

Ans 1)

Initialize:

$\pi(s) \in A(s)$ (arbitrarily); $\forall s \in S$

$Q(s, a) \in \mathbb{R}$ (arbitrarily); $\forall s \in S, \forall a \in A$

Returns-sum(s, a) $\in \mathbb{R}$ (initialize all values to 0)

Returns-count(s, a) $\in \mathbb{R}$ (initialize all values to 0)

loop forever (for each episode):

choose $s_0 \in S, A_0 \in A(s_0)$ randomly such that
all pairs have probability > 0

generate an episode from s_0, A_0 following

$\pi: s_0, A_0, R_1, \dots, s_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

loop for each step of the episode, $t = T-1 \dots 0$

$G \leftarrow \gamma G + R_{t+1}$

Returns-sum(s_t, A_t) $+= G$

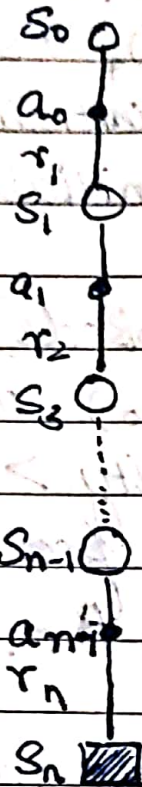
Returns-count(s_t, A_t) $+= 1$

$Q(s_t, A_t) = \text{Returns-sum}(s_t, A_t) / \text{Returns-count}(s_t, A_t)$

$\pi(s_t) = \underset{a}{\operatorname{argmax}} (Q(s_t, a))$

Ans 2

$$q_{\pi}(s, a) = r + V_{\pi}(s', a) \\ = q_{\pi}(s, a) + \alpha (G_t - q_{\pi}(s, a))$$



Ans 3

$$q_{\pi}(s, a) = E_{\pi}(G_t | S_t = s, A_t = a)$$

Now, $Pr(A_{t+1}, S_{t+1}, \dots, S_T | S_t = s, A_t = a, A_{t+1:T-1} \sim \pi)$

$$= \pi(A_{t+1} | S_{t+1}, A_t = a)$$

$$= p(s_{t+1} | s, a) \pi(A_{t+1} | s_{t+1}) \dots \pi(A_{T-1} | s_{T-1}) \\ \times p(s_T | A_{T-1}, s_{T-1})$$

$$\Rightarrow P_T = \prod_{k=t+1}^{T-1} \pi(A_k | S_k) P(S_{k+1} | S_k, A_k)$$

$$\Rightarrow \rho_{t+1:T-1} = \prod_{k=t+1}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$$

Also, $t+1:T-1$ can be written as $t \in \tau(s, a)$ where $\tau(s, a)$ is trajectory after starting from s, a

$$\Rightarrow q_{\pi}(s, a) = \frac{\sum_{t \in \tau(s, a)} \rho_{t:T(t)-1} G_t}{\sum_{t \in \tau(s, a)} \rho_{t:T(t)-1}}$$