Q1] Importance Sampling

(a)  Dataset $= \{(a, r)\}$

$$v^{\pi}(s) = E_{\pi}[r \mid a \sim \pi]$$

$$= E_{\pi_b}\left[ \frac{\pi(a|s)}{\pi_b(a|s)} \cdot r \mid a \sim \pi_b \right]$$

$$= \frac{\pi(a|s)}{\underbrace{\pi_b(a|s)}_{\rho}} \cdot r$$

[∵ There is only data pt in the dataset]

It is an unbiased est.

$$= \rho r$$

(b)  $$E_{\pi_b}\left[ \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} \right] = \sum_{a_i \in A} \frac{\pi(a_i|\cdot) \cdot \pi_b(a_i|\cdot)}{\pi_b(a_i|\cdot)}$$

$$= \sum_{a_i \in A} \pi(a_i|a) = 1$$

(∵ $\pi_b$ fully supports $\pi$, $\forall a \in A$ if $\pi(a) > 0 \Rightarrow \pi_b(a) > 0$]

(c)  $\pi_b$ is a uniformly random policy
$$a \sim U \Rightarrow \pi_b(a|s) = \frac{1}{k} \quad \text{(If there are total } k \text{ actions)}$$
$$|A| = k$$
$\pi$ is a deterministic policy

Imp sampling ratio $\rho = \dfrac{\pi(a|\cdot)}{\pi_b(a|\cdot)} = \dfrac{\pi(a|\cdot)}{1/k}$

$$\pi(a|s) = \begin{cases} 1 & a = \pi(s) \\ 0 & \text{otherwise} \end{cases}$$

$$\rho(s) = \begin{cases} k & a = \pi(s) \\ 0 & \text{otherwise} \end{cases}$$

$$= \mathbb{1}_{a = \pi(s)} \cdot k$$

(d) Reward function is deterministic ~~re~~

ie $R(a) = r$

$\pi_b$ — Uniform behaviour policy $\qquad a \sim U$

$\pi$ — Deterministic target policy.

~~Variance~~

$$\text{Var}[v^\pi] = \text{Var}\left[\rho r \mid a \sim U\right]$$

$$= r^2 \, \text{Var}\left[\rho \mid a \sim U\right]$$

$$= r^2 \left[ \text{Var}\left( \frac{\pi(a\mid)}{\pi_b(a)} \mid a \sim U \right) \right]$$

$$= r^2 \left[ E_U\left[ \frac{\pi(a\mid)}{\pi_b^2(a)} \mid a \sim U\right] - \underbrace{E_U\left[ \frac{\pi(a\mid)}{\pi_b(a)} \mid a \sim U\right]^2}_{=1 \text{ from } 1(b)} \right]$$

$$= r^2 \left[ E\left( \left[ \mathbb{1}_{a = \pi(s)} \cdot k\right]^2 \mid a \sim U\right) - 1 \right]$$

$$\underbrace{\sum_a \frac{1}{k} \mathbb{1}_{a=\pi(s)} k^2 = \frac{1}{k} \cdot k^2 + 0 \cdot k^2 = k}_{}$$

$$= r^2(k-1)$$

(e)

$$\text{Var}[v^\pi] = \text{Var}[\rho r] = E_{\pi_b}[\rho^2 r^2] - \left(E[\rho r]\right)^2$$

$$\leq E(\rho^2 r^2)$$

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

$$= E_\pi\left[ \frac{\pi(a)}{\pi_b(a)} \cdot r^2 \right] \qquad \leq E_\pi\left[ \frac{\pi(a)}{\pi_b(a)} \right] \qquad \text{for } r \in [0,1]$$

$$= \sum_a \pi(a) \left[ \frac{\pi(a)}{\pi_b(a)} \right] \cdot \cancel{R^2 r^2} = k$$

(f) Trajectory: $\tau: s_0, a_0, s_1, a_1, \ldots s_t, a_t$

$$s_0 \sim \mu(s_0)$$

$$P \rightarrow \pi$$
$$Q \rightarrow \pi_b$$

Prob of occurrence of $\tau$ using policy $\pi$

$$p(\tau; \pi) = \mu(s_0) \prod_{t=1}^{\infty} \pi(a_t | s_t) \, p(s_{t+1} | a_t, s_t)$$

Is wt: $\dfrac{P(\tau)}{Q(\tau)} = \dfrac{p(\tau; \pi)}{p(\tau; \pi_b)} = \dfrac{\mu(s_0) \prod_{t=1}^{\infty} \pi(a_t | s_t) \, p(s_{t+1} | a_t, s_t)}{\mu(s_0) \prod_{t=1}^{\infty} \pi_b(a_t | s_t) \, p(s_{t+1} | a_t, s_t)}$

$$= \prod_{t=1}^{\infty} \dfrac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)} \quad //.$$