Problem 1)

a) Update rule in the algorithm is

$$V(s) \leftarrow V(s) + \alpha_t \left[ \underbrace{\frac{\pi(a|s)}{\pi_b(a|s)} r - V(s)}_{\text{horizon length = 1}} \right]$$

Assuming initialization $V(s) = 0$

$$V(s) = \alpha_t \left[ \frac{\pi(a|s)}{\pi_b(a|s)} r \right] = \alpha_1 \left[ \frac{\pi(a|s) r}{\pi_b(a|s)} \right]$$

assuming $\alpha_1 = 1$

$V(s) = \frac{\pi(a|s)}{\pi_b(a|s)} r$ , there is only 1 state

$$\therefore \hat{V}^{\pi} = \frac{\pi(a|\bullet)}{\pi_b(a|\bullet)} r$$

$$E[\hat{V}^{\pi}] = E_{\pi_b} \left[ \frac{\pi(a|s)}{\pi_b(a|s)} r \right] := \sum_{a \in A} \pi_b(a|s) \frac{\pi(a|s)}{\pi_b(a|s)} r$$

actions are sampled from $\pi_b$

$$= E_{\pi}[r]$$

$\therefore \hat{V}^{\pi}$ is an unbiased estimator of $V^{\pi}$

**b)**

$$\mathbb{E}_{\pi_b}\left[\frac{\pi(a|s)}{\pi_b(a|s)}\right] = \sum_{a \in A} \pi_b(a|s) \frac{\pi(a|s)}{\pi_b(a|s)}$$

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

**ⓒ**

$$\pi_b(a|s) = 1/k$$

$$\pi(a|s) = \begin{cases} 1 & \text{if } a = a' \\ 0 & \text{o/w} \end{cases}$$

$$IS = \frac{\pi(a|s)}{\pi_b(a|s)} = \begin{cases} K & \text{if } a = a' \\ 0 & \text{o/w} \end{cases}$$

**ⓓ**

$$\hat{v}^\pi = \frac{\pi(a|s)}{\pi_b(a|s)} r$$

$$Var(\hat{v}^\pi) = \mathbb{E}_{\pi_b}\left[(\hat{v}^\pi)^2\right] - \mathbb{E}_{\pi_b}\left[\hat{v}^\pi\right]^2$$

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

$$= Kr^2 - (r)^2 = r^2(K-1)$$

(d)

From part – a | Given

$$\hat{v}^{\pi} = \frac{\pi(a|\cdot)}{\pi_b(a|\cdot)} r \qquad \Big| \quad R(a) = r \;\; \forall a$$

we know that

$$Var(x) = E[x^2] - E[x]^2$$

$$Var(\hat{v}^{\pi}) = E[(\hat{v}^{\pi})^2] - E[\hat{v}^{\pi}]^2$$

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

from last subpart

$$\pi(a|\cdot) = \begin{cases} 1 & a = a' \\ 0 & o/w \end{cases}$$

$$\pi_b(a|\cdot) = \{ 1/k \quad \forall a \in A$$

$$= K r^2 - [r]^2 = \underbrace{r^2(K-1)}$$

This is the variance of $\hat{v}^{\pi}$

e)

Now $r$ is bounded between $[0,1]$

$$Var_{\pi_b}(\hat{V}^\pi) = E_{\pi_b}\left[\frac{\pi^2(a|s)}{\pi_b^2(a|s)} r_a^2\right] - E_{\pi_b}\left[\frac{\pi(a|s)}{\pi_b(a|s)} r_a\right]^2$$

$\pi$ is a deterministic policy

$$\pi(a|\cdot) = \begin{cases} 1 & a = a' \\ 0 & o/w \end{cases}$$

$$\pi_b(a|\cdot) = \{1/k \quad \forall a \in A$$

$$Var_{\pi_b}(\hat{V}^\pi) = \sum_{a \in A} \frac{\pi^2(a|s)}{\pi_b(a|s)} r_a^2 - \left(\sum_{a \in A} \pi(a|s) r_a\right)^2$$

$$= K r_{a'}^2 - r_{a'}^2 \qquad [r_{a'} \text{ is a } rv]$$

$$= r_{a'}^2 (K-1)$$

$$\leq (K-1)$$

$$\therefore r_{a'} \in [0,1]$$

f)

Let ~~the~~ trajectory length $(\tau)$ ~~be~~ of some ~~size t~~

Let us consider the trajectory $\tau$ till length $\ell$

$$P(\tau_\ell) = \begin{cases} 1 & \text{if } \forall\, s_i, a_i : \pi(a_i | s_i) = 1 \\ 0. & o/w \end{cases}$$

assuming a deterministic environment

$$Q(\tau_i) = K^{-\ell} .$$

$$\therefore \frac{P(\tau_\ell)}{Q(\tau_\ell)} = \begin{cases} 0 & \text{if } \tau_\ell \text{ isn't possible under } \pi \\ K^\ell & o/w \end{cases}$$

$$\frac{P(\tau)}{Q(\tau)} = \lim_{\ell \to \infty} \frac{P(\tau_\ell)}{Q(\tau_\ell)} = \begin{cases} \infty & \text{if } \tau \sim \pi \\ 0 & o/w \end{cases}$$