

NAME: KARMA TARAP
CSCI E-89c Deep Reinforcement Learning
Part I of Assignment 4

Please consider a Markov Decision Process with $\mathcal{S} = \{s^A, s^B, s^C\}$.

Given a particular state $s \in \mathcal{S}$, the agent is allowed to either try staying there or switching to one of the “neighboring” states. Let’s denote an intention to stay by 0, an attempt to move to the left by -1 , and an intention to move to the right by $+1$. The agent does not know transition probabilities, including the distributions of rewards.

Suppose the agent chooses the policy $\pi(+1|s) = 1$ for $s \in \{s^A, s^B\}$ and $\pi(0|s^C) = 1$ and runs the First-Visit Monte Carlo (MC) prediction algorithm for estimating $q_\pi(s, a)$. If the agent observes the following two episodes:

episode 1:

$S_0 = s^B, A_0 = 0, R_1 = 20, S_1 = s^B, A_1 = +1, R_2 = 30, S_2 = s^B, A_2 = +1, R_3 = 20, S_3 = s^B, A_3 = +1, R_4 = 10, S_4 = s^C, A_4 = 0, R_5 = 130,$

episode 2:

$S_0 = s^B, A_0 = -1, R_1 = 10, S_1 = s^A, A_1 = +1, R_2 = 20, S_2 = s^B, A_2 = +1, R_3 = 10, S_3 = s^B, A_3 = +1, R_4 = 30, S_4 = s^B, A_4 = +1, R_5 = 10,$

assuming $\gamma = 0.9$, find the First-Visit MC estimates of

- (a) $q_\pi(s^B, 0)$
- (b) $q_\pi(s^B, +1)$
- (c) $q_\pi(s^B, -1)$

SOLUTION:

$$q_\pi(s^B, 0) \text{ episode 1} = 20 + 0.9 * 30 + 0.9^2 * 20 + 0.9^3 * 10 + 0.9^4 * 130 = 155.783$$

$$q_\pi(s^B, +1) \text{ episode 1} = 30 + 0.9 * 20 + 0.9^2 * 10 + 0.9^3 * 130 = 150.87$$

$$\text{episode 2} = 10 + 0.9 * 30 + 0.9^2 * 10 = 45.1$$

$$= (150.87 + 45.1) / 2 = 97.985$$

$$q_\pi(s^B, -1) \text{ episode 1} = 10 + 0.9 * 20 + 0.9^2 * 10 + 0.9^3 * 30 + 0.9^4 * 10 = 70.351$$