

NAME: KARMA TARAP  
CSCI E-89c Deep Reinforcement Learning  
Part I of Assignment 6

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  $n$ -step Temporal-Difference (TD) prediction algorithm with  $\alpha = 0.1$ :

$$V_{t+n}(s) = \begin{cases} V_{t+n-1}(s) + \alpha [G_{t:(t+n)} - V_{t+n-1}(s)] , & \text{if } s = S_t, \\ V_{t+n-1}(s), & \text{if } s \neq S_t, \end{cases}$$

where  $G_{t:(t+n)} \doteq R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$ .

Assume  $\gamma = 0.9$ . If the agent observes the following sequence of states, actions, and rewards:

$S_0 = s^A, A_0 = +1, R_1 = 20, S_1 = s^B, A_1 = +1, R_2 = 30, S_2 = s^B, A_2 = +1, R_3 = 20, S_3 = s^B, A_3 = +1, R_4 = 10, S_4 = s^C, A_4 = 0, R_5 = 130, S_5 = s^C$ ,

find  $V_{t+n}(s)$  for (a)  $n = 1$  and (b)  $n = 2$ . Assume that the state-values are initialized at 0 for all  $s \in \mathcal{S}$  and fill the tables in below.

SOLUTION:

(a) 1-step TD, assuming  $V_0(s) = 0$  for all  $s \in \mathcal{S}$ :

	$t + n = 0$	$t + n = 1$	$t + n = 2$	$t + n = 3$	$t + n = 4$	$t + n = 5$
$V_{t+n}(s^A)$	0	2	2	2	2	2
$V_{t+n}(s^B)$	0	0	3	4.7	5.23	5.23
$V_{t+n}(s^C)$	0	0	0	0	0	13

(b) 2-step TD, assuming  $V_0(s) = V_1(s) = 0$  for all  $s \in \mathcal{S}$ :

	$t + n = 0$	$t + n = 1$	$t + n = 2$	$t + n = 3$	$t + n = 4$	$t + n = 5$
$V_{t+n}(s^A)$	0	0	4.7	4.7	4.7	4.7
$V_{t+n}(s^B)$	0	0	0	4.8	7.22	19.198
$V_{t+n}(s^C)$	0	0	0	0	0	0