

RL Assignment 3

Question 1 -

Pseudo code for monte Carlo ES -

Prerequisites :-

T states ranging from 0 to (T-1)

An episode is generated from s_0, A_0 following -

$s_0, A_0, R, s_1, A_1, \dots, s_T, A_T, R_T,$
 Returns $(s_t, A_t) \leftarrow \emptyset$ (Empty list)

Pseudocode :-

For each episode i.e., from $t = T-1, T-2, \dots, 0$.

$$G_t \leftarrow r G_t + R_{t+1}$$

returns

Append G to Returns (s_t, A_t)

$Q(s_t, A_t) \leftarrow \text{average}(\text{Returns}(s_t, A_t))$

$\pi(a_t) \leftarrow \underset{a}{\text{argmax}} Q(s_t, a)$

Since, $\frac{1}{T} \sum_{t=0}^{T-1} G_t$

So, rewritten as -

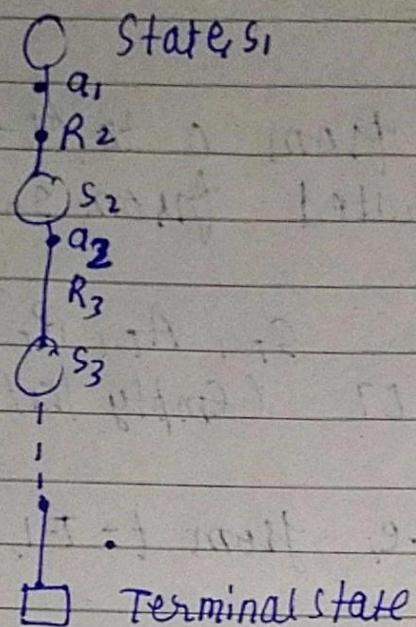
$$\frac{1}{T} \left[G_{T-1} + \sum_{t=0}^{T-2} G_t \right] = \frac{1}{T} \left[G_{T-1} + \frac{1}{T-1} \sum_{t=0}^{T-2} G_t \right]$$

$$= \frac{1}{T} \left[G_{T-1} + (T-1) Q(s_{T-1}, A_{T-1}) \right]$$

Recursively, $Q(s_t, A_t) = Q(s_{t-1}, A_{t-1}) + \frac{1}{T} [G_t - Q(s_{t-1}, A_{t-1})]$

Question 2 -

Backup Diagram for Monte-Carlo -



It includes entire episode, only one choice can be made at each state. It doesn't estimates on the basis of other estimates, & estimates are independent.

equation given as -

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

for $Q(s, a)$, it will be -

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

where, $P_t : T(t)-1$ is important sampling ratio
 G_t is expected return at time t .

Question 3

Equation given as -

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

for $Q(s, a)$, it will be -

$$Q(s, a) = \frac{\sum_{t \in J(s, a)} P_t : T(t)-1 q_t}{\sum_{t \in J(s, a)} P_t : T(t)-1}$$

where, $P_t : T(t)-1$ is important sampling ratio
 q_t is expected return at time t .

Question 5:-

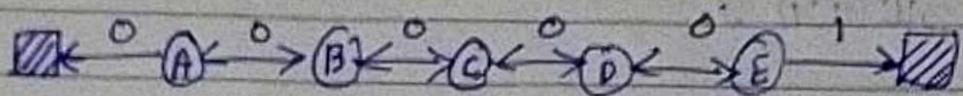
TD learning would be better than MC learning in the case where we are moved to a new building since it is just a change in our initial route & some of the state encountered during general episode will be same.

For example:- On our drive home many of states are the same once we enter the highway & value function estimates for these states obtained when we worked in original building should be very close to what we will compute starting from our new building.

Starting with a very good initial guess should result in faster convergence.

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

Question 6.3



Since, $\alpha = 0.1$, & $r = 1$.

Then,

TD expression for state value function will be -

$$\begin{aligned} V(S_t) &= V(S_t) + \alpha (r_{t+1} + V(S_{t+1}) - V(S_t)) \\ &= V(S_t) + 0.1 \cdot (1 + V(S_{t+1}) - V(S_t)) \end{aligned}$$

Initial value function begin with constant value S_0 , in first update, there is no change in value function.

If we terminate on left the reward is zero, so state A will be updated.

$$\begin{aligned} V(A) &= V(A) + 0.1 (r_{t+1} + V(A_{t+1}) - V(A)) \\ &= V(A) + 0.1 (0 + 0 - V(A)) \\ &= 0.5 + 0.1 (-V(A)) \\ &= 0.45. \end{aligned}$$

Thus, we go to left of random walk, then we decrease the state value function by 0.05.

Question 6.4

The Right figure (graph) in Example 6.2 is dependent on value of α .

Since, when α will be small, the convergence of Monte Carlo & Temporal Difference is satisfied. Small value of α is better than large value in long run, And so, learning will be more & finally reach to lowest error.

Therefore, we cannot conclude which algorithm will perform better if wider range of α is used.

Question 6.5

It is given that α is large enough, so, large value of α causes more change in state-value function for each timestamp.

Thus, TD(0) is heavily dependent on specific returns. So, initially graph of RMS error goes down because learning rate is high & then up.

At smaller value of α , learning will take longer & it is less sensitive to specific random step as compare to large value of α . While on other hand, if initialization of initial value to state generates linear relations about updates during transition, there may chance it over estimates result around terminal state.

Question 6

No, it will not take same selection because Q-learning will update Q-function value first. As Q-learning is off-policy algorithm.

In Q-learning, the next action is to perform to selected in next iteration derived from q-function while SARSA is on-policy algorithm, it chooses a' , s' & then update Q-value.

Q-learning Algorithm :-

Initialise Q-table with ones
for each episode

Initialise s

for each step of Episode

choose a from s using greedy policy

Take action a , observe r, s' .

$$Q[s, a] \leftarrow Q[s, a] + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

$s \leftarrow s'$

until s is terminal

SARSA Algorithm :-

Initialise Q-table with ones
for each episode

Initialise state, s

choose action from state using policy derived from Q (greedy approach)

For each step of episode

Take action, observe r, s' .

choose a' from s' & update Q-table

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma(Q(s', a') - Q(s, a)))$$

repeat until $s \leftarrow s'$ and $r = 0$
 $a \leftarrow a'$ taking max a'
until s is terminal.