

## Assignment - 02

### Reinforcement Learning

K. Surya Prakash

EE18 BTECH11026

Q1) a) Evaluate  $V(s)$  : first visit MC:

$$\times \text{ State A: } \frac{14+15+17+16+15}{5} = 15.4$$

$$\times \text{ State B: } \frac{13+14+16+15+14}{5} = 14.4$$

$$\times \text{ State C: } \frac{12+13+15+14+13}{5} = 13.4$$

$$\times \text{ State D: } \frac{12+12+12+11}{4} = 11.75$$

$$\times \text{ State E: } \frac{11+11+11+10+9}{5} = 10.4$$

$$\text{State F: } \frac{11+10+10+10+9}{5} = 8$$

$$\text{State G: (Terminal State): } 0$$

b) States {A, D} will remain the same  
{B, C, E, F} are likely to change.

Reason: {A, D} only occur only once

in the sequence in all cases when they appear. ~~mat~~

\* Rest other states are occurring more than once in an episode (atleast once).

(c)  $\pi_f$ : moves right or jump

$$V^{\pi_f}(s) = E_{\pi_f} \left( \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \right) : \gamma=1$$

\*  $G$ : terminal state :  $\boxed{V^{\pi_f}(G) = 0}$

$$* V^{\pi_f}(F) = \frac{10}{1} = 10$$

$$* V^{\pi_f}(E) = \frac{(10+1)}{1} = 11$$

$$* V^{\pi_f}(D) = \frac{(1+1+10)}{1} = 12$$

$$* V^{\pi_f}(C) = \frac{((1+12) + (4+11))}{2} = 14$$

$$* V^{\pi_f}(B) = \frac{1+14}{1} = 15$$

$$* V^{\pi_f}(A) = \frac{1+15}{1} = 16$$

d) Consider trajectories  $\{2, 3\}$

$V^{\pi_f}(s)$  : from MLE estimates

$$\hat{P}(s'|s, a) = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{L_k-1} \mathbb{1}(s_{k,t} = s, a_{k,t} = a, s_{k,t+1} = s')$$

$$\hat{R}(s'|s, a) = \frac{1}{N(s, a, s')} \sum_k \sum_t \mathbb{1}(s_{k,t} = s, a_{k,t} = a, s_{k,t+1} = s') \cdot r_{t,k}$$

\* State-G : Terminal State :  $\emptyset$

$$V^{\pi_f}(G) = 0$$

\* State-F :  $\hat{P}_{s_1, s_2}^A \Rightarrow \hat{P}(s_2 | s_1, A)$  :  $A$ : R  $\rightarrow$  right  
J  $\rightarrow$  jump

$$\hat{P}_{GF}^R = \frac{2}{2} = 1$$

$$\hat{R}_{GF}^R = \frac{10}{2} + \frac{10}{2} = 10$$

$$V^{\pi_f}(F) = \hat{P}_{GF}^R \left( R_{GF}^R + \gamma V^{\pi_f}(G) \right)$$

$$= 1(10 + \gamma(0))$$

$$= 10$$

$$\boxed{V^{\pi_f}(F) = 10}$$

Similarly:

\*

State-E:

$$\begin{aligned} V^{\pi_f}(E) &= \hat{P}_{EF}^R \left( \hat{R}_{EF}^R + \gamma V^{\pi_f}(F) \right) \\ &= 1(1+10) = \underline{\underline{11}} \end{aligned}$$

\* State-D:

$$\begin{aligned} V^{\pi_f}(D) &= \hat{P}_{DE}^R \left( \hat{R}_{DE}^R + \gamma V^{\pi_f}(E) \right) \\ &= \frac{1}{1} \left( \frac{1}{1} + 11 \right) = \underline{\underline{12}} \end{aligned}$$

\* State-C: (Both ways Right & Jump)

$$\begin{aligned} V^{\pi_f}(C) &= \hat{P}_{CD}^R \left( \hat{R}_{CD}^R + \gamma V^{\pi_f}(D) \right) + \\ &\quad \hat{P}_{CE}^J \left( \hat{R}_{CE}^J + \gamma V^{\pi_f}(E) \right) \end{aligned}$$

$$\hat{P}_{CD}^R = \frac{1}{2}, \quad \hat{P}_{CE}^J = \frac{1}{2}$$

$$\hat{R}_{CD} = \frac{1}{1} > 1, \quad \hat{R}_{CE} = \frac{4}{4} > 4$$

$$V^{\pi_f}(C) = \frac{1}{2}(1+12) + \frac{1}{2}(4+11) = \underline{\underline{14}}$$

\* State-B:

$$\hat{R}_{BC} = \frac{2}{2} > 1; \quad \hat{R}_{BE} = \frac{2}{2} > 1$$

$$V^{\pi_f}(B) = 0.1(1 + \overset{14}{14}) = \underline{\underline{15}}$$

\* State-A:

$$\begin{aligned} V^{\pi_f}(A) &= \hat{R}_{AB} \left( \hat{R}_{AB} + V^{\pi_f}(B) \right) \\ &= \frac{9}{2} \left( \frac{2}{2} + 15 \right) = \underline{\underline{16}} \end{aligned}$$

(e) Since  $V^{\pi_f}(s)$  is a unique value func.  
with infinite trajectories  
 $\rightarrow$  MC converges to  $V^{\pi_f}(s)$

while

By satisfying Robin Munoe's convergence

$$\rightarrow TD(0) \rightarrow V^{\pi_f}(s)$$

Both will converge to same value.

(f) Q-learning.

$$Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma (\max_{a'} Q(s',a')) - Q(s,a))$$

$$s \leftarrow s'$$

C jump 4:  $Q(C,J) = 0 + \frac{1}{2} (4 + 1(0) - 0) = 2$

E right 1:  $Q(E,R) = 0 + \frac{1}{2} (1 + 1(0) - 0) = 0.5$

F left -2:  $Q(F,L) = 0 + \frac{1}{2} (-2 + 0.5) = -0.75$

E right +1:  $Q(E,R) = 0.5 + \frac{1}{2} (1 + 0.5 - 0.5) = 0.75$

Q-table:

	<del><math>Q(C,L)</math></del>	<del><math>Q(C,J)</math></del>	<del><math>Q(E,L)</math></del>	<del><math>Q(E,R)</math></del>	<del><math>Q(F,L)</math></del>	<del><math>Q(F,R)</math></del>
Init	0	0	0	0	0	0
T1	0	2	0	0	0	0
T2	0	2	0	0.5	0	0
T3	0	2	0	0.5	-0.75	0
T4	0	2	0	0.75	-0.75	0



(9) Greedy policy.  $\pi(s) = \arg \max_a Q(s|a)$

$$\pi(C) = \text{Jump}$$

$$\pi(E) = \text{Right}$$

$$\pi(F) = \text{Right}$$

Q2) a)  $\sum_{t=0}^{\infty} \alpha_t = \infty$

b)  $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$

1)  $\alpha_t = 1/t$

$$\begin{aligned} \text{a. } \sum 1/t &= 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} + \dots \\ &> 1 + \frac{1}{2} + \underbrace{\frac{1}{4} + \frac{1}{4} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}_{\frac{1}{2}} + \dots \\ &> 1 + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \dots \rightarrow \infty \\ &> \infty \end{aligned}$$

$$\therefore \boxed{\sum \alpha_t = \infty}$$

$$\begin{aligned} \text{b) } \sum 1/t^2 &= \frac{1}{1} + \frac{1}{2^2} + \dots \\ &= \frac{\pi^2}{6} < \infty \end{aligned}$$

$\therefore a_t = \frac{1}{t}$  will converge

(b)  $a_t = 1/t^2$

$$\sum a_t = \sum \frac{1}{t^2} = \frac{\pi^2}{6} < \infty \neq \infty$$

# Does not satisfy condition

# Does not converge

(c)  $a_t = \frac{1}{t^{4/3}}$  - ff

ff  $\sum 1/t^p \rightarrow$  converges if  $p > 1$   
 $\rightarrow$  diverges if  $p \leq 1$

$\therefore \because p = 2/3 < 1 \rightarrow$  diverges

$$\therefore \sum a_t = \infty \quad \checkmark$$

$$\sum a_t^2 = \sum \frac{1}{t^{4/3}} \Rightarrow p = 4/3 \rightarrow \text{converges}$$

$$\sum a_t^2 < \infty \quad \checkmark$$

#  $a_t = \frac{1}{t^{4/3}}$  converges



$$(4) \quad \alpha_t = \frac{1}{t^{1/2}}$$

$$\Rightarrow \sum \alpha_t^2 = \sum \frac{1}{t} = \infty \quad \left. \vphantom{\sum \frac{1}{t}} \right\} \text{proved in (1).}$$

Does not satisfy conditions

$\therefore$  Does not converge

$\Rightarrow$  By observation

let  $\alpha_t = \frac{1}{t^p}$  to converge

~~if  $p \geq 1$  &  $p \leq \frac{1}{2}$~~   
 ~~$p > 1$  &  $p \leq \frac{1}{2}$~~   
 ~~$p \in [\frac{1}{2}, 1]$~~

~~if~~  $p \leq 1$  ,  $2p > 1$   
 $p \leq 1$  ;  $p > \frac{1}{2}$

$p \in (\frac{1}{2}, 1]$

Q4)  $|A| = K$

horizon-length = 1

$a \in A$ ;  $R^g(x) \rightarrow$  distributic of reward [0,1]

$a \sim \pi_b$  } behaviour

$\pi$  } Target

(a)  $V^\pi(s) = \mathbb{E}_\pi (r | a \sim \pi)$

$V^{\pi_b}(s) = \mathbb{E}_{\pi_b} (r | a \sim \pi_b) = r$  } from data.

$= \pi_b(a) \cdot r = 1 \cdot r = r$

$V^\pi(s) = \frac{\pi(a)}{\pi_b(a)} \cdot V^{\pi_b} = \frac{\pi(a) \cdot r}{\pi_b(a)}$

$$\boxed{\hat{V}^\pi = \frac{\pi(a)}{\pi_b(a)} \cdot r}$$

$\therefore$  Only one observation:  $\pi_b(a) = 1$

$$\boxed{\hat{V}^\pi = \pi(a) \cdot r}$$

\* Consider

$$E_{\pi}(\hat{V}^{\pi}) = E_{\pi_b} \left[ \frac{\pi(a)}{\pi_b(a)} \cdot r \right]$$

$$= \sum_a \sum_{\pi} \frac{\pi(a)}{\pi_b(a)} \cdot R^a(\pi) \cdot \pi$$

$$= \sum_a \sum_{\pi} \pi(a) R^a(\pi) \cdot \pi$$

$$= \sum_a \pi(a) \sum_{\pi} R^a(\pi) \cdot \pi$$

$$= \sum_a \pi(a) E_{\pi}(r|a)$$

$$= V^{\pi}$$

$$\boxed{E_{\pi}(\hat{V}^{\pi}) = V^{\pi}}$$

unbiased estimate

$$\begin{aligned} \text{(b)} \quad E_{\pi_b} \left[ \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} \right] &= \sum_a \pi_b(a|\cdot) \cdot \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} = \sum_a \pi(a|\cdot) \\ &= \underline{\underline{1}} \end{aligned}$$

$$\text{(c)} \quad \pi_b: \text{uniform}; \quad \pi_b(a) = \frac{1}{K}$$

$$\text{importance-sampling ratio: } \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)}$$

Since  $\pi$ : deterministic policy

$$\text{let } \pi(a) = 1 \quad \text{if } a = a' \\ = 0 \quad \text{else}$$

$$\therefore \frac{\pi(a)}{\pi_b(a)} = \frac{\pi(a \cdot)}{\pi_b(a \cdot)} = \frac{1}{1/K} = K \quad \text{if } a = a' \\ = 0/1/K = 0 \quad \text{else}$$

$$(d) \quad \text{Var}_{\pi} [V^{\pi}] = \text{Var}_{\pi_b} \left[ \frac{\pi(a \cdot)}{\pi_b(a \cdot)} \cdot r \right]$$

as  $\pi$  is deterministic

$$\sum_a R(a') \cdot \pi(a') = r$$

$$\frac{\pi(a \cdot)}{\pi_b(a \cdot)} = \begin{cases} 1/K & \text{if } a = a' \\ 0 & \text{else} \end{cases}$$

$$= E_{\pi_b} \left( \left( \frac{\pi(a \cdot)}{\pi_b(a \cdot)} \cdot r \right)^2 \right) - \left( E_{\pi_b} \left( \frac{\pi(a \cdot)}{\pi_b(a \cdot)} \cdot r \right) \right)^2$$

$$= \sum_{a \in A} \pi_b(a) \cdot \left( \frac{\pi(a \cdot)}{\pi_b(a \cdot)} \cdot r \right)^2 - \left( \sum_{a \in A} \pi_b(a) \cdot \left( \frac{\pi(a \cdot)}{\pi_b(a \cdot)} \cdot r \right) \right)^2$$

$$= \pi_b(a') \cdot \left( \frac{\pi(a' \cdot)}{\pi_b(a' \cdot)} \cdot r \right)^2 - \left( \pi_b(a') \cdot \left( \frac{\pi(a' \cdot)}{\pi_b(a' \cdot)} \cdot r \right) \right)^2$$

$$= \frac{1}{K} (K \cdot n)^2 - \left( \frac{1}{K} (K \cdot n) \right)^2$$

$$= K n^2 - n^2$$

$$= (K-1) n^2$$

$$\boxed{\text{Var}(\bar{v}) = (K-1) n^2}$$

$$(c) \quad \text{Var}_Q(M_Q) = \text{Var}_Q \left[ \frac{P_Q(n)}{Q(n)} \cdot f(n) \right]$$

$$= E_P \left[ \left( \frac{P(n)}{Q(n)} \cdot f(n) \right)^2 \right] - E_P(f(n))^2$$

$$\text{Var}_{\pi_b}(\hat{v}_b) = E_{\pi_b} \left( \frac{\pi(a)}{\pi_b(a)} \cdot x^2 \right) - (E_{\pi_b}(x))^2$$

$$= \sum_a \pi(a) \left( \frac{\pi(a)}{\pi_b(a)} \cdot x^2 \right)$$

$$- \left( \sum_a \left( \frac{1}{\pi_b(a)} \cdot x \right) \right)^2$$

$$\text{Var}_{\pi_b}(\hat{v}_b) = \left( \sum_a \pi(a) \left( \frac{\pi(a)}{\pi_b(a)} - 1 \right) \right) \left( \sum_a x^2 R^a(x) \right)$$

~~max bound~~

$$= \left( \frac{\sum \pi^2(a)}{a \pi_b(a)} - \frac{\sum \pi(a)}{a} \right) \left( \sum \pi^2 R^q(n) \right)$$

$$= \left( \frac{\sum \pi^2(a)}{a \pi_b(a)} - 1 \right) \left( \sum \pi^2 R^q(n) \right)$$

~~$\sum \pi(a) = 1$  / max value of  $\sum \pi^2$~~

~~Max bound occurs when  $R^q(n)$  is deterministic~~

---

~~max bound.~~

$$= \left( \sum \frac{1}{\pi_b(a)} - 1 \right) \left( \sum \pi^2 R^q(n) \right)$$


---

(f) If  $\pi$  completely diverges from  $\pi_b$   
 Then both trajectories does not match

Given  $P(T)$  ;  $Q(T) \approx 0$

$$\Rightarrow \frac{P(T)}{Q(T)} \rightarrow \infty$$

Else if follows



$$P(T) = P(S_0).$$

$$P(S_0, a_0, S_1, \dots) = P(S_0) \cdot \frac{1}{K} \cdot P_{S_0, S_1}^{a_0}$$

$$= P(S_0) \cdot P(a_0) \cdot P_{S_0, S_1}^{a_0} \cdot \prod_{k=1}^N (P_{a_k} \cdot P_{S_k, S_{k+1}}^{a_k})$$

$$Q(S_0, a_0, \dots) = Q(P(S_0) \cdot P(a_0) \cdot P_{S_0, S_1}^{a_0} \dots)$$

$\therefore P$  is followed by  $\pi_b$  which is uniform random.  $\Rightarrow P(a_0) = 1/K$

$\Rightarrow Q$  is deterministic.  $P(a_0) = 1$

$$\Rightarrow \frac{P(T)}{Q(T)} = \frac{1}{K^N} \quad \because N : \text{horizon length}$$

$$\Rightarrow \text{if } N \rightarrow \infty \Rightarrow \frac{P(T)}{Q(T)} \rightarrow 0$$

$$\frac{P(T)}{Q(T)} = \begin{cases} 0 & \text{if } P \text{ matches } Q \\ \infty & \text{else} \end{cases}$$

### Q3) Q-learning

$$S = \{s\} \quad A = \{a_1, a_2\}$$

$$r = 1$$

$$E(r|a_1) = E(r|a_2) = c$$

$$r \sim \mathbb{R}^{a_i} \quad i \in \{1, 2\}$$

$$a) \quad Q(s, a_1) = E(r|s, a_1) = E(r|a_1) = c$$

$$Q(s, a_2) = E(r|s, a_2) = E(r|a_2) = c$$

$$V^*(s) = \max(Q(s, a_1), Q(s, a_2)) = \underline{c}$$

(b) Prove that estimated value of  $\hat{V}^{\hat{\pi}}$  is a biased estimate of  $V^*$

$\{r_1, r_2, \dots, r_{n/2}\} \rightarrow$  Rewards obtained by choosing  $a_1$

$\{r_{n/2+1}, \dots, r_n\} \rightarrow$  By choosing  $a_2$

$$\hat{V}^{\hat{\pi}} = \max(\hat{Q}(s, a_1), \hat{Q}(s, a_2))$$

$$\hat{Q}(s, a_1) = \frac{r_1 + \dots + r_{n/2}}{n/2}$$

$$\hat{Q}(s, a_2) = \frac{r_{n/2+1} + \dots + r_n}{n/2}$$

$$\hat{V}^{\hat{n}}, \max f$$

Bias of Estimator:

$$\begin{aligned} \text{Bias}(\hat{V}^{\hat{n}}) &= E(\hat{V}^{\hat{n}}) - V^* \\ &= E(\max(\hat{Q}(s_1, a_1), \hat{Q}(s_1, a_2))) - V^* \end{aligned}$$

Since (max.) is a convex operator

$$E(\max(Q)) \geq \max E(Q)$$

$$\geq \max[E(\hat{Q}(s_1, a_1)), E(\hat{Q}(s_1, a_2))] - V^*$$

For  $n \rightarrow \infty$

$$\geq \max(c, c) - V^*$$

$$\geq c - V^*$$

$$\geq c - c$$

$$\geq 0$$

$$\boxed{\text{Bias}(\hat{V}^{\hat{n}}) \geq 0} \quad (=0 \text{ when } n \rightarrow \infty)$$

$\hat{V}^{\hat{n}}$  is a biased estimate of  $V^*$

(C)  $a_1 \rightarrow$  constant reward  $c$

$a_2 \rightarrow c + N(-0.2, 1)$

~~$E(a_1)$~~   $E(r|a_1) = c$  ;  $E(r|a_2) = E\left(\frac{c + N(-0.2, 1)}{a_L}\right)$

$$E(r|a_2) = \underline{c - 0.2}$$

$$\therefore E(r|a_1) > E(r|a_2)$$

\*  $a_1$  is a better action

b) No . TD algorithms might not always favour the best in expectation especially when trained on finite samples.

$\rightarrow$  It can happen that in every episode trained we might earn a reward of  $c + \epsilon$  ; ( $\epsilon > 0$ ) for action ' $a_2$ '

$\therefore$  Here the

Although the chance for this to happen in all episodes is less

but is not impossible

→ In this case  $Q(s, a_1) = c$   
 $Q(s, a_2) = c + \epsilon > c$

The algorithm might prefer ' $a_2$ '.

→ X ←