of three different actions (red, green, and blue) using three state features, resulting in the following weights and features:

$$\mathbf{x}(s) = \begin{bmatrix} x_0(s) \\ x_1(s) \\ x_2(s) \end{bmatrix} \quad \mathbf{x}(s, a) = \begin{bmatrix} \begin{matrix} x_0(s) \\ x_1(s) \\ x_2(s) \end{matrix} \bigg\} a_0 \\ \begin{matrix} x_0(s) \\ x_1(s) \\ x_2(s) \end{matrix} \bigg\} a_1 \\ \begin{matrix} x_0(s) \\ x_1(s) \\ x_2(s) \end{matrix} \bigg\} a_2 \end{bmatrix} \quad \mathbf{w}(s) = \begin{bmatrix} \begin{matrix} w_0(s) \\ w_1(s) \\ w_2(s) \end{matrix} \bigg\} a_0 \\ \begin{matrix} w_0(s) \\ w_1(s) \\ w_2(s) \end{matrix} \bigg\} a_1 \\ \begin{matrix} w_3(s) \\ w_4(s) \\ w_5(s) \end{matrix} \bigg\} a_2 \end{bmatrix}$$

Which of these statements is true about our function approximation scheme?

☒ The action-value of the green action will be exactly the same as the action-value of the red action for any state.

✓ **Correct**

Correct: Since the red and green actions share the exact same weights, computing the dot product between the feature vector of each action and their corresponding weights will result in the same value.

☐ The action-value of the red action will be exactly the same as the action-value of the blue action.

☐ The green action will never be selected.

☐ There will not be any generalization across states.

7. Consider the following feature and weight vectors corresponding to three state features and three actions (red, green, and blue):

0 / 1 point

$$\mathbf{x}(s) = \begin{bmatrix} x_0(s) \\ x_1(s) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.2 \end{bmatrix} \quad \mathbf{x}(s, a) = \begin{bmatrix} x_0(s) \\ x_1(s) \end{bmatrix} \begin{matrix} \left. \vphantom{\begin{matrix} x_0(s) \\ x_1(s) \end{matrix}} \right\} a_0 \\ \left. \vphantom{\begin{matrix} x_0(s) \\ x_1(s) \end{matrix}} \right\} a_1 \end{matrix}$$

$$\begin{bmatrix} \tilde{x}_1(s) \\ x_2(s) \end{bmatrix} = \begin{bmatrix} \tilde{x}_1(s) \\ 0.3 \end{bmatrix} + \underbrace{\begin{bmatrix} \tilde{x}_1(s) \\ x_2(s) \\ x_0(s) \\ x_1(s) \\ x_2(s) \end{bmatrix}}_{a_2}$$

$$\mathbf{w}(s) = \begin{bmatrix} w_0(s) \\ w_1(s) \\ w_2(s) \\ w_3(s) \\ w_4(s) \\ w_5(s) \\ w_6(s) \\ w_7(s) \\ w_8(s) \end{bmatrix} = \begin{bmatrix} 30 \\ 60 \\ 50 \\ 4 \\ 4 \\ 4 \\ 8 \\ 15 \\ 30 \end{bmatrix}$$



What is the approximate action-value for the green action?

- ☐ 4
- ☒ 42
- ☐ 140
- ☐ 16

**✗ Incorrect**

Incorrect: Remember that to compute the approximate action-value of an action you need to compute the dot product between the state features and the weights corresponding to the action.

8. Which of the following statements about epsilon-greedy policies and optimistic initialization are true:

**1 / 1 point**

- ☒ Optimistic initialization results in a more systematic exploration in the tabular case because the agent takes actions it has not taken as often in a state.

**✓ Correct**

Correct: Optimistic initialization is like the agent imagining that unexplored actions will result in more reward than they would in reality. Systematically taking actions in each state causes the agent to reach states and areas of the state space it has also seen less often.

- ☐ Implementing optimistic initialization is always straightforward and simple to implement regardless of the function approximation technique.
- ☒ Epsilon-greedy can be easily combined with any function approximation technique because it only needs to be able to query for the approximate action-values, without needing to know how they are initialized or computed.

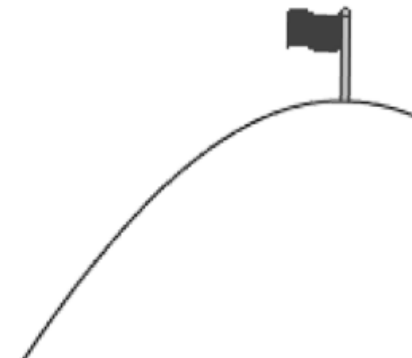
✓ **Correct**

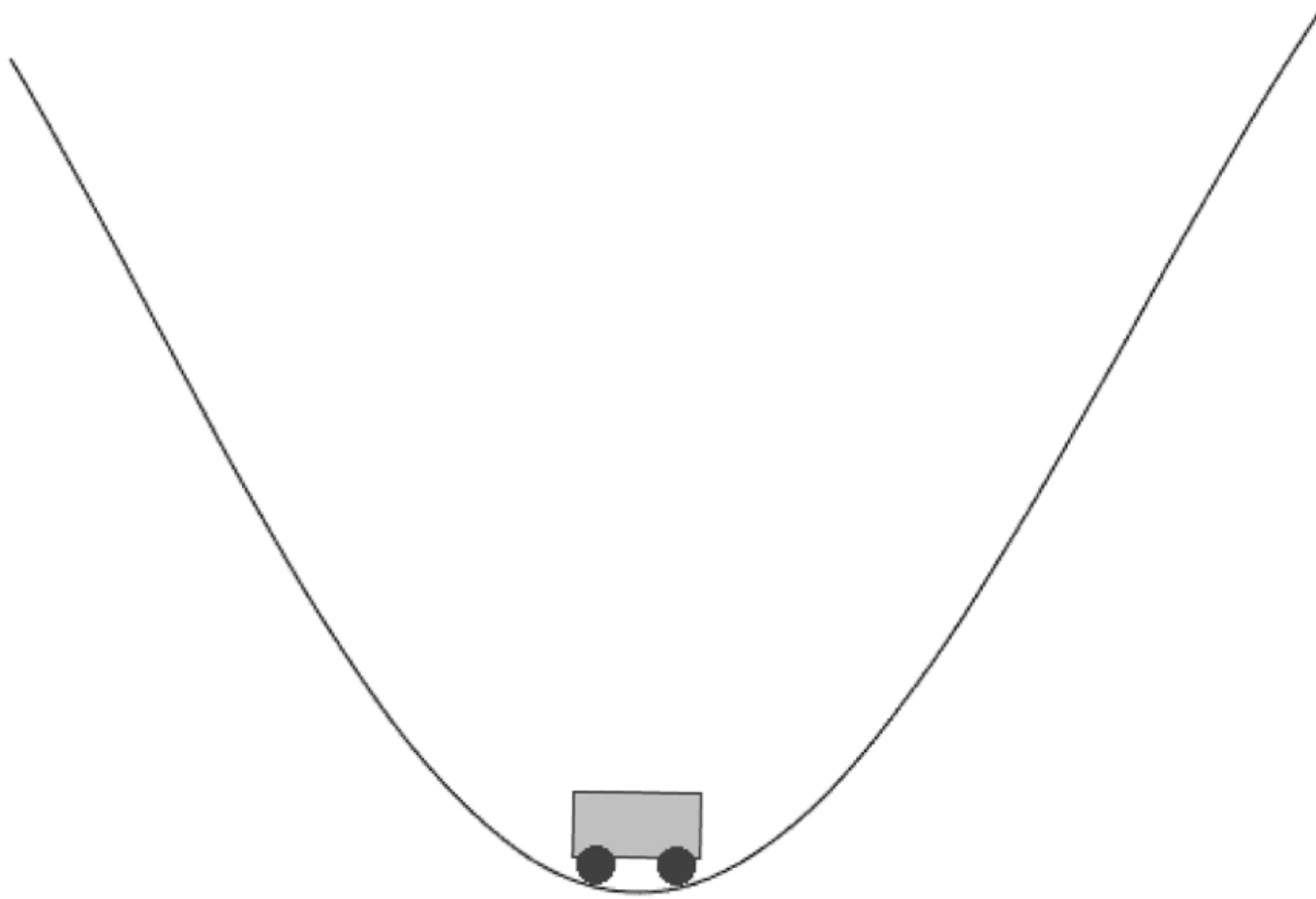
Correct: Epsilon-greedy is simple because it does not involve any special scheme for randomizing the actions or initializing the weights of the function approximation. It's as simple as rolling a die!

- ☐ Optimistic initialization is always preferred over using an epsilon-greedy policy.

9. Consider the mountain car environment where all the rewards are -1 after every action until reaching the flag on top of the right hill:

1 / 1 point





Which of the following are examples of optimistic initialization when using linear function approximation?

- ☐ Initialize all the weights at random using a uniform distribution between -5 and 5.

☒ Initialize all the weights so that the action-value function is 10 everywhere.

✓ **Correct**

Correct: Since the true value of the action-value function is less than or equal to -1 in the mountain car environment, this would result in an initial value considerably greater than the true value.

☒ Initialize all the weights so that the action-value function is zero everywhere.

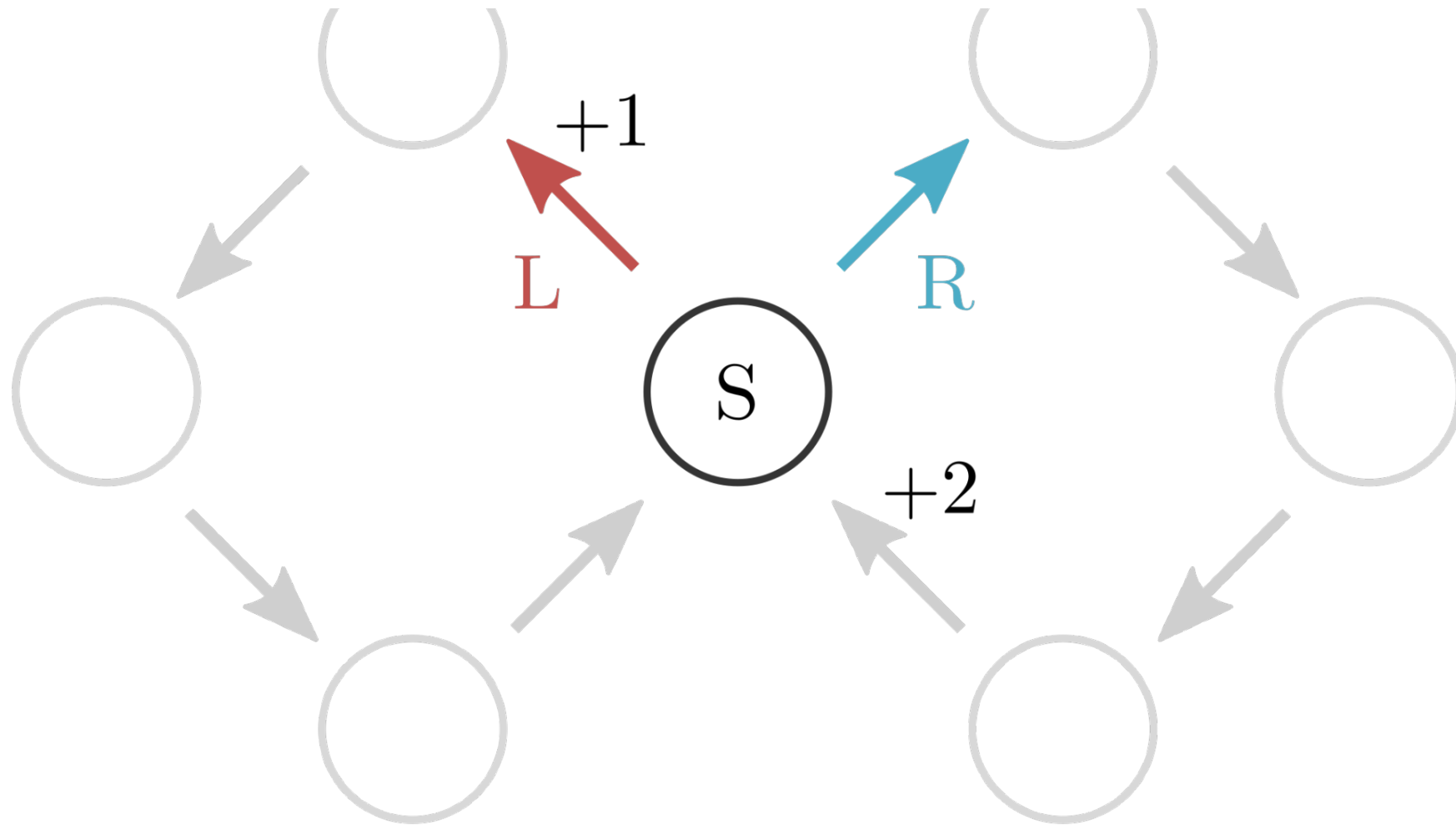
✓ **Correct**

Correct: Since the true value of the action-value function is less than or equal to -1 in the mountain car environment, this would result in an initial value slightly greater than the true value.

☐ Initialize all the weights so that the action-value function is -1 everywhere.

10. Consider the following MDP where the red action (L) at state S results in an immediate reward of +1 and 0 afterwards, and the blue action (R) results in 0 immediate reward, but a final reward of +2 when returning to state S.

1 / 1 point



In this MDP, what is the optimal policy at state  $S$  when using discounted return with a value of  $\gamma$  of 0.5? When using discounted return with a value of  $\gamma$  of 0.9? When using average reward? (Hint: You can modify the equations used for the task from lesson 3 of this module to get the solution)

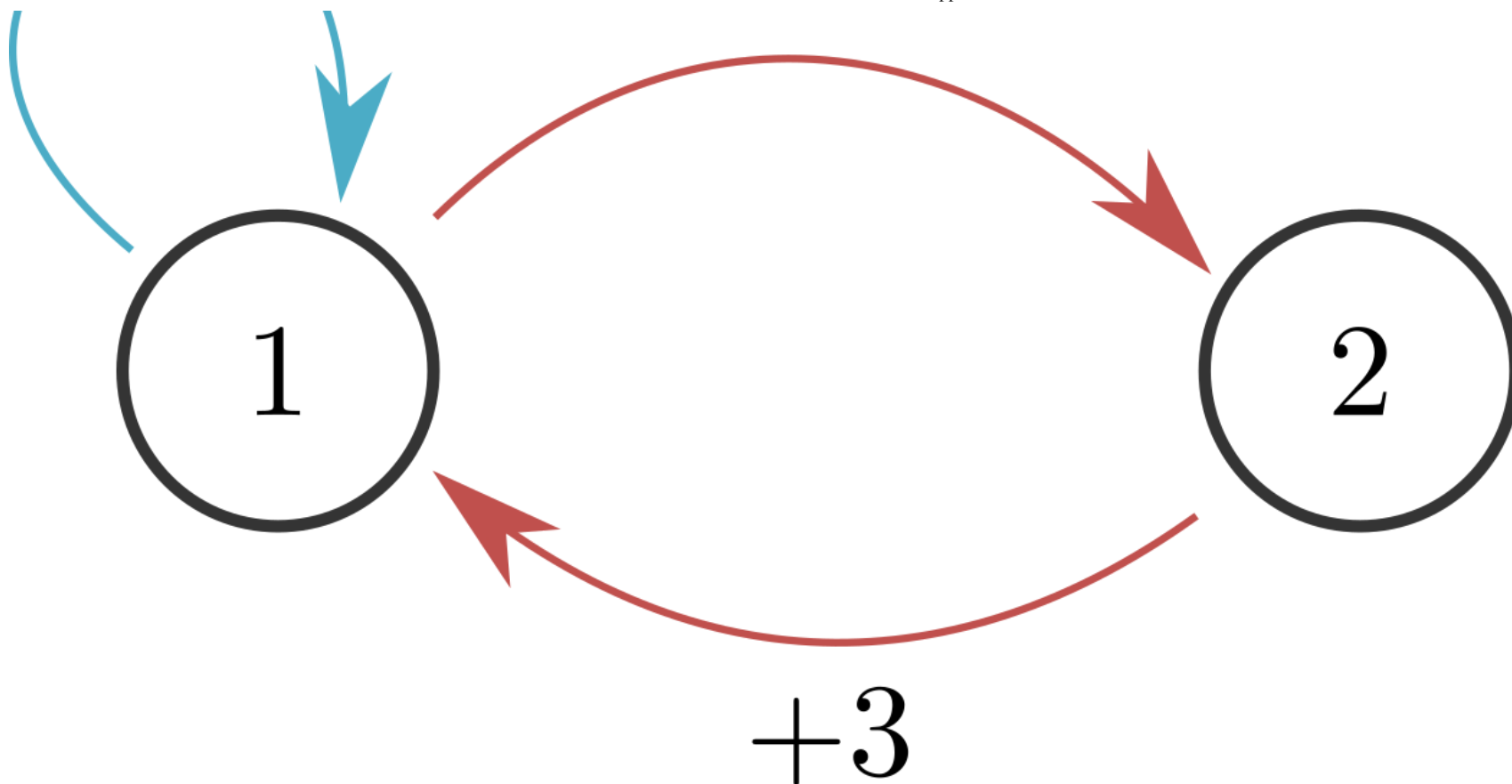
- ☐ The optimal policy is:  $\gamma = 0.5$ , blue action;  $\gamma = 0.9$ , red action; average reward, red action.
- ☒ The optimal policy is:  $\gamma = 0.5$ , red action;  $\gamma = 0.9$ , blue action; average reward, blue action.
- ☐ The optimal policy is:  $\gamma = 0.5$ , blue action;  $\gamma = 0.9$ , red action; average reward, blue action.
- ☐ The optimal policy is:  $\gamma = 0.5$ , blue action;  $\gamma = 0.9$ , blue action; average reward, red action.

✓ **Correct**  
Correct!

11. Consider the following MDP with two states. State 1 has two actions: staying in the same state (blue) and switching to the other state (red). State 2 has a single (red) action leading to State 1. The rewards are listed next to each arrow/transition.

1 / 1 point





If the agent is following a policy that takes the blue and red actions with equal probability, what is the average reward? (Hint: the state visitation probability is  $\frac{2}{3}$  for state 1 and  $\frac{1}{3}$  for state 2, the target policy probability  $\pi$  is the random policy, and all the transitions  $p$  are deterministic) Recall the formula is:

$$r(\pi) = \sum_s \mu_\pi(s) \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) r$$

☒ Average reward =  $4/3$

☐ Average reward =  $3/4$

☐ Average reward = 2

☐ Average reward = 1

✓ **Correct**  
Correct!

12. For which of the following tasks is the average reward setting preferable than the discounted return setting if we are interested in maximizing the total amount of reward?

1 / 1 point

☐ An Atari 2600 game where the agent can keep playing until it loses all its lives, which are a finite amount.



- ☒ The job scheduling task from Section 10.3 of Sutton and Barto's RL textbook. The agent manages how jobs are scheduled into three different servers according to the priority of each job. When the agent accepts a job, the job runs in one of the servers and the agent gets a positive reward equal to the priority of the job, whereas when the agent rejects a job, the job returns to the queue and the agent receives a negative reward proportional to the priority of the job. The servers become available as they finish their jobs. Jobs are continually added to the queue and the agent accepts or rejects them.

✓ **Correct**

Correct: Since this is a continuing task, the agent will maximize the total amount of reward in the average reward setting, whereas in the discounted setting, the agent will maximize the total reward only for some values of discount factor less than one.

- ☐ The episodic mountain car environment.

- ☒ The nearsighted MDP from Week 3. The agent starts at an initial state  $S$  that splits into two paths, each leading to a ring that leads back to  $S$  after some time. If the agent chooses the left ring, it observes an immediate reward of +1 followed by 0 rewards, whereas choosing the right ring results in an immediate reward of 0, but a final reward of +2 when returning to state  $S$ .

✓ **Correct**

Correct: The agent will always maximize the total reward it receives when using the average reward setting, whereas for the discounted reward setting, the agent will maximize the total reward received only for some values of the discount factor