

## Assignment 3,

Answer 1,

Pseudo code

Initialise

$$\left. \begin{array}{l} \pi(s) \in A(s) \text{ , for all } s \in S \\ Q(s,a) \in R \text{ for all } s \in S, a \in A(s) \end{array} \right\} \text{arbitrarily}$$

$$N(s) \leftarrow 0 \quad \forall s$$

Loop infinite (for each episode)

choose  $s_0 \in S$  ,  $A_0 \in A(s_0)$  randomly such that all pairs have probability  $> 0$

Generate an episode from  $s_0, A_0$  following policy  $\pi$ ,

$$s_0, A_0, R_1, \dots, s_{T-1}, A_{T-1}, R_T$$

$$Q \leftarrow 0$$

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

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

$\gamma \leftarrow$  discount factor.

unless the pair  $s_t, A_t$  appears in  $s_0, A_0, s_1, A_1, \dots, s_{t-1}, A_{t-1}$

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

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \frac{1}{N(s_t)} [Q - Q(s_t, A_t)]$$

$$\pi(s_t) \leftarrow \arg\max_a q(s_t, a)$$

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

Above step computes avg. Return for  $(s_t, A_t)$

We can do the above using the incremental mean. For example consider  $\mu_k \leftarrow \text{avg}$ ;  $k \leftarrow \text{no. of samples}$ ,  
 $x_1, x_2 \dots x_k$  are samples,

$$\begin{aligned} \mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\ &= \frac{1}{k} \left[ x_k + \sum_{j=1}^{k-1} x_j \right] \quad \text{--- (1)} \end{aligned}$$

$$\sum_{j=1}^{k-1} x_j = \underbrace{\frac{1}{k-1} \sum_{j=1}^{k-1} x_j}_{\mu_{k-1}} (k-1)$$

eq (1) can be written as,

$$\mu_k = \frac{1}{k} \left[ x_k + (\mu_{k-1}) \cdot k-1 \right]$$

$$\mu_k = \mu_{k-1} + \frac{1}{k} \left[ x_k - \mu_{k-1} \right]$$

↪ equivalent to

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \frac{1}{N(s_t)} \left[ r - Q(s_t, A_t) \right]$$

Answer 2,

Backup diagram can be drawn as,



$$Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A') - Q(s, A)]$$

$\gamma \leftarrow$  discounting factor.

Ans.

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

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

$$q(s, a) = \frac{\sum_{t \in T(s)} P_{t:T-1} G_t}{|T(s)|}$$

In case of weighted importance sampling

$$q(s, a) = \frac{\sum_{t \in T(s)} P_{t:T-1} G_t}{\sum_{t \in T(s)} P_{t:T-1}}$$

Ans 5.

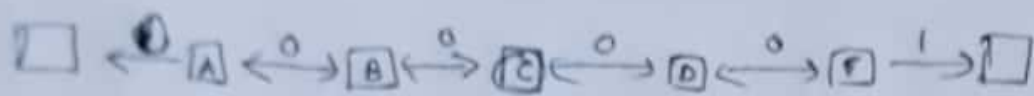
TD learning will be much better than MC learning  
in such case.

Since only the initial initial route is changed and  
other states encountered in new building will be  
same such as entry highway, secondary road.

Thus the value function estimates for the state of  
new building will be very close to the old building,  
so if we guess the initial route then function  
convergence will be faster in case of TD estimate.

Ans 6.

exercise 6.3.



$$\lambda = 1, \quad \alpha = 0.1$$

If we take TD(0) update.

$$V(s_t) = V(s_t) + 0.1 (R_{t+1} + V(s_{t+1}) - V(s_t))$$

since initial fn is const.,

first update is  $V(s_t) = V(s_t)$ .

$$V(A) = V(A) + 0.1 (0 + 0 - V(A))$$

$$V(s_{t+1}) = 0$$

$$R_{t+1} = 0.$$

$$V(A) = 0.5 \quad (\text{initial value})$$

$$V(A) = 0.45.$$

Thus in first episode state value is decreased by 0.05.

exercise 6.4.

TD performs better for wide range of values of  $\alpha$ .

exercise 6.5

$\alpha \uparrow$  (higher)  $\rightarrow$  more value  $V(x)$  update for each step.

Temporal diff. depends on return received on each step.

Thus the going down and up for RMS error may be due to randomness in the reward.

Learning takes a lot of time for smaller values of  $\alpha$ .

Ans 8,

$\phi$  learning is off-policy, because target and behaviour policy are different.

$\phi$  learning and SARSA (which is on-policy) becomes same if action selection is greedy.

We continually estimate  $q_{\pi}$  and at the same time change behaviour policy  $\pi$  towards greediness w.r.t  $q_{\pi}$ .

Thus since  $\phi$ -learning becomes on-policy as SARSA they will make same action-selection and  $Q$ -updates.