

1. 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  $G_0$ ?

2. 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.

3. Prove that the discounted sum of rewards is always finite, if the rewards are bounded:  $|R_{t+1}| \leq R_{\max}$  for all  $t$  for some finite  $R_{\max} > 0$ .

$$\left| \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \right| < \infty \quad \text{for } \gamma \in [0, 1)$$

Hint: Recall that  $|a + b| < |a| + |b|$ .