

RL

Homework-3

Ans 1. Monte Carlo Exploring starts, Modified:

Initialize:

1. $\pi(s) \in A(s)$ (arbitrarily), for all $s \in S$
2. $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in S$, $a \in A(s)$
3. $N(s, a) = 0 \quad \forall s \in S, a \in A(s)$:

Loop forever (for each episode):

4. Choose $s_0 \in S$, $A_0 \in A(s_0)$ randomly such that all pairs have probability > 0
5. Generate an episode from s_0, A_0 following $\pi: s_0, A_0, R_1, \dots, s_{T-1}, A_{T-1}, R_{T-1}$
6. $G \leftarrow 0$
7. Loop for each step of episode, $t = T-1, T-2, \dots, 0$:
8. $G \leftarrow \gamma G + R_t$
9. $N(s_t, A_t) = N(s_t, A_t) + 1$
 unless the pair s_t, A_t appears in
10. $Q(s_t, A_t) = Q(s_t, A_t) + \frac{1}{N(s_t, A_t)} (G - Q(s_t, A_t))$
11. $\pi(s_t) \leftarrow \arg\max_a Q(s_t, a)$

(State, action pairs)

line 3 initializes the count of $N(s, a) \quad \forall s \in S, a \in A$

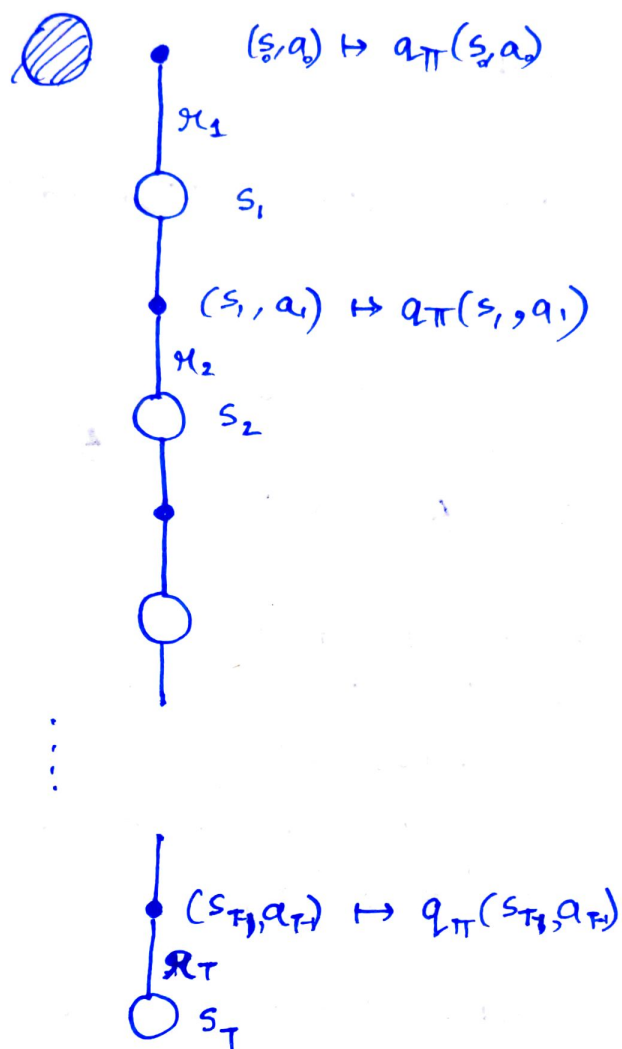
line 9 increments the count of s_t, A_t by 1

line 10 updates the Q value by incremental mean approach



Ans 2

Backup diagram for Monte Carlo Estimation of q_π .



~~Ans 6.1~~ Ans 6.3 In the first episode, the terminal state is the left terminal state ~~not~~ because of which the value of state ~~of~~ ^{only} decreases by $\alpha * (0.5)$ which is 0.45. The estimate of state A changed because the estimated value of terminal state is 0 while of others is 0.5 and also the reward is 0 in this episode.

Ans 6.4