

### Assignment-3

#### Monte Carlo and TD

Q1 Rather than calculating the mean of return repeatedly, we could follow an iterative approach. <sup>update</sup> Here we wish to calculate the state value  $V$  or action-value function in each episode depending upon the state, action and type of visit approach.

The updated <sup>sub-part</sup> pseudo-code is as follows:

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

updated step:

$$G_t = \gamma G_{t+1} + R_{t+1}$$

$$Q_{\pi}(S_t, A_t) = (1-\alpha) Q_{\pi}(S_t, A_t) + \alpha G_t$$

$$\pi(S_t) = \arg\max_{a \in A(S_t)} Q_{\pi}(S_t, a)$$

Note: Here  $(1-\alpha)$  and  $\alpha$  weights are used in case of Exponentially weighted average.

For Sample mean average, we can replace  $\alpha$  by  $\frac{1}{n}$  where  $n$  = the

$n^{\text{th}}$  instance for state  $S_t$ , action  $A_t$  among all episodes.

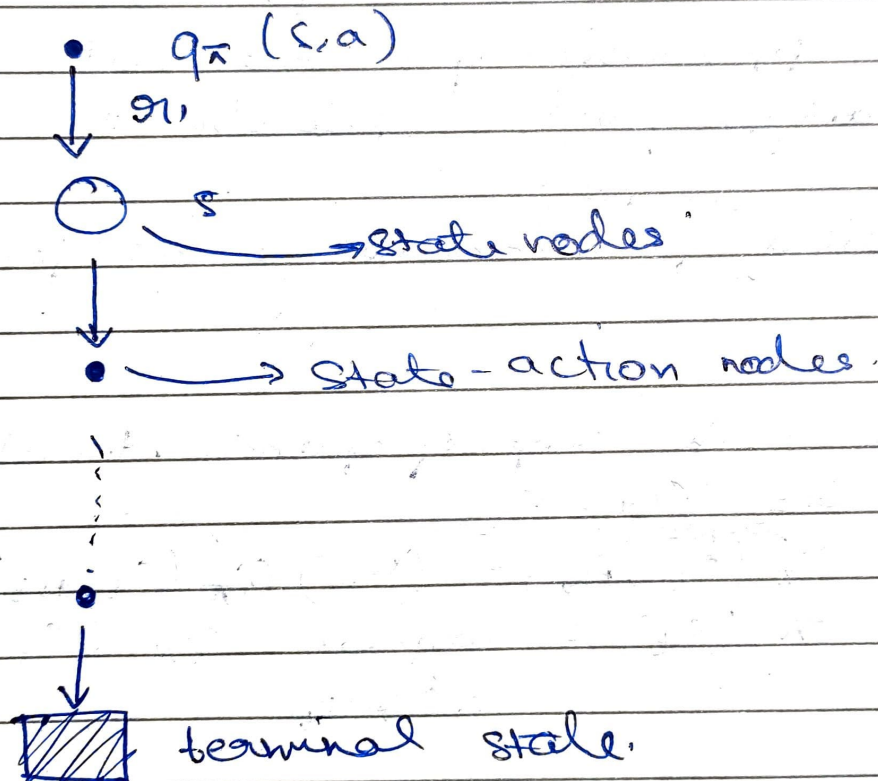
# The intuition behind this is that:

$$Q_{\pi}(S_t, A_t) = Q_{\pi}(S_t, A_t) + \alpha (G_t - Q_{\pi}(S_t, A_t))$$

So ' $\alpha$ ' is the weighting factor to improve the current action-state value function <sup>for the error</sup> eg.

Q2 we initially choose a random state action pair which is then followed by reward and proceeding towards the next state.

Root node is the top most node which is an action - state pairs.



Q3. Equation analogous to action-values  $Q(s, a)$  would be :

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

Here  $Q(s, a)$  is the state-action pair and  $T(s, a)$  is the set of all ~~max~~ time steps when  $S_t = s$  and  $A_t = a$ .



Q5 TD learning would be very efficient in the case mentioned ~~as~~ because:

- Since only a small part of the route has changed, thus keeping most of the states similar in the new as well as old problem. ~~Also~~
- Since the device has prior experience on the initial route this state-value func. of TD for old problem would have already converged to optimal values. Also TD bootstraps the original state values <sup>as obtained</sup>.
- We can update the state values for other / prev. states during the episode unlike Monte Carlo, wherein we ~~need~~ need to wait for the episode to end.

Even if the initial state values are close to the true values, the same could be seen and they would converge faster.

Q6

Left graph: Estimated values Graph

Right graph: Empirical Rms error, avg over all states

6.3

We have the following hyperparameters:

-  $\alpha = 0.1$ 

- episodes = 100

- Reward for each state for transition to each state except terminating right state = 0.

- Reward for transition to terminating right = 1

→ Updates in state-value func are as follows.

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

→ Initially,  $V = [0, 0.5, 0.5, 0.5, 0.5, 0.5, 0]$ 

→ we start at state: C

Thus

$$V(C) = 0.5 + 0.1 (0 + 0.5 - 0.5)$$

→ both right and left action.

$$= 0.5$$

here  $V(C)$  is unchanged

Now we can either go left or go right

~~Now  $V(C)$  is also unchanged.~~In the 1<sup>st</sup> ep, when we go to A, i.e. left terminal state, reward is 0.



Also as value function for left terminal state is 0 thus,

~~$$V(A) = 0.5 + 0.5$$~~

$$V(A) = V(A) + 0.1 (0 + 0 - V(A))$$

$$V(A) = 0.5 + 0.1(-0.5)$$

$$\Rightarrow 0.45$$

Thus, we can conclude that only only state value of A gets changed.

The decrease is 0.05 i.e. there is a 10% deduction in the state value of A.

6.4

The conclusions about the performance of each algorithm are dependent on the <sup>range</sup> value of  $\alpha$  that we look to observe in.

TD(0) would perform poorly for larger values of  $\alpha$  as then larger jumps/oscillations would start to occur and it would be hard to get a better result ~~through~~ as convergence would become difficult in this case.

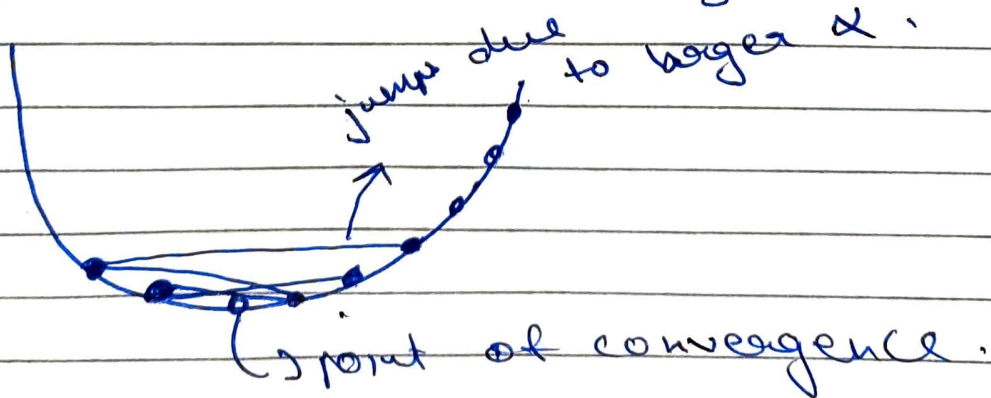
whereas TD(0) would always perform better for smaller learning rate ( $\alpha$ ) relative to Monte-Carlo algorithm.

Also a specific alpha could be found where these algo might perform relatively well.

### 6.5

As ' $\alpha$ ' increases, oscillation might start to occur in the case of TD learning algo. These ~~cases~~ oscillations or jumps would prevent the algorithm from converging.

This phenomena is true for all larger values of  $\alpha$  and is not a function of how the value function is initialized.





Q3 Even if the action selection is greedy, Q-learning and SARSA would be different as. In Q-learning,  $Q(s,a)$  is updated using greedy approach i.e by finding the best action  $a'$  for which target is max.  $(\arg\max (r + \gamma Q(s',a')))$ . here  $Q$  values are updated and then the next action is picked whereas In SARSA, we pick the next action from  $a$  given the  $Q$  values according to the next state. Further, the  $Q$ -values can now be updated according to the next-action.

Thus, due to these reasons Q-learning is different from TD learning.