

# **Course 1, Module 3**

# **Markov Decision Processes**

# **Worksheet Questions**

CMPUT 397  
Fall 2019

# Admin

- You had to submit three MDPs by last night
- You need to grade your peers by Sunday night
- As usual, also have to complete Practice Quiz and Discussion Question
- Come to Instructor Office hours
- Worksheet posted for Week 2 (MDPs); we will sometimes put extra questions on the worksheet, since some like to have more problems to work through.

# Worksheet Question 1

(Exercise 2.2 from S&B 2nd edition) Consider a  $k$ -armed bandit problem with  $k = 4$  actions, denoted 1, 2, 3, and 4. Consider applying to this problem a bandit algorithm using  $\epsilon$ -greedy action selection, sample-average action-value estimates, and initial estimates of  $Q_1(a) = 0$ , for all  $a$ . Suppose the initial sequence of actions and rewards is  $A_1 = 1, R_1 = 1, A_2 = 2, R_2 = 1, A_3 = 2, R_3 = 2, A_4 = 2, R_4 = 2, A_5 = 3, R_5 = 0$ . On some of these time steps the  $\epsilon$  case may have occurred, causing an action to be selected at random. On which time steps did this definitely occur? On which time steps could this possibly have occurred?

# Worksheet Question 2

Suppose  $\gamma = 0.9$  and the reward sequence is  $R_1 = 2, R_2 = -2, R_3 = 0$  followed by an infinite sequence of 7s. What are  $G_1$  and  $G_0$ ?

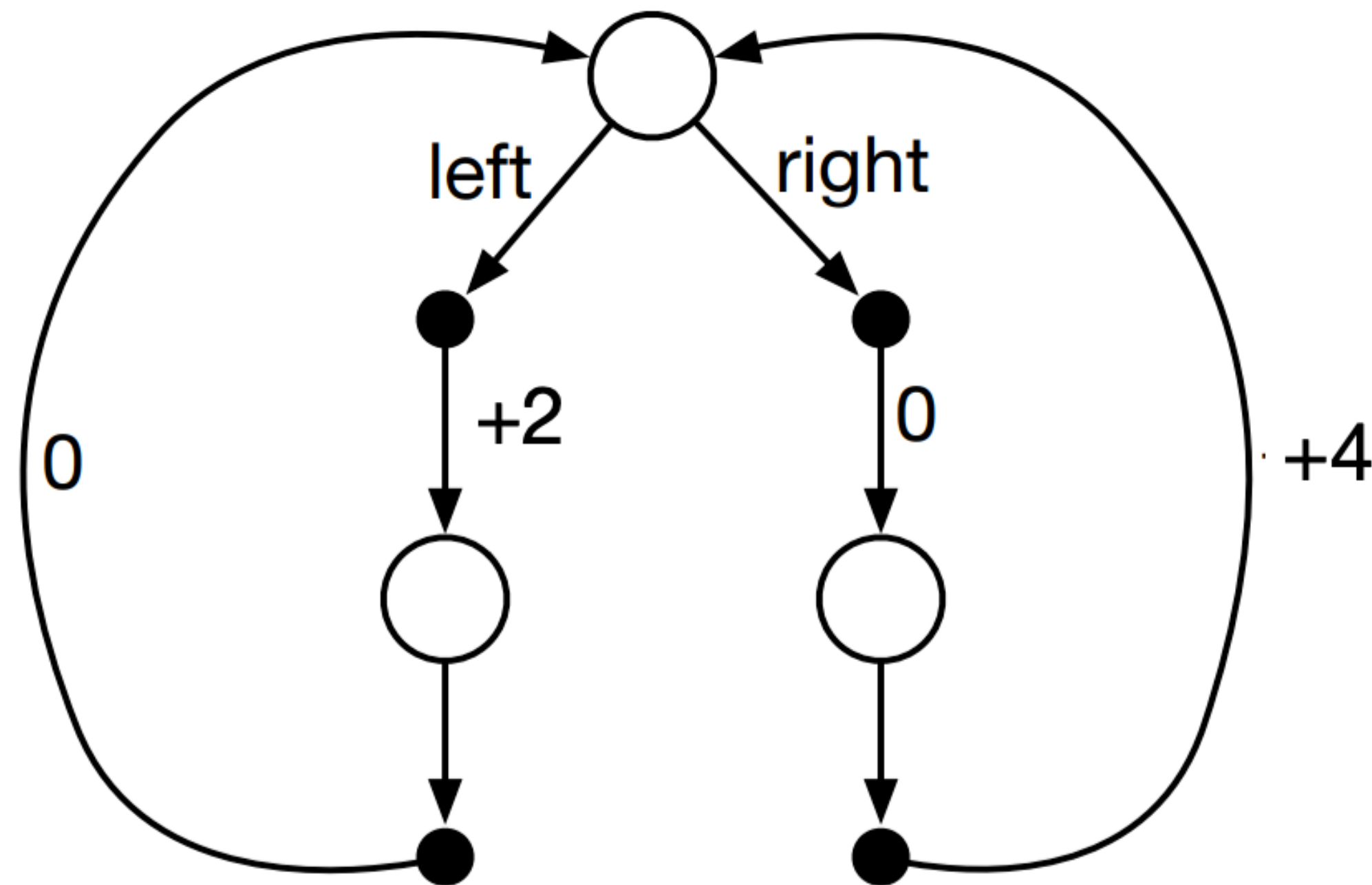
# Worksheet Question 3

Assume you have a bandit problem with 4 actions, where the agent can see rewards from the set  $\mathcal{R} = \{-3.0, -0.1, 0, 4.2\}$ . Assume you have the probabilities for rewards for each action:  $p(r|a)$  for  $a \in \{1, 2, 3, 4\}$  and  $r \in \{-3.0, -0.1, 0, 4.2\}$ . How can you write this problem as an MDP? Remember that an MDP consists of  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, P, \gamma)$ .

**More abstractly**, recall that a Bandit problem consists of a given action space  $\mathcal{A} = \{1, \dots, k\}$  (the  $k$  arms) and the distribution over rewards  $p(r|a)$  for each action  $a \in \mathcal{A}$ . Specify an MDP that corresponds to this Bandit problem.

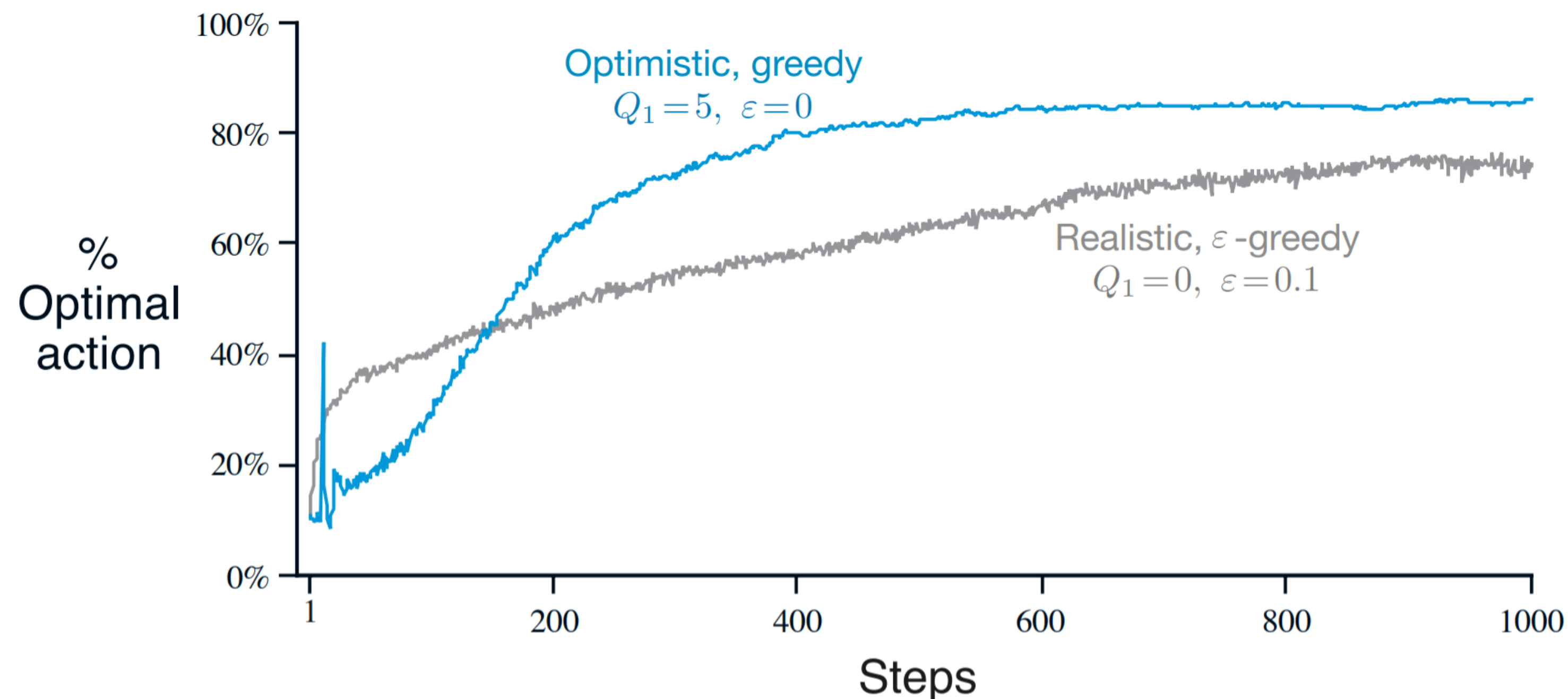
# Worksheet Question 5

Consider the continuing MDP shown on the bottom. The only decision to be made is that in the top state, where two actions are available, left and right. The numbers show the rewards that are received deterministically after each action. There are exactly two deterministic policies,  $\pi_{\text{left}}$  and  $\pi_{\text{right}}$ . What policy is optimal if  $\gamma = 0$ ? If  $\gamma = 0.9$ ? If  $\gamma = 0.5$ ?



# Worksheet Challenge Question

(Exercise 2.6 from S&B 2nd edition) The results shown in Figure 2.3 should be quite reliable because they are averages over 2000 individual, randomly chosen 10-armed bandit tasks. Why, then, are there oscillations and spikes in the early part of the curve for the optimistic method? In other words, what might make this method perform particularly better or worse, on average, on particular early steps?



**alpha = 0.1**