

\* Assignment - 2 \*

Nelakuditi Rahul Naga

AI20BTECH11029

1) a) We know that the algorithm terminates when the following condition is satisfied:

$$\|V_{k+1} - V_k\|_{\infty} \leq \epsilon \rightarrow (1)$$

→ Now consider the norm given by  $\|V_k - V^{\pi}\|_{\infty}$ . Using triangle inequality we have:

$$\|V_k - V^{\pi}\|_{\infty} \leq \|V_k - V_{k+1}\|_{\infty} + \|V_{k+1} - V^{\pi}\|_{\infty} \rightarrow (2)$$

→ Let  $L^{\pi}$  be the Bellman evaluation operator. We know that  $L^{\pi}(V_k) = V_{k+1}$ . Also,  $L^{\pi}$  is a  $\gamma$ -contraction map. Hence we have:

$$\begin{aligned} \|V_{k+1} - V^{\pi}\|_{\infty} &= \|L^{\pi}(V_k) - L^{\pi}(V^{\pi})\|_{\infty} \quad (\because L^{\pi}(V^{\pi}) = V^{\pi}) \\ &\leq \gamma \cdot \|V_k - V^{\pi}\|_{\infty} \rightarrow (3) \end{aligned}$$

From (1), (2), (3) we have:

$$\|V_k - V^{\pi}\|_{\infty} \leq \epsilon + \gamma \cdot \|V_k - V^{\pi}\|_{\infty}$$

$$\|V_k - V^{\pi}\|_{\infty} \leq \frac{\epsilon}{1-\gamma} \rightarrow (4)$$

From (3), (4) we can conclude that:

$$\|V_{k+1} - V^{\pi}\|_{\infty} \leq \frac{\epsilon \cdot \gamma}{1-\gamma}$$

b) From equation (3) in the previous part of this question we have:

$$\|V_{k+1} - V^{\pi}\|_{\infty} \leq \gamma \cdot \|V_k - V^{\pi}\|_{\infty}$$

→ Similarly we have,  $\|V_k - V^{\pi}\|_{\infty} \leq \gamma \cdot \|V_{k-1} - V^{\pi}\|_{\infty} \Rightarrow \|V_{k+1} - V^{\pi}\|_{\infty} \leq \gamma^2 \cdot \|V_{k-1} - V^{\pi}\|_{\infty}$

Hence, applying the above inequality iteratively we obtain the following ineq-

- uality :

$$\|v_{k+1} - v^\pi\|_\infty \leq \sqrt{k} \cdot \|v_1 - v^\pi\|_\infty$$

c) Consider a state  $s \in S$ . There exists actions  $a$  and  $a'$  that satisfy the following relations:

$$L(u(s)) = R(s, a) + \sqrt{\gamma} \cdot \sum_{s'} P(s'|s, a) \cdot u(s') \rightarrow \textcircled{1}$$

$$L(v(s)) = R(s, a') + \sqrt{\gamma} \cdot \sum_{s'} P(s'|s, a') \cdot v(s') \rightarrow \textcircled{2}$$

It is easy to infer that:

$$L(v(s)) \geq R(s, a) + \sqrt{\gamma} \cdot \sum_{s'} P(s'|s, a) \cdot v(s') \rightarrow \textcircled{3}$$

The operation  $\textcircled{1} - \textcircled{3}$  gives:

$$L(u(s)) - L(v(s)) \leq \sqrt{\gamma} \cdot \sum_{s'} P(s'|s, a) \cdot (u(s') - v(s')) \rightarrow \textcircled{4}$$

→ Since  $u \leq v$  we have  $u(s') - v(s') \leq 0$ . Hence from equation  $\textcircled{4}$  we can infer that:

$$L(u(s)) - L(v(s)) \leq 0 \quad \text{for an arbitrary } s \in S$$

→ As state ' $s$ ' is chosen arbitrarily, the above equation holds for all  $s \in S$  which implies that  $L(u) \leq L(v)$ . Hence, Bellman Optimality operator is monotonic.

2) a) It is given that  $P, Q$  are two contractions defined on a normed vector space  $\langle V, \|\cdot\| \rangle$ . By definition we have:-

$$\|P(u) - P(v)\| \leq \sqrt{p} \|u - v\| \quad \forall u, v \in V \rightarrow \textcircled{1}$$

$$\|Q(u) - Q(v)\| \leq \sqrt{q} \|u - v\| \quad \forall u, v \in V \rightarrow \textcircled{2}$$

Here,  $\sqrt{p}, \sqrt{q} \in [0, 1]$ . Now consider  $P \circ Q$  map:

$$\begin{aligned} \|P \circ Q(u) - P \circ Q(v)\| &= \|P(Q(u)) - P(Q(v))\| \leq \sqrt{p} \cdot \|Q(u) - Q(v)\| \\ &\leq \sqrt{p} \cdot \sqrt{q} \cdot \|u - v\| \quad \forall u, v \in V \end{aligned}$$

Since  $\sqrt{p}, \sqrt{q} \in [0, 1)$ , we can conclude that  $P \circ Q$  is a contraction map on normed vector space  $\langle V, \|\cdot\| \rangle$ . Now consider  $Q \circ P$  map:

$$\begin{aligned} \|Q \circ P(u) - Q \circ P(v)\| &= \|Q(P(u)) - Q(P(v))\| \\ &\leq \sqrt{q} \cdot \|P(u) - P(v)\| \\ &\leq \underbrace{\sqrt{q} \cdot \sqrt{p}}_{\in [0, 1)} \|u - v\| \quad \forall u, v \in V \end{aligned}$$

Therefore,  $Q \circ P$  is also a contraction map on normed vector space  $\langle V, \|\cdot\| \rangle$ .

b) It is easy to see from part (a) of this question that the contraction coefficient for  $P \circ Q$  and  $Q \circ P$  is given by  $\sqrt{p} \cdot \sqrt{q}$ . Since  $\sqrt{p} \in [0, 1)$  and  $\sqrt{q} \in [0, 1)$ , their product i.e;  $\sqrt{p} \cdot \sqrt{q}$  also belongs to  $[0, 1)$ .

c) It is required that the value iteration algorithm converge to a unique solution. For this to happen, the operator  $B: F \circ L$  must be a contraction map i.e; both  $F$  &  $L$  should be contraction maps under  $\infty$ -norm. We have:

$$\|B(u) - B(v)\|_{\infty} \leq \sqrt{r} \cdot \|u - v\|_{\infty} \quad \exists \sqrt{r} \in [0, 1); \forall u, v \in V$$

→ Also, let  $V_*$  be the optimal value function. If the map  $B$  converges to  $V_*$ , then  $V_*$  must be a fixed point of  $B$  i.e;

$$B(V_*) = F(L(V_*)) = V_*$$

→ Under the above conditions, the value iteration algorithm converges to the unique solution  $V_*$  if we use operator  $B$  instead of  $L$ .



3) a) It is easy to see that state A is an absorbing state. Therefore the general form for a trajectory starting from state 'S' is given by  $SSSS... (K \text{ times}) A$  where  $K: 1, 2, \dots$

b) We will consider a trajectory starting from state 'S' as follows to estimate  $v(s)$  using first visit MC:

$$\text{Trajectory} = SS... (K \text{ times}) A$$

Using the above single trajectory, the estimate for  $v(s)$  using first visit MC is given by:-

$$v(s) = \frac{1+1+1+\dots (K \text{ times})}{1} \quad (\because \text{Reward for being in state } S = 1)$$

$$v(s) = K$$

c) We will consider a trajectory starting from state 'S' as follows to estimate  $v(s)$  using every visit MC:

$$\text{Trajectory} = SS... (K \text{ times}) A$$

Using the above single trajectory, the estimate for  $v(s)$  using every visit MC is given by:-

$$v(s) = \frac{K + (K-1) + (K-2) + \dots + 1}{K}$$

$$v(s) = \frac{\frac{K \cdot (K+1)}{2}}{K}$$

$$v(s) = \frac{K+1}{2}$$

d) We know that in order to calculate the true value function, we have to use the Bellman Evaluation equation which is given by:-

$$V = (I - \gamma P)^{-1} R$$

It is given that  $r = 1$ . Also,  $P$  and  $R$  are given by :-

$$P = \begin{matrix} & \begin{matrix} S & A \end{matrix} \\ \begin{matrix} S \\ A \end{matrix} & \begin{bmatrix} 1-p & p \\ 0 & 1 \end{bmatrix} \end{matrix}$$

$$R = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{matrix} S \\ A \end{matrix}$$

→ As we desire to calculate  $v(s)$ , we will exclude the terminal/absorbing state 'A' from  $P$  and  $R$ . Then we have:-

$$V(s) = (1 - (1-p))^{-1} \cdot (1)$$

$$V(s) = (p)^{-1}$$

$$V(s) = \frac{1}{p}$$

→ We have obtained the above matrix equation to find  $v(s)$  since  $P = [1-p]$  and  $R = [1]$  if we exclude the absorbing state 'A'.

e) Consider  $E[v(s)]$  for every visit MC. It is given by :-

$$\begin{aligned} E[v(s)] &= E\left[\frac{k+1}{2}\right] \\ &= \sum_{k=1}^{\infty} \left(\frac{k+1}{2}\right) \cdot (1-p)^{k-1} \cdot p \\ &= \frac{p}{2} \cdot \left( \sum_{k=1}^{\infty} (1-p)^{k-1} + \sum_{k=1}^{\infty} k \cdot (1-p)^{k-1} \right) \end{aligned}$$

Upon evaluating the summations we have:-

$$\sum_{k=1}^{\infty} (1-p)^{k-1} = \frac{1}{p}$$

$$\sum_{k=1}^{\infty} k \cdot (1-p)^{k-1} = \frac{1}{p^2}$$

$$\therefore E[v(s)] = \frac{p}{2} \left( \frac{1}{p} + \frac{1}{p^2} \right) = \frac{1 + \frac{1}{p}}{2}$$

→ We have derived in part (d) of this question that  $v(s) = \frac{1}{P}$ .

∴ for every visit MC,  $E[v(s)] \neq v(s) \Rightarrow$  Every visit

MC estimate is biased.

f) The first visit MC method converges to  $v(s)$  by the law of large numbers as the number of trajectories/experiences increase.

→ Every visit MC method also converges to  $v(s)$  as the number of trajectories increase (by the law of large numbers). But the convergence of every visit MC is not as straight-forward as convergence of first visit MC, since there might be multiple samples from the same trajectory that contribute to the mean calculation in every visit MC unlike first visit MC.

4) a) We have  $\delta_t = r_t + \gamma v^\pi(s_{t+1}) - v^\pi(s_t)$ . We have to calculate  $E_\pi(\delta_t | s_t = s)$  when  $\delta_t$  uses the true state value function  $v^\pi$ .

$$\begin{aligned} E_\pi(\delta_t | s_t = s) &= E_\pi(r_t + \gamma v^\pi(s_{t+1}) - v^\pi(s_t) | s_t = s) \\ &= E_\pi(r_t + \gamma v^\pi(s_{t+1}) | s_t = s) - E_\pi(v^\pi(s_t) | s_t = s) \\ &= v^\pi(s) - v^\pi(s) \\ &= 0. \end{aligned}$$

→ While evaluating  $E_\pi(\delta_t | s_t = s)$  above, we made use of the fact that

$$E_\pi(r_t + \gamma v^\pi(s_{t+1}) | s_t = s) = v^\pi(s) \text{ for true state value function } v^\pi.$$

b) We have  $\delta_t = r_t + \gamma v^\pi(s_{t+1}) - v^\pi(s_t)$ . We have to calculate  $E_\pi(\delta_t | s_t = s, A_t = a)$  when  $\delta_t$  uses the true state value function  $v^\pi$ .

$$\begin{aligned} E_\pi(\delta_t | s_t = s, A_t = a) &= E_\pi(r_t + \gamma v^\pi(s_{t+1}) - v^\pi(s_t) | s_t = s, A_t = a) \\ &= E_\pi(r_t + \gamma v^\pi(s_{t+1}) | s_t = s, A_t = a) - E_\pi(v^\pi(s_t) | s_t = s, A_t = a) \end{aligned}$$



Using the fact that  $E_{\pi}(\delta_{t+1} + \gamma \cdot v^{\pi}(s_{t+1}) | s_t: s, A_t: a) = Q^{\pi}(s, a)$  we have:

$$E_{\pi}(\delta_t | s_t: s, A_t: a) = Q^{\pi}(s, a) - v^{\pi}(s)$$

c) The weight corresponding to  $n$ -step return i.e;  $G_t^{(n)}$  in the expression for  $G_t^{\lambda}$  is given by  $(1-\lambda) \cdot \lambda^{n-1}$ . Clearly,  $w_1$  = Weight corresponding to  $G_t^{(1)} = 1-\lambda$ . It is given that after  $\eta(\lambda)$  time-steps, the weight would have fallen to half of the initial value i.e;  $w_1$ .

$$\therefore w_{\eta(\lambda)} = \cancel{(1-\lambda)} \cdot \lambda^{\eta(\lambda)-1} = \frac{w_1}{2} = \frac{\cancel{1-\lambda}}{2}$$

$$\Rightarrow (\eta(\lambda)-1) \cdot \ln(\lambda) = \ln\left(\frac{1}{2}\right) \quad (\text{taking logarithm on both sides})$$

$$\Rightarrow \eta(\lambda) = \frac{\ln\left(\frac{1}{2}\right)}{\ln(\lambda)} + 1$$

→ Now it is given that  $\eta(\lambda) = 3$ . Therefore we have:

$$3 = \frac{\ln\left(\frac{1}{2}\right)}{\ln(\lambda)} + 1$$

$$\Rightarrow \ln(\lambda) = \frac{\ln\left(\frac{1}{2}\right)}{2} = \ln\left(\frac{1}{\sqrt{2}}\right)$$

$$\therefore \lambda = \frac{1}{\sqrt{2}}$$

5) In order to check the divergence and convergence of  $\sum_{t=1}^{\infty} \alpha_t$  and  $\sum_{t=1}^{\infty} \alpha_t^2$  respectively, we will use the "Integral-Test". It states that:

→ If  $f(x)$  is a decreasing positive function defined on  $[1, \infty)$ , then the series  $\sum_{n=1}^{\infty} f(n)$  converges if and only if the integral  $\int_1^{\infty} f(x) dx$  converges.

→ It is easy to see that all  $\alpha_t$  given in the problem are decreasing positive functions defined on  $[1, \infty)$ . This holds true even for  $\alpha_t^2$ .

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

→ To check the convergence/divergence of  $\sum_{t=1}^{\infty} \frac{1}{t}$ , consider  $\int_1^{\infty} \frac{1}{t} dt$

$$: [\ln(t)]_1^{\infty} = \infty \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t} = \infty.$$

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t^2 = \sum_{t=1}^{\infty} \frac{1}{t^2}. \text{ Consider } \int_1^{\infty} \frac{1}{t^2} dt : \left[ -\frac{1}{t} \right]_1^{\infty} = 1 \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^2} < \infty.$$

$\therefore \alpha_t = \frac{1}{t}$  satisfies Robbins-Monroe condition  $\Rightarrow$  TD

Algorithm converges to true  $v(s)$ .

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

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t = \sum_{t=1}^{\infty} \frac{1}{t^2}. \text{ Consider } \int_1^{\infty} \frac{1}{t^2} dt = 1 \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^2} < \infty.$$

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t^2 = \sum_{t=1}^{\infty} \frac{1}{t^4}. \text{ Consider } \int_1^{\infty} \frac{1}{t^4} dt = \frac{1}{3} \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^4} < \infty.$$

$\therefore \alpha_t = \frac{1}{t^2}$  does not satisfy Robbins-Monroe condition  $\Rightarrow$  TD

Algorithm doesn't converge to true  $v(s)$ .

$$iii.) \alpha_t = \frac{1}{t^{2/3}}$$

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t = \sum_{t=1}^{\infty} \frac{1}{t^{2/3}}. \text{ Consider } \int_1^{\infty} \frac{1}{t^{2/3}} dt = \infty \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^{2/3}} = \infty.$$

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t^2 = \sum_{t=1}^{\infty} \frac{1}{t^{4/3}}. \text{ Consider } \int_1^{\infty} \frac{1}{t^{4/3}} dt = 3 \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^{4/3}} < \infty.$$

$\therefore \alpha_t = \frac{1}{t^{2/3}}$  satisfies Robbins-Monroe condition  $\Rightarrow$  TD

Algorithm converges to true  $v(s)$ .

$$iv.) \alpha_t = \frac{1}{t^{1/2}}$$

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t = \sum_{t=1}^{\infty} \frac{1}{t^{1/2}}. \text{ Consider } \int_1^{\infty} \frac{1}{t^{1/2}} dt = \infty \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^{1/2}} = \infty.$$



$$\rightarrow \sum_{t=1}^{\infty} \alpha_t^2 = \sum_{t=1}^{\infty} \frac{1}{t} \quad . \quad \text{consider } \int_1^{\infty} \frac{1}{t} dt = \infty \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t} = \infty$$

$\therefore \alpha_t = \frac{1}{t^{1/2}}$  does not satisfy Robbins-Monroe condition

which implies that TD algorithm doesn't converge to true  $v(s)$ .

\* Now we have  $\alpha_t = \frac{1}{t^p}$  where  $p > 0$ . for the Robbins-Monroe condition to be satisfied we should have the following:

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t = \sum_{t=1}^{\infty} \frac{1}{t^p} = \infty$$

$$\rightarrow \sum_{t=1}^{\infty} \alpha_t^2 = \sum_{t=1}^{\infty} \frac{1}{t^{2p}} < \infty$$

$\rightarrow$  It is easy to see that if  $p \leq 1$  we have:

$$\sum_{t=1}^{\infty} \frac{1}{t^p} \geq \sum_{t=1}^{\infty} \frac{1}{t} = \infty \Rightarrow \sum_{t=1}^{\infty} \frac{1}{t^p} = \infty$$

$$\therefore \sum_{t=1}^{\infty} \alpha_t = \infty \text{ if } p \leq 1 \rightarrow \textcircled{1}$$

$$\rightarrow \text{Now for } \sum_{t=1}^{\infty} \frac{1}{t^{2p}} \text{ we need the integral } \int_1^{\infty} \frac{1}{t^{2p}} dt = \left[ \frac{-1}{(2p-1) \cdot t^{2p-1}} \right]_1^{\infty}$$

to converge. It is easy to see that the integral converges if and only if  $2p-1 > 0 \Rightarrow p > \frac{1}{2}$ .

$$\therefore \sum_{t=1}^{\infty} \alpha_t^2 < \infty \text{ if } p > \frac{1}{2} \rightarrow \textcircled{2}$$

$\therefore$  From  $\textcircled{1}$  &  $\textcircled{2}$  we can conclude that TD algorithm converges to true  $v(s)$  only if  $\frac{1}{2} < p \leq 1$ .