

NAME:

CSCI S-89c Deep Reinforcement Learning
Part I of Midterm

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 any of the other states. Let's denote an intention to move to state s^A by a^A , to state s^B by a^B , and to state s^C by a^C . The agent does not know transition probabilities, including the distributions of rewards.

Suppose the agent uses the following behavior policy $b(a|s)$:

$$\begin{aligned} b(a|s^A) &= \begin{cases} 0.5, & \text{if } a = a^A, \\ 0.25, & \text{if } a = a^B, \\ 0.25, & \text{if } a = a^C, \end{cases} \\ b(a|s^B) &= \begin{cases} 0.25, & \text{if } a = a^A, \\ 0.5, & \text{if } a = a^B, \\ 0.25, & \text{if } a = a^C, \end{cases} \\ b(a|s^C) &= \begin{cases} 0.25, & \text{if } a = a^A, \\ 0.25, & \text{if } a = a^B, \\ 0.5, & \text{if } a = a^C, \end{cases} \end{aligned}$$

to generate two episodes:

episode 1:

$S_0 = s^A, A_0 = a^B, R_1 = 20, S_1 = s^B, A_1 = a^C, R_2 = 10, S_2 = s^C, A_2 = a^A, R_3 = 90, S_3 = s^A, A_3 = a^C, R_4 = 30$;

episode 2:

$S_0 = s^C, A_0 = a^A, R_1 = 50, S_1 = s^C, A_1 = a^C, R_2 = 30, S_2 = s^C, A_2 = a^B, R_3 = 10, S_3 = s^B, A_3 = a^B, R_4 = 20$.

Using the Every-Visit Monte Carlo (MC) prediction algorithm for estimating $v_\pi(s)$, please estimate

(a) $v_\pi(s^A)$,

(b) $v_\pi(s^B)$,

(c) $v_\pi(s^C)$,

where the target policy is $\pi(a^C|s) = 1$ for $s \in \{s^A, s^B\}$ and $\pi(a^C|s^C) = 1$. Assume $\gamma = 0.9$.

SOLUTION: