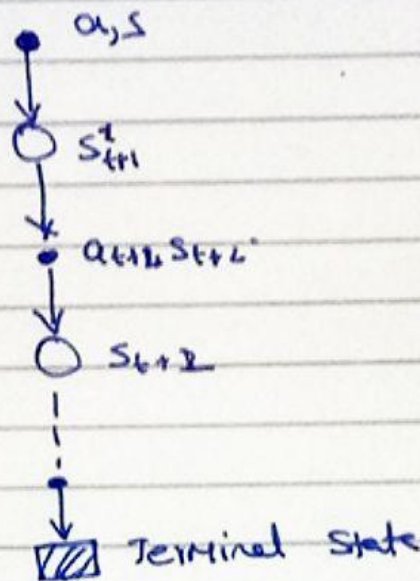


Q2. Monte Carlo Estimation of $q_{\pi}(a, s)$



Q1. The idea to make MC more efficient is by incrementally $V(s)$ after each episode $s_1, a_1, r_1, \dots, s_T$ instead of storing the entire history requires only keeping track of previous $V(s)$ value.

For Incremental,

$$V(s_{t+1}) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t))$$

Pseudo code:

Here, G_t can also be computed incrementally,
 $G_{t+1} = R_{t+1} + \gamma G_t$

Similarly, for $Q(a, s)$, we have

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G_t - Q(s_t, a_t))$$

Pseudo-code:

Initialize

$$\begin{array}{l} \pi(s) \in A(s) \\ Q(s, a) \geq 0 \end{array} \quad \left. \vphantom{\begin{array}{l} \pi(s) \in A(s) \\ Q(s, a) \geq 0 \end{array}} \right\} \text{Arbitrarily}$$

Loop for each episode:

(a) Choose s_0, a_0 such that all pairs have prob. > 0

(b) Generate episode from s_0, a_0 : s_0, a_0, \dots, L

(c) $G \leftarrow 0$

Loop for each $t = T-1, T-2, \dots, 0$

(d) $G \leftarrow \gamma G + R_{t+1}$

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

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

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (G_t - Q(s_t, a_t))$$

(e) $\pi(s_t) \leftarrow \arg \max_a Q(s_t, a)$

Previously,

Execute (e), (f), (g) unless s_t, a_t appears in $s_0, a_0, \dots, s_{t-1}, a_{t-1}$,

Using the new pseudocode, we can reduce one iteration of the (s, a) pairs in the episode.

Q3
$$P_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)}$$

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

Just like $E[P_{t:T-1} G_t | S_t = s] = V_\pi(s)$

by definition we can write,

$$E[P_{t:T-1} G_t | S_t = s, A_t = a] = q_\pi(s, a)$$

$$P_{t:T-1} G_t = P_{t:T-1} (R_{t+1} + \gamma V_{t+2} + \gamma^2 V_{t+3} + \dots + \gamma^{T-t} R_{t+T})$$

$$q_\pi(s, a) = E_\pi [P_{t:T-1} G_t | S_t = s, A_t = a]$$

$$= E_\pi [P_t R_{t+1} + P_{t:T-1} G_{t+1} | S_t = s, A_t = a]$$

$$= P_t R(s, a) + P_t \gamma V_\pi(s)$$

$$= P_t R(s, a) + P_t \gamma \frac{\sum P_{t+1} G_{t+1}}{\sum P_{t+1}}$$

Q5. TD Learning Vs MC Learning

Case:

Consider I have a lot of experience driving home from work. Then I move to a new building & new parking lot but still use the same highway.

In this case, TD learning would be better than MC learning where we move to new building as this is just a change in our initial route and some of the states encountered during the rest of the episode will be same.

Example, on the drive home, many of the states are the same since we use the same highway & the value estimates for these states is very close to what we will compute for new building.

This is true for case where initial guess of value function is close to true function in which case convergence is faster.

Q8. Exercise 6.12

On choosing a greedy policy, the performance of SARSA is a little different from that of Q -learning.

In Q -learning: $Q(s_{t+1}, a_{t+1}) = \max_a Q(s_{t+1}, a)$

$$\text{SARSA: } Q(s_{t+1}, a_{t+1}) = \epsilon \cdot \text{mean } Q(s_{t+1}, a) + (1-\epsilon) \max_a Q(s_{t+1}, a)$$

Since, Q -learning always chooses the optimal solution, it will always have the best (shortest) path to the goal.

The rate of convergence however can vary.

For instance, in the cliffhanger problem, SARSA will choose the sub-optimal path while Q -learning will choose the optimal path. However, the cumulative reward (performance) of SARSA might be better over the short run compared to Q -learning.

Over the long run, however, Q -learning performs better.

→ Exercise 6.3, 6.4, 6.5

Q6. Assuming $\alpha = 0.1$

For TD(0),

$$V(S_t) \leftarrow V(S_t) + 0.1(r_{t+1} + V(S_{t+1}) - V(S_t))$$

Computing $V(A)$ starting from terminal state,

$$\begin{aligned} V(A) &\leftarrow V(A) + 0.1(0 + 0 - V(A)) = 0.9 V(A) \\ &= 0.45 \end{aligned}$$

\uparrow
 $V(S_{t+1})$
↳ Terminal State

On the first episode, the random walk results in 'left' action & we decreased the state value function by 0.05.

Growth of RMS Error.

large values of α indicates a larger $V(s)$ update at each time step. A larger value allows to converge to optima quickly compared to smaller value of α .

The up & down movement of RMS is due to the stochastic nature of updates made at each time step of the random walk.