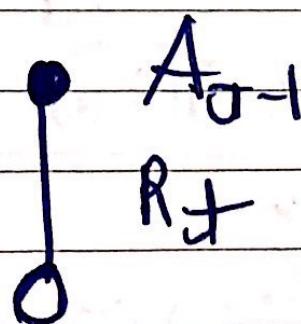
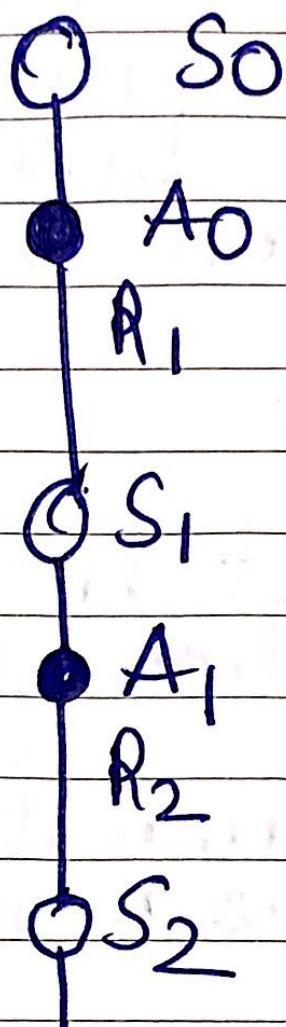


Date

Question 2: The backup diagram for Monte Carlo will show all states till the end of the episode.



Question 3

Following notation given in the
TB,

$T(s,a)$ \Rightarrow set of all time steps
in which we arrived
at state s and
picked action.

~~$T(t)$~~ Denote

$T(t)$ \Rightarrow first time of terminal
following time t .

$g_t \Rightarrow$ return after t with
through $T(t)$.

So do

$$Q(s,a) = \sum_{t \in T(s,a)} p_t \cdot T(t) + g_t$$

$$\sum_{t \in T(s,a)} p_t \cdot T(t) -$$

Question 5:

The given diagram states that there is only some part of the path that is charged. And since we have a lot of experience driving in the path that is constant, we just have to learn from parking lot to highway.

Td algorithm will take advantage of this, and hence by using state values of the constant path, and here be better off initially.

In the original scenario, if our initial estimates are close to the true values, then also TD update would be steeper on average than Monte Carlo.

Date / /

Question 8:

~~Question~~

No, Q learning is not same as Sarsa.

At its best, Q-learning is an off-policy and Sarsa is on-policy.

Sarsa chooses a' and s' and then updates Q-functions whereas Q-learning first updates Q-function, and next action is selected after. These may not be the same.

Date

Question 1:

What are state-action hairs hardness as follows:

$$Q_{n+1}(s, a) = \frac{1}{n} \left(\sum_{j=1}^n q_j \right)$$

$$= \frac{1}{n} \left(\sum_{j=1}^{n-1} q_j + q_n \right)$$

$$= \frac{1}{n} \left((n-1) \times \left[\frac{1}{n-1} \sum_{j=1}^{n-1} q_j \right] + q_n \right)$$

$$Q_{n+1}(s, a) = \frac{1}{n} \left((n-1) \cdot Q_n(s, a) + q_n \right)$$

If we use this what we do,
we just need to store the
previous value for each
state action pair instead
of n reward.



Date

The Q-function code will be:

- initializing ~~array(s,a)~~. array(s,a)

Look For next (each episode):

choose $s_0 \in S$, $a_0 \in A(s_0)$

randomly such that all
this have non zero prob.

Generate an episode starting
at s_0, a_0 and following

71.

$$q = 0$$

Look for each step, $t = T-1, T-2, 0$:

$$q = \sqrt{q + R_{t+1}}$$

$$\text{array}(s,a)_{T=1}$$

unless final s, a address
previou.

$$Q(s,a) = Q(s,a) + \frac{1}{\text{array}(s,a)} \times (q - Q(s,a))$$

$$A(s) \leftarrow \arg \max_a Q(s,a)$$



Date

We have only changed the update rule in the code and rest everything remains the same. Here both codes are evaluated.

Question 6:

6.3

Under the given conditions, update rule 1

$$w(S_t) = w(S_t) + 0.1 [R_{t+1} + u(S_{t+1}) - b(S_t)]$$

We initial $w(S) = 0.5 \quad \forall s \in S$,

Any transition from $S_t \rightarrow S_{t+1}$

does not result in an update because,

$$u(S_t) = w(S_t) + 0.1 [0 + 0.5 - 0.5]$$

However, in the first step
we went from state A to
left terminal state value is 0

$$\begin{aligned} v(A) &= v(A) + 0.1 [0 + 0 - v(A)] \\ &= 0.9 v(A) \end{aligned}$$

$$v(A) = 0.45.$$

Hence, only $v(A)$ got updated.

6.4

Yes, if we use large values
of α , then TD will perform
worse than monte carlo.

Because large values of α
will lead to higher
sensitivity towards a
receiving rewards in a
time step. This might result
in estimates that are way
off.

6.5

If we increase the value of α greatly, the update is much larger at each timestep.

This makes the T_f algorithm much more sensitive to the rewards received at that particular step.

What we notice in the graph is because of the randomly selected step in the random walk initially. The disturbance are because of this high sensitivity. If we use a lower value of α , these disturbances are reduced.

Forget the mistake, remember the lesson.