

## Assignment - 3

①①  $\pi$  is the target policy  
 $\pi_b$  is the behaviour policy,  
 the importance sampling ratio can be  
 written as  $\frac{\pi(a|s)}{\pi_b(a|s)}$

$$v^\pi(s) = \mathbb{E}_{a \sim \pi} (v) = \mathbb{E}_{a \sim \pi_b} \left[ \frac{\pi(a|s)}{\pi_b(a|s)} v \right]$$

The unbiased IS estimate of  $v^\pi$  is given  
 by  $\frac{\pi(a|s)}{\pi_b(a|s)} v$ .

$$\begin{aligned} \textcircled{1} \textcircled{2} \quad \mathbb{E}_{a \sim \pi_b} \left[ \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} \right] &= \sum_{a \in A} \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} \pi_b(a|\cdot) \\ &= \sum_{a \in A} \pi(a|\cdot) \end{aligned}$$

$$\mathbb{E}_{a \sim \pi_b} \left[ \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} \right] = 1$$

$$\textcircled{1} \textcircled{3} \quad \pi_b(a|\cdot) = \frac{1}{K} \quad (\text{as there are } K \text{ actions which are equiprobable})$$

$$\pi(a|\cdot) = \begin{cases} 1 & \pi(s)=a \\ 0 & \text{otherwise} \end{cases} \quad (\text{we know that } \pi(s)=a)$$

so the probability that from state  $s$  with policy  $\pi$  the action taken  $a$  is 1.

$$\text{importance sampling ratio} = \frac{1_{\pi(s)=a}}{1/K}$$

(1) d

$$\text{var}_{\pi}(\bar{y}) = \text{var}_{\pi_b}(\bar{y})$$

$$\text{where } \rho = \frac{\pi(a_i)}{\pi_b(a_i)} = \frac{1}{K} = k$$

$$= \sigma^2 \text{var}_{\pi_b}(\rho)$$

$$= \sigma^2 \left[ \mathbb{E}_{\pi_b}(\rho^2) - (\mathbb{E}_{\pi_b}(\rho))^2 \right]$$

$$= \sigma^2 \left[ \frac{K^2}{K} - 1 \right]$$

$$= \sigma^2 (K - 1)$$

(1) e

$$\text{var}_{\pi}(\bar{y}) \leq \text{var}_{\pi_b}(\bar{y}) \quad \left( \rho = \frac{\pi(a_i)}{\pi_b(a_i)} \right)$$

$$= \sigma^2 \text{var}_{\pi_b}(\rho)$$

$$\leq \sigma^2 \mathbb{E}_{\pi_b}(\rho^2)$$

(maximum  
value  
for  
possible  
w)

$$= \mathbb{E}_{\pi_b}(\rho^2)$$

$$= \frac{K^2}{K}$$

$$\text{var}_{\pi}(\bar{y}) < K$$



①. given  $\tau$  is the state action sequence given by  $s_0, a_0, \dots, s_t, a_t, \dots$  we can write the probability of finding the trajectory  $\tau$  with policy  $\pi$  as

$$P(\tau/\pi) = \mu(s_0) \prod_{t=0}^{\infty} \pi(a_t/s_t) P(s_{t+1}/s_t, a_t)$$

similarly, with policy  $\pi_b$ , the probability can be written as

$$P(\tau/\pi_b) = \mu(s_0) \prod_{t=0}^{\infty} \pi_b(a_t/s_t) P(s_{t+1}/s_t, a_t)$$

we can see from the above expressions  $\mu(s_0)$  and  $P(s_{t+1}/s_t, a_t)$  do not depend on the policy.

importance sampling ratio =  $\frac{P(\tau/\pi)}{Q(\tau/\pi_b)}$

$$= \frac{\mu(s_0) \prod_{t=0}^{\infty} \pi(a_t/s_t) P(s_{t+1}/s_t, a_t)}{\mu(s_0) \prod_{t=0}^{\infty} \pi_b(a_t/s_t) P(s_{t+1}/s_t, a_t)}$$

$$= \prod_{t=0}^{\infty} \frac{\pi(a_t/s_t)}{\pi_b(a_t/s_t)}$$