

Reinforcement Learning

Homework 3

1/1

Quest 8 Exercise 6.12

Q learning is not same as SARSA if the action selection is greedy. This is because Q learning is off policy while SARSA is on policy. In Q learning the update equation is:

$$Q(s,a) = Q(s,a) + \alpha [R + \max_{a'} Q(s',a') - Q(s,a)]$$

here, a' is found using ϵ greedy (or any other behaviour policy)

s' is then observed from s,a . s is updated to s' but

a is not updated to greedy a' selection. Although we take the greedy action in each update, our underlying policy is still the same. On the other hand in SARSA, update eqn is

$$Q(s,a) = Q(s,a) + \alpha [R + Q(s',a') - Q(s,a)]$$

where $a' \rightarrow \epsilon$ greedy ^(here greedy) policy. Hence, both s and a are updated to s' and a' . \therefore the Q values used for updation are different \therefore the α are different.

Quest 1. For calculating the average value, we do not need to store all the returns for a given s,a and recompute mean every time. The mean can be easily calculated given the current mean and number of times the (s,a) has been visited.

Say current mean $\bar{x} = \frac{v_1 + v_2 + \dots + v_n}{n}$ (1)

new mean $\bar{x}' = \frac{v_1 + v_2 + \dots + v_n + v_{n+1}}{n+1}$ (2)

From (1), $v_1 + v_2 + \dots + v_n = \bar{x} \cdot n$, put into (2),

$$\bar{x}' = \frac{\bar{x} \cdot n + v_{n+1}}{n+1}$$

\therefore the line $Q(s_t, A_t) \leftarrow \text{Average (Returns } (s_t, A_t))$
can be modified as $Q(s_t, A_t) \leftarrow \frac{Q(s_t, A_t) \cdot (\text{count}(s, A) - 1) + G}{\text{count}(s, A)}$

Pseudocode:

Initialize

$$\pi(s) \in A(s) \quad \forall s \in S$$

$$Q(s, a) \in \mathbb{R} \quad \forall s \in S, a \in A(s).$$

$$\text{count}(s, a) = 0. \quad \forall s \in S, a \in A(s).$$

loop forever:

Choose $s_0 \in S, a_0 \in A(s_0)$ randomly

Generate episode s_0, a_0 following π ?

$$G \leftarrow 0.$$

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

get s_t, a_t of episode

$$\text{count}(s_t, a_t) += 1$$

$$G \leftarrow \gamma G + R_{t+1}$$

unless s_t, a_t appears in $s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1}$:

$$Q(s_t, a_t) = (Q(s_t, a_t) \times (\text{count}(s_t, a_t) - 1) + G) / \text{count}(s_t, a_t)$$

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

Ques 3

Exercise 5-6

$$V(s) = \sum_{t \in \tau(s)} \gamma^{t: \tau(s)-1} G_t \quad \text{in terms of } Q(s, a).$$

$$\sum_{t \in \tau(s)} \gamma^{t: \tau(s)-1}$$

$\tau(s)$ determines the staying time / first time a ~~part~~ state s was visited in the sequence of states.

Now we need $\tau(s, a)$, for first visit of tuple s, a

$$Q(s, a) = \sum_{t \in \tau(s, a)} \gamma^{t: \tau(s, a)-1} G_t$$

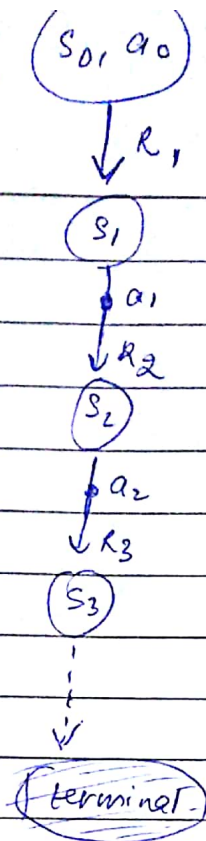
→ modified expression.

$$\sum_{t \in \tau(s, a)} \gamma^{t: \tau(s, a)-1}$$

Ques 2

Backup diagram for MC estimation of q_π

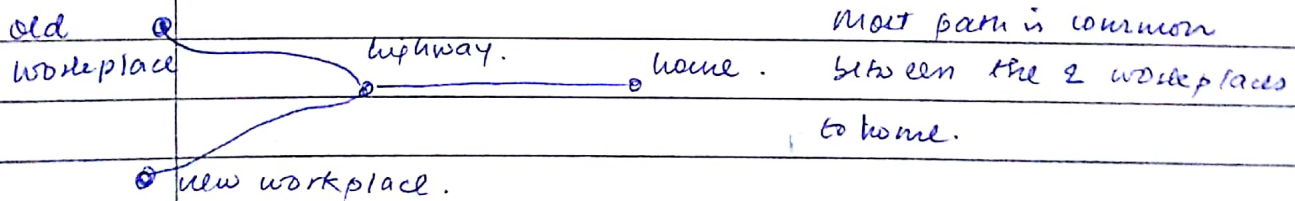
$Q(s, a)$ → start with s, a tuple



Ques 5

Yes a same thing happens for the original scenario.

Next, if we move to a new work place, the situation is somewhat like:



Since TD updates values at each step and does not need to wait for the episode to end, there will not be a large change in values of states common to both paths. As a result, optimal values can be found faster in case of TD.