

RL-A3

Abhishek Maity

2016005

Q1. To make the update function efficient, we need a step update

function

$$\theta_n(S_t, A_t) = \frac{1}{n} \sum_{i=1}^n G_i$$

$$= \frac{1}{n} \left( G_n + \sum_{i=1}^{n-1} G_i \right)$$

$$= \frac{1}{n} \left( G_n + (n-1) \cdot \frac{1}{(n-1)} \sum_{i=1}^{n-1} G_i \right)$$

$$= \frac{1}{n} \cdot (G_n + (n-1) \theta_{n-1})$$

$$= \cancel{\frac{G_n}{n}} + \cancel{\theta_{n-1}} + \frac{\theta_{n-1}}{n}$$

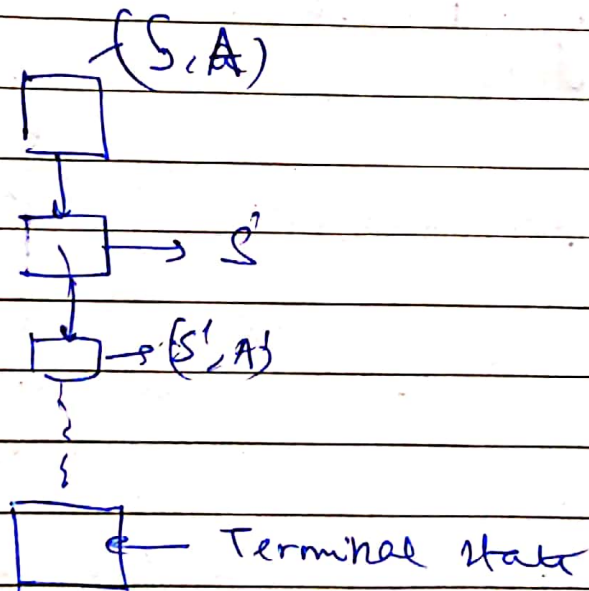
$$= \theta_{n-1}(S_t, A_t) + \frac{1}{n} (G_n - \theta_{n-1}(S_t, A_t))$$

Now we only need to maintain the count of each time considering  $(S_t, A_t)$

We don't need to keep track of the previous updates.

Pseudo code done at the end

Q2 Ans. Backup diagram for  $G_t(S, A)$ , since



Monte-carlo only follows one particular episode, there are no other branches.



Q3 Ans:

~~Not~~

Suppose all the ~~next~~ time steps in which state 's' is visited and action 'a' is taken be denoted by  $f(s, a)$ .

$$\text{Then } Q(s, a) = \frac{\sum_{t \in f(s, a)} \gamma^{t-T-1} G_t}{\sum_{t \in f(s, a)} \gamma^{t-T-1}}$$

where  $G_t$  are the returns and

$\gamma^{t-T-1}$  is the importance sampling parameter.

~~the strategy from left to right~~

Q Same:

When we ~~start~~ <sup>coming</sup> from a new building, but the part after highway remains same, then we just need to update the portion before reaching the highway as rest of the estimates remain same.

Q Same: Even if the action is greedy,  $Q$ -learning

is different from Sarsa because

in  $Q$ -learning, we change the value (update it) and then select the action in the next time step but

in Sarsa, we ~~first~~ choose the action based on old values and then update it.



## Q6: Ex. 6.3

Given that  $\alpha = 0.1$ , and  $\gamma = 1$

We know that

$$V(S_t) = V(S_t) + \alpha [R_{t+1} + V(S_{t+1}) - V(S_t)]$$

~~known~~  
We can see that all the transitions give 0 reward apart from when going to transition state.

Also, we initialised our value function with 0.5.

When we ended up in the ~~state~~  
left terminal state we got reward 1.  
But we were supposed to get

0.5

so the update looks like

$$V(A) \leftarrow V(A) + 0.1 [0 + 0 - 0.5]$$

$$:= 0.9 V(A)$$

$$:= 0.9 \times 0.5$$

$$:= 0.45$$

so the new value of <sup>left</sup> the terminal state is 0.45 and decreased by 0.05.

For all the other states

$$\alpha [R_{t+1} + V(S_{t+1}) - V(S_t)]$$

is 0.

hence no updates.



## Exercise 6.5

$\alpha'$  can be seen as the jump parameter or the step-size, if we overshoot the optimal estimation, it tries to come back hence the observed ups and downs in the graph.

## Exercise 6.4

Yes, if the  $\alpha$ 's are bigger we will get results for TD(0) method which are worse performing than ~~the~~ monte carlo. This is because,  $\alpha'$  represents the jump size, and if the jump size is large it is harder to optimize.

Q 3

Pseudo code

Input : a policy  $\pi$  to be evaluated

Initialize :

$V(s) \leftarrow \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$

$\text{Return}(s) \leftarrow$  an empty list for all  $s \in \mathcal{S}$

Loop Forever (for each episode):

Generate an episode  $\pi : s_0, a_0, \dots$

$G \leftarrow 0$

Loop for each step of episode  $t = T-1, \dots, 0$

$G \leftarrow \gamma G + R_{t+1}$

Unless  $s_t$  appears in  $s_0, s_1, \dots, s_{t-1}$ :

$V(s_t) \leftarrow V(s_t) + \frac{1}{n} (G - V(s_t))$