

1) a) It is given that  $\pi$  is fully supported on  $\pi_b$ . Hence, importance sampling can be used to estimate  $v^\pi$ . As it is given that the dataset consists of a single sample  $(a, \delta)$ , the estimate of  $v^\pi$  is given by:

$$\text{Estimate of } v^\pi : \frac{\pi(a|s)}{\pi_b(a|s)} \cdot \delta$$

The estimate of  $v^\pi$  is unbiased since  $E_{a \sim \pi_b} \left[ \frac{\pi(a|s)}{\pi_b(a|s)} \cdot \delta \right] = E_{a \sim \pi} [\delta]$ .

As  $E_{a \sim \pi} [\delta]$  is equal to true  $v^\pi$  ( $\because v^\pi = E_\pi[\delta | a \sim \pi]$ ), we can conclude that the IS estimate of  $v^\pi$  is unbiased.

$$\begin{aligned} \text{b) } E_{a \sim \pi_b} \left[ \frac{\pi(a|.)}{\pi_b(a|.)} \right] &= \sum_{a \in A} \frac{\pi(a|.)}{\pi_b(a|.)} \cdot \pi_b(a|.) \\ &= \sum_{a \in A} \pi(a|.) \\ &= 1 \end{aligned}$$

c) It is given that  $\pi_b(a|s) = \frac{1}{K} \forall a \in A$ . Also, as policy  $\pi$  is deterministic,  $\pi(a|s) = 1$  if  $\pi(s) = a$  and  $\pi(a|s) = 0$  if  $\pi(s) \neq a$ .

$$\therefore \text{IS} = \frac{\pi(a|s)}{\pi_b(a|s)} = \begin{cases} K, & \text{if } a = \pi(s) \\ 0, & \text{otherwise} \end{cases}$$

d) We have already calculated in part (a) that IS estimate of  $v^\pi$  is given by  $\frac{\pi(a|s)}{\pi_b(a|s)} \cdot \delta$ . Say  $\alpha = \frac{\pi(a|s)}{\pi_b(a|s)} \Rightarrow \text{Estimate of } v^\pi = \alpha \cdot \delta$ . As  $\pi_b$  is a uniformly random policy, we can assume that action 'a' is drawn from uniform distribution  $U$  i.e;  $a \sim U$ . We need to calculate  $V[\alpha \delta | a \sim U]$  where  $V[\cdot]$  denotes variance.

$$V[\alpha \delta | a \sim U] = \delta^2 \cdot V[\alpha | a \sim U] \quad (\because \delta \text{ is a constant})$$

→ Now we know that for a random variable  $X$ ,  $V[X] = E[X^2] - (E[X])^2$ .

Therefore we have:

$$V[\alpha\delta|a \sim U] = \delta^2(E[\alpha^2|a \sim U] - (E[\alpha|a \sim U])^2)$$

→ We have already proved in part (b) that  $E_{a \sim \pi_b} \left[ \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} \right] = E_{a \sim \pi_b}(\alpha) = 1$ .

Here we have  $\pi_b = U \Rightarrow E_{a \sim U}(\alpha) = E[\alpha|a \sim U] = 1$ .

$$\therefore V[\alpha\delta|a \sim U] = \delta^2(E[\alpha^2|a \sim U] - 1)$$

→ Now,  $\alpha^2 = \begin{cases} K^2, & \text{if } a = \pi(s) \\ 0, & \text{otherwise} \end{cases}$  (from part (c)). Therefore we

have:

$$E[\alpha^2|a \sim U] = K^2 \cdot \frac{1}{K} + 0 \cdot \frac{1}{K} + \dots + 0 \cdot \frac{1}{K} = K.$$

$$\therefore V[\alpha\delta|a \sim U] = \delta^2(K-1).$$

e) We now have a reward distribution which is bounded in the range  $[a, 1]$ .

The variance of IS estimate is given by  $V[\alpha\delta|a \sim U]$ . As  $V[X] = E[X^2] - (E[X])^2$

we have:

$$V[\alpha\delta|a \sim U] = E[\alpha^2\delta^2|a \sim U] - (E[\alpha\delta|a \sim U])^2 \leq E[\alpha^2\delta^2|a \sim U]$$

→ As  $\delta \in [0, 1]$ , we have:

$$E[\alpha^2\delta^2|a \sim U] \leq E[\alpha^2|a \sim U] = K \text{ (from part (d))}$$

$$\therefore V[\alpha\delta|a \sim U] \leq K.$$

f) As  $P(T)$  is a joint distribution over the trajectory  $T$  induced by target policy  $\pi$ , we have:

$$P(T) = \mu(s_0) \cdot \prod_{t=0}^{\infty} (P'(s_{t+1}|s_t, a_t) \cdot \pi(a_t|s_t))$$

where  $\mu$  is the initial start state distribution and  $P'(s_{t+1}|s_t, a_t)$  denotes the state transition probability from  $s_t$  to  $s_{t+1}$  via action  $a_t$ .

Similarly,  $Q(T)$  is given by:

$$Q(T) = u(s_0) \cdot \prod_{t=0}^{\infty} (P'(s_{t+1}|s_t, a_t) \cdot \pi_b(a_t|s_t))$$

Therefore IS weight  $\frac{P(T)}{Q(T)}$  is given by:

$$\frac{P(T)}{Q(T)} = \frac{u(s_0) \cdot \prod_{t=0}^{\infty} (P'(s_{t+1}|s_t, a_t) \cdot \pi(a_t|s_t))}{u(s_0) \cdot \prod_{t=0}^{\infty} (P'(s_{t+1}|s_t, a_t) \cdot \pi_b(a_t|s_t))} = \prod_{t=0}^{\infty} \left( \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \right)$$

2) I have submitted my solution to this question in a ipynb notebook named Assignment 3 - AI20BTECH11029. This file is included in the zip file that I have submitted.