

## Q 2. Exercise 2.6

(along with others)

In optimistic greedy approach, the agent tends to explore more in the initial steps. This is due to the dissatisfying rewards it receives i.e. the actual rewards received are much less than the optimistic action value estimates and hence different actions are tried. As optimal action would also be chosen, hence ~~in~~ no. of times it is chosen (across several runs) in the initial steps gives spikes.

Q 3. The incremental rule for estimated reward  $Q$  is given by :

$$Q_{n+1} = Q_n + \underset{\substack{\downarrow \\ \text{step size}}}{\alpha} [R_n - Q_n]$$

with step size  $\beta_n = \frac{\alpha}{\bar{Q}_n}$ , we have  $\bar{Q}_n = \bar{Q}_{n-1} + \alpha(1 - \bar{Q}_{n-1})$

$$\begin{aligned} Q_{n+1} &= Q_n + \beta_n (R_n - Q_n) \\ &= \beta_n R_n + \cancel{Q_n} - \beta_n Q_n + Q_n \\ &= \beta_n R_n + (1 - \beta_n) Q_n \\ &= \beta_n R_n + (1 - \beta_n) (Q_{n-1} + \beta_{n-1} (R_{n-1} - Q_{n-1})) \\ &= \beta_n R_n + (1 - \beta_n) \beta_{n-1} R_{n-1} + (1 - \beta_n) Q_{n-1} - \beta_{n-1} Q_{n-1} (1 - \beta_n) \\ &= \beta_n R_n + (1 - \beta_n) \beta_{n-1} R_{n-1} + (1 - \beta_n) (1 - \beta_{n-1}) Q_{n-1} \\ &\vdots \\ &= \beta_n R_n + (1 - \beta_n) (\beta_{n-1} R_{n-1}) - \dots - + \prod_{i=1}^n (1 - \beta_i) Q_1 \end{aligned}$$

Now as  $\bar{Q}_0 = 0 \therefore \bar{Q}_1 = 0 + \alpha(1 - 0) = \alpha$

$$\Rightarrow \boxed{\beta_1 = 1}$$

$\therefore \prod_{i=1}^n (1 - \beta_i) = 0 \therefore$  Initial bias is eliminated