

(1A).

$$(a): d^{\pi_{\theta}}(s) = \sum_{k=0}^{\infty} \Pr\{s_0 \rightarrow s, k, \pi_{\theta}\}.$$

$Q^{\pi_{\theta}}(s, a)$  = total rewards when action  $a$  is taken in state  $s$  and  $\pi_{\theta}$  is followed subsequently.

(b). Both are fine as:

$$\max_a Q^{\pi}(s, a) \neq \sum_b \frac{Q^{\pi}(s, b)}{1A(s)} \quad \text{do not depend on given action } a.$$

$$(2A). (a): 0.5 \cdot 0.5 \cdot 0.2 \left( \frac{1}{0.5 \cdot 0.5} \right) (1) + 0.5 \cdot 0.5 \cdot 0.8 \cdot 0.5 \cdot 0.5 \cdot 0.2 \left( \frac{1}{(0.5)^4} \right) + \dots$$
$$= 0.2 \left[ \sum_{t=0}^{\infty} (0.8)^t \right] = 1.$$

$$(b). \text{Variance} = 0.2 \sum_{t=0}^{\infty} (0.8)^t 2^{2t+2} = \infty$$

(c). Weighted I.S works well.

Reason: Consider trajectories  $\tau_1, \tau_2, \dots, \tau_n$  and the objective is to estimate  $V^{\pi}(s_1)$ . Let  $I \subseteq \{1, 2, \dots\}$  &  $G_I = 1$ .

$$\text{using weighted I.S: } V^{\pi}(s_1) = \frac{\sum_{i \in I} s_i}{\sum_{i \in I} 1} = 1. \text{ exactly } V^{\pi}(s_1).$$

$$\text{using ordinary I.S: } V^{\pi}(s_1) = \frac{\sum_{i \in I} s_i}{n}, \text{ need not be 1.}$$

(3A). (a). In Monte Carlo, Value function is updated after entire trajectory is generated. Whereas, in TD(0), update is done sample-by-sample.

As a result, MC has large variance, but less bias. If there is a very large availability of data (trajectories), one can prefer MC over TD(0).

(b). For MC,  $V(B) = 1$  &  $V(A) = 0$ . [Took first visit MC. Every visit MC is also fine!].

$$\text{TD}(0): V(B) = (0.5)(0) + 0.5(1) = 0.5$$

$$V(B) = (0.5)(0.5) + 0.5(1) = 0.75$$

$$V(B) = (0.5)(0.75) + 0.5(1) = 0.875$$

$$V(A) = 0.5(0) + 0.5[0 + 0.875] = 0.4375$$

$$V(B) = (0.5)(0.875) + 0.5(0 + 0) = 0.4375$$

Note: Took  $\gamma = 1$ .

Any other value of

$\gamma$  is also fine!

(4A) Refer textbook.

(5A).

(a). Let ' $w$ ' be the parameter.

For the sample  $(S_t, a_t, r_t, S_{t+1})$  we have:

$$w_{t+1} \doteq w_t + \alpha [r_t + \gamma w_t^T \chi(S_{t+1}) - w_t^T \chi(S_t)] \chi(S_t).$$

(b). First, let us compute stationary distribution:

$$[\pi_1 \ \pi_2] \begin{bmatrix} 0.7 & 0.3 \\ 1 & 0 \end{bmatrix} = [\pi_1 \ \pi_2]$$

$$0.3 \pi_1 = \pi_2 \quad \& \quad \pi_1 + \pi_2 = 1$$

$$\Rightarrow 1.3 \pi_1 = 1 \Rightarrow \pi_1 = \frac{10}{13} \quad \& \quad \pi_2 = \frac{3}{13}.$$

$$\begin{aligned}
 \text{Now, } b &= \frac{10}{13} [0.7 \times 5 + 0.3 \times 3] \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \\
 &\quad \frac{3}{13} [3] \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 53/13 \\ 9/13 \end{bmatrix}
 \end{aligned}$$

4

$$\begin{aligned}
 A &= \frac{10}{13} \left[ 0.3 \left[ \begin{bmatrix} 1 \\ 0 \end{bmatrix} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \right)^T \right] + \right. \\
 &\quad \left. 0.7 \left[ \begin{bmatrix} 1 \\ 0 \end{bmatrix} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.5 \\ 0 \end{bmatrix} \right)^T \right] \right] + \\
 &\quad \frac{3}{13} \left[ \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0.5 \\ 0 \end{bmatrix} \right)^T \right] + \\
 &= \frac{10}{13} \left[ 0.3 \begin{bmatrix} 0.5 & -0.5 \\ 0 & 0 \end{bmatrix} + 0.7 \begin{bmatrix} 0.5 & 0 \\ 0 & 0 \end{bmatrix} \right] + \\
 &\quad \frac{3}{13} \begin{bmatrix} 0.5 & 1 \\ 0.5 & 1 \end{bmatrix} \\
 &= \frac{10}{13} \begin{bmatrix} 0.5 & -0.15 \\ 0 & 0 \end{bmatrix} + \frac{3}{13} \begin{bmatrix} 0.5 & 1 \\ 0.5 & 1 \end{bmatrix} \\
 &= \begin{bmatrix} 0.5 & 1.5/13 \\ 1.5/13 & 3/13 \end{bmatrix}
 \end{aligned}$$

$\therefore$  The point of convergence is  $\bar{A}^{-1}b$ .

$$\approx \begin{bmatrix} 2.26 & -1.12 \\ -1.12 & 4.59 \end{bmatrix} \begin{bmatrix} 53/13 \\ 9/13 \end{bmatrix}$$

12

$$\begin{bmatrix} 8.434 \\ -1.184 \end{bmatrix}$$