

Reinforcement Learning :-

A-3 - Theory

Donal Doitani (AI21BTECH11009)

Ans
1(a)

Given the samples, we observe that [Using first-visit MC]

$$N(A) = 5$$

$$N(D) = 4 \quad (\text{Not in 3rd sample})$$

$$N(B) = 5$$

$$N(E) = 5 = N(F)$$

$$N(C) = 5$$

Let, $G_j(A) :=$ denote discounted sum of rewards starting from the first-visit of (state A) in sample (trajectory) j .

$S(A) :=$ denote the running sum of total returns

$$S(A) = G_1(A) + G_2(A) + G_3(A) + G_4(A) + G_5(A)$$

$$= 14 + 15 + 17 + 16 + 15 = 77$$

$$\text{Sly, } S(B) = \sum_{i=1}^5 G_i(B)$$

$$= 13 + 14 + 16 + 15 + 14 = 72$$

$$S(C) = 12 + 13 + 15 + 14 + 13 = 67$$

$$S(D) = 12 + 12 + 0 + 12 + 11 = 47$$

$$S(E) = 11 + 11 + 11 + 10 + 10 = 52$$

$$S(F) = 10 + 10 + 10 + 10 + 9 = 49$$

$$\therefore V(A) = \frac{77}{5}, \quad V(B) = \frac{72}{5}, \quad V(C) = \frac{67}{5}$$

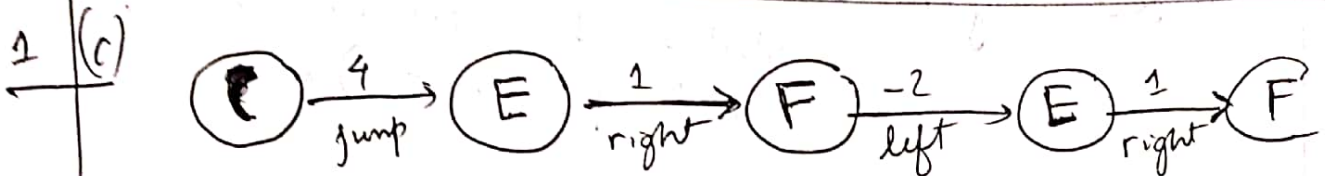
$$V(D) = \frac{47}{4}, \quad V(E) = \frac{52}{5}, \quad V(F) = \frac{49}{5}$$

Ans

1) (a) $V(G) = \frac{0}{5} = 0$ ($\therefore S(G) = 0$)

Ans

1 (b) Since, states C, E, F are visited more than once in one (or) more given episodes / samples, they are likely to have different value estimates. But, we may still get the same value.



~~Q(2, j)~~ Given, $\alpha = 0.7$, $\gamma = 1$

We'll be using the Q-learning update rule,

$$\begin{aligned}
 \text{(i)} \quad Q(C, \text{jump}) &= Q(C, \text{jump}) + 0.7(4 + 1 \max_a Q(E, a) - Q(C, \text{jump})) \\
 &= -10 + 0.7(4 + 1(-10) - (-10)) \\
 &= \boxed{-7.2}
 \end{aligned}$$

In general, the formula :-

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t [r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

$$\begin{aligned}
 \text{(ii)} \quad Q(E, \text{right}) &= -10 + 0.7[1 + 1(-10) - (-10)] \\
 &= \boxed{-9.3}
 \end{aligned}$$

$$\begin{aligned}
 \text{(iii)} \quad Q(E, \text{left}) &= -10 + 0.7[-2 + 1 \max(-9.3, -10) - (-10)] \\
 &= -10 + 0.91 = \boxed{-10.91}
 \end{aligned}$$

Contd

1(d) iv $Q(E, \text{right}) = -9.3 + 0.7 [1 + 1 * \max(-10.9, -10) - (-9.3)]$
 $= -9.3 + 0.21 = \boxed{-9.09}$

$Q(c, \text{left})$	$Q(c, \text{jump})$	$Q(E, \text{left})$	$Q(E, \text{right})$	$Q(F, \text{left})$	$Q(E, \text{right})$
-10	-10	-10	-10	-10	-10
-10	<u>-7.2</u>	-10	-10	-10	-10
-10	-7.2	-10	<u>-9.3</u>	-10	-10
-10	-7.2	-10	-9.3	<u>-10.9</u>	-10
-10	-7.2	-10	<u>-9.09</u>	-10.91	-10

1(d) $\pi(c) = \operatorname{argmax}_a Q(c, a) = \text{jump}$ $\left\{ \begin{array}{l} \therefore \text{jump} \rightarrow -7.2 \\ \text{left} \rightarrow -10 \end{array} \right.$

1(e) Robins Monroe condⁿ are:-
 $\sum \alpha_t = \infty$ and $\sum \alpha_t^2 < \infty$ (finite)

(i) $\alpha_t = \frac{1}{t}$

From calculus, we know that $\sum_{t=1}^{\infty} \frac{1}{t^p}$ converges if $p > 1$.

$\sum \alpha_t = \sum \frac{1}{t} = \infty$ [Harmonic series] [$\because p=1$]

$\sum \alpha_t^2 = \sum \frac{1}{t^2} = C < \infty$ [$\because p=2 > 1$]

Note: p-series test :- P.T.O \rightarrow
 The series $\sum_{n=1}^{\infty} \left(\frac{1}{n}\right)^p$ converges iff $\boxed{p > 1}$

$$\sum = \sum_{t=1}^{\infty}$$

$C = \text{const}$

1 (e) (ii) $\alpha_t = \frac{1}{t^2}$

$$\sum \alpha_t = \sum \frac{1}{t^2} = C < \infty \quad [\because p=2 > 1]$$

~~MA~~

So, $\alpha_t = \frac{1}{t}$ follows the condⁿ

while $\alpha_t = \frac{1}{t^2}$ violates it.

1(f) Q-learning follows Generalised Policy Iteration to find the optimal policy. In the specific case of Q-learning, we use TD approximation method (~~the off-policy one~~) to evaluate optimal state-action function.

TD Appx:

$$Q(s_t, a_t) := Q(s_t, a_t) + \alpha_t * [R(s_t, a_t, s_{t+1}) - Q(s_t, a_t) + \gamma \max_a Q(s_{t+1}, a)]$$

Where, the trajectory segment $(s_t, a_t, r_{t+1}, s_{t+1})$ is generated by the ϵ -greedy policy.

while, $\max_a Q(s_{t+1}, a) := a_{t+1}$ is greedy

w.r.t Q . [being greedy at s_{t+1}]

So, our target is greedy w.r.t $Q(s, a)$

In the f(i), all the convergence criteria for Policy EVALUATION are satisfied [Given]

1(j)

Also, the only difference in the question is that we are sampling a_t [from s_t] by following a fixed policy π with $\text{prob} = 0.5$ (or) Choose an action uniform randomly.

But, because our target policy is greedy, $Q(s_t, a_t)$ will improve for all (s_t, a_t) as there is randomness in choosing action and all s - a pairs are visited only. So, it will eventually converge to the optimal Q value f^* .

(ii)

SARSA is an on-policy, implies it uses same policy to interact with environment and learn from it too. In the given scenario, the SARSA agent will not reach the optimal Q -value f^* , since there would be no improvement in the Q -values [since the target is not improving]

Because we were following ϵ -greedy over Q to choose action ' a_{t+1} ' in SARSA, but here we are following a fixed policy or random [no-greedy] & we already proved that For an policy π , the ϵ -greedy policy π' w.r.t Q^π is an improvement over π . [Here, no sign of greedy]

Ans