

## R2 ASSIGNMENT - 3

### Question 1

The code  $\left[ \begin{array}{l} \text{Append } G \text{ to List } R. \\ Q(s_t, A_t) \leftarrow \text{avg}(R) \end{array} \right]$  is highly inefficient.

This stores all returns in a list. And the list needs to be stored for every iteration and the Action-Value function  $Q(s_t, A_t)$  is updated by taking an average of all returns stored in the list. We update the update rule for  $Q(s_t, A_t)$  as follows.

$$Q_n(s_t, A_t) = \frac{1}{n} \sum_{i=1}^n G_i(s_t, A_t) \quad \text{--- (i)}$$

= We split it as follows.

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

$$= \frac{1}{n} \left[ G_n(s_t, A_t) + \frac{(n-1)}{(n-1)} \sum_{i=1}^{n-1} G_i(s_t, A_t) \right]$$

= Using (i)

$$= \frac{1}{n} \left[ G_n(s_t, A_t) + (n-1) Q_{n-1}(s_t, A_t) \right]$$

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

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

The update rule is finally as follows

$$Q_n(S_t, A_t) = Q_{n-1}(S_t, A_t) + \frac{1}{n} [G_n(S_t, A_t) - Q_{n-1}(S_t, A_t)]$$

The new updated and equivalent pseudo-code is:-

~~According~~ As the only change is in the update rule which is equivalent to the one used before we can claim that the code itself is equivalent.

Also

we need to calculate,

$n \rightarrow$  the count of occurrences for all  $(s, a)$  (Every visit  $n \leftarrow 1$ )

Code:

~~Initialization~~ Steps

$Q(s, a), \pi(s) \leftarrow$  arbitrarily  $\forall A \in \mathcal{A}(s) \forall s$ .

~~Initialize~~  $n \leftarrow 0$ .

while (True)

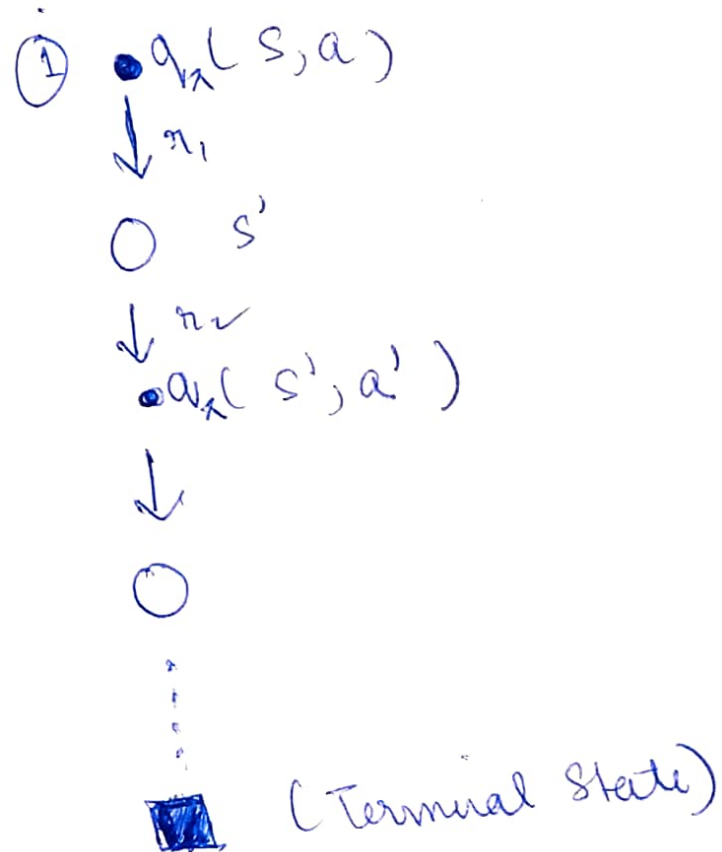
- Select  $s_0, A_0 \in \mathcal{S}, \mathcal{A}$  randomly with  $P(E) > 0$ .
- Generate episodes
- Return  $\leftarrow 0$

for every timestep in Episode

- $\text{Return} \leftarrow \text{gamma} \times \text{Return} + R_{t+1}$
- ~~unless~~  $(s, a)$  not in (episode pairs):
  - ~~increment~~  $n = n + 1$
  - $Q = Q + \frac{1}{n} [G - Q]$
  - $\pi(s) \leftarrow \text{argmax } Q$

## Question 2

The backup diagram for  $Q \rightarrow$  is ..



The ① i.e. the topmost node shows the initial (State-action) pair. The rest of the nodes represent all intermediate nodes

Index

● → (State-action) Nodes

○ → State nodes

### Question 3

Let  $J(s, a)$  represent the set of all state action pair  $(s, a)$  which have been visited by an agent.

Given the equation (5.6)

$$V(s) = \frac{\sum_{t \in J(s)} e_{t:T(t)-1} G_t}{\sum_{t \in J(s)} e_{t:T(t)-1}}$$

We rewrite  $J(s)$  as follows

$$J(s, a) = \frac{\sum_{t \in J(s, a)} e_{t+1:T(t)-1} G_t}{\sum_{t \in J(s, a)} e_{t+1:T(t)-1}}$$

where  $G_t$  is the return at steps  $t$  to  $T(t)$

&  $e_{t:T(t)-1}$  &  $e_{t+1:T(t)-1}$  are the imp. sampling ratios.



## Question 5

The mentioned example scenario in the text would show significantly better performance of a TD Based approach, as many of the initial states in the problem remain unchanged, i.e. For the route to home while the home itself is changed the highway remains the same.

→ TD learning would be able to make quicker updates and as the number of ~~optimal~~ common states is comparatively higher we'll reach an optimal solution in quicker time.

→ MC. methods on the other hand would have to wait for an whole episode to finish as they are not bootstrapped and would end up being inefficient as compared to TD methods.

~~Also,~~

## Question 8

Even if the action selection policy is greedy, Q-learning is not exactly the same as SARSA learning. Q-learning algorithms first update the values and then the next action is picked on the basis of the updated values.

On the other hand Q-learning algorithms first select an action and then update the Q values for every (state-action) pair.

x A v

Q6 Exercise accompanying the generated graphs in the code.

Graph 1 : Estimated Values Graph

Graph 2 : Empirical RMS Error.

6.3

We start from the middle i.e. state 'C'. We also know that transition to any state has a reward '0', except +1 for terminating right.

We follow the following update rule for updating our state value function

$$V(s) = V(s) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s)]$$

According to our initialization

$$V(B) = 0.5, V(D) = 0.5$$

$R=0$  for  $B \neq D$ .

$$\begin{aligned} \hookrightarrow V(C) &= 0.5 + 0.1 (0 + 0.5 - \cancel{0.5} \rightarrow 0) \\ &= 0.5 \quad (\text{remains unchanged}) \end{aligned}$$

Now we can either go to E or A (in the direction of). Since  $V(E)$  is unchanged the agent must have gone towards 'A', as a result of which rest all values are unchanged.

Also,

As left terminal state's value function is 0 by definition.

$$\begin{aligned} V(A) &= 0.5 + 0.1 (0.3 + 0 - 0.5) \\ &= 0.45 \end{aligned}$$

Change = 0.05 i.e. a 10% deduction.

6.4

The conclusions about the performance of algorithms would somewhat depend on the range of alpha that we choose. TD Learning would always perform better than the MC methods given  $\alpha$  is reasonably small.



However taking larger values of  $\alpha$  might affect the TD approaches as they <sup>may</sup> result in more sudden and large jumps / oscillations which might stop the algorithm from converging.

6.5

Larger values of ' $\alpha$ ' can cause such oscillations, in the error of TD learning algorithms. As the algorithm tries to reach an Optimal State, a larger value of ' $\alpha$ ' may stop the algorithm from converging completely.

This is true for all larger values of ' $\alpha$ ' independent of ~~the~~ how the value function is initialised.