

PART-1

EXERCISE - 5.4

MC-ES estimates

Initialize:

$\pi(s) \in A(s)$ arbitrarily, for $\forall s \in S$

$Q(s, a) \in \mathbb{R}$ arbitrarily, for $\forall s \in S, a \in A(s)$

$N(s, a) \leftarrow 1$ for $\forall s \in S, a \in A(s)$

Loop forever (for each episode)

Choose $s_0 \in S, A_0 \in A(s_0)$ randomly
such that all pairs have probability > 0 .

Generate an episode from s_0, A_0 following

TT: $s_0, A_0, R_1, \dots, s_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

loop for each step of episode, $t = T-1, T-2, \dots, 0$

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

Unless the pair s_t, A_t appears in
 $s_0, A_0, s_1, A_1, \dots, s_t, A_{t-1}$:

$N(s_t, A_t) \leftarrow N(s_t, A_t) + 1$

$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \frac{1}{N(s_t, A_t)} \times$

$(G_t - Q(s_t, A_t))$

Basically Equation:-

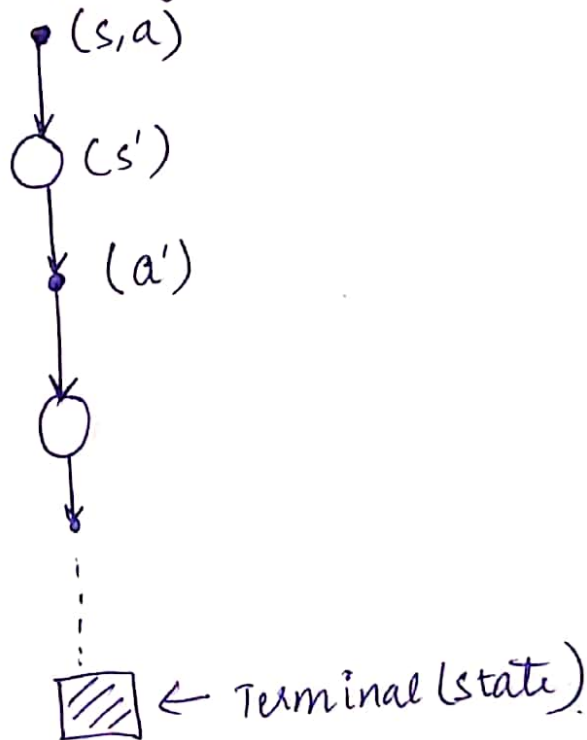
$Q(s, a) \leftarrow Q(s, a) + \frac{1}{N(s, a)} [G - Q(s, a)]$

PART-2

Ex 5.3.

Monte-Carlo updates the state-action pairs values only after generating the entire episode.

Hence, the Back-up diagram for this would be:-



where ~~we~~ we want to update root node (s, a) which is a state-action pair & followed by all the transitions of the states & actions appear in an episode till termination which contribute to the update.

PART-3

EX 5.6

According to Monte-Carlo, importance sampling ratio is given by:-

$$J_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$$

& this ratio transforms the returns to have right expected value as:-

$$V_{\pi}(s) = E[J_{t:T-1} G_t | s_t = s]$$

Now, when we have ^{stat.} action-value function

$$\text{then } q_{\pi}(s, a) = E[J_{t:T-1} G_t | s_t = s, A_t = a]$$

$$Q(s, A) = \frac{\sum_{t \in T(s, a)} J_{t:T(t)-1} G_t}{\sum_{t \in T(s, a)} J_{t:T(t)-1}}$$

Part-5

EX 6.2

⇒ The ^{new} scenario ^{has} shifts to the new building for work. In this case too, TD learning is expected to be much better than MC learning, ~~indeed~~ since there is only a change in the initial route. Once we enter highway, many of the states after highway are common to reach home. Therefore, the value function estimates for these common states should be very close to the values that were there in the original problem/building.

This would happen same⁹ in our original learning task if our initial guess at the value function is very close to that of the true value function.

Part 6.

Ex 6.3.

Given the equation for TD(0) update for a ^{state} value function is :-

$$V(s_t) \leftarrow V(s_t) + 0.1 [R_{t+1} + V(s_{t+1}) - V(s_t)]$$

$V(A)$ update after constant first episode $V(s_t)$ update will be:-

$$V(A) \leftarrow V(A) + 0.1 (0 + 0 - V(A))$$

$$= 0.9 V(A) = 0.45$$

~~which~~ which is also there in graph 2 happens only when the left action is chosen from state A.

So, the change is by $0.5 - 0.45 = 0.05$.

Part 6.

Ex 6.3.

Given the equation for TD(0) update for a state value function is :-

$$V(s_t) \leftarrow V(s_t) + 0.1 [R_{t+1} + V(s_{t+1}) - V(s_t)]$$

$V(A)$ update after constant first episode $V(s_t)$ update will be:-

$$V(A) \leftarrow V(A) + 0.1 (0 + 0 - V(A))$$

$$= 0.9 V(A) = 0.45$$

~~which~~ which is also there in graph 2 happens only when the left action is chosen from state A.

So, the change is by $0.5 - 0.45 = 0.05$.

Exercise 6.4.

No, I don't think that the wide range of alphas would affect ~~the~~ concluding about which algorithm is better. As can be observed, that for different values of α , the plots were deflecting a bit for higher alphas & ~~same~~ same goes for MC. but overall everytime TD performs better than MC.

It cannot be fixed i.e. value of alpha since it will depend on the task to be learnt. And also, it is seen that TD performs better than MC for these tasks like there where one step Reward is known & ~~don't need the entire episode to be generated.~~

Exercise 6.5

Since the equation for TD is :-

$$V(s) = V(s) + \alpha [R + \gamma V(s') - V(s)]$$

Where $R + \gamma V(s') - V(s)$ is the error. So, the large values of alpha help converging faster but this actually impacts the random step taken at each time step during the algorithm. ~~where~~. ~~the algorithm~~ and values become highly sensitive to these steps with large alpha values. Also, with large alpha, the error part during updation will be higher which initially helps converging but later will result in more value state function.

Part 8.

Exercise 6.12

Since action selection is greedy, then Q-learning will no more be off-policy method, since initially in Q-learning, we select actions on the basis of ϵ -greedy & value updation is done using greedy policy. Hence the same policy is used for selecting action as well as value updation. Therefore it converts to on-policy method.

However, the action selection & updates for this changed Q-learning will not be same as SARSA because although both are now on-policy methods but SARSA selects ^{action} updates ~~values~~ ^{weights} according to ϵ -greedy policy than greedy policy used ~~here~~ ^{changed} in Q-learning. Unless SARSA also uses the same greedy policy.