

Question 3

When we use sample average as our action value method, ~~we~~ Q_n is independent of Q_1 .

However, ~~we~~ if we use a constant step size parameter, we have

$$Q_n = Q_n + \alpha [R_n - Q_n]$$

$$\begin{aligned} Q_{n+1} &= \alpha R_n + (1-\alpha) Q_n \\ &= \alpha R_n + (1-\alpha) (\alpha R_{n-1} + (1-\alpha) Q_{n-1}) \\ &= (1-\alpha)^n Q_1 + \sum_{j=1}^n \alpha (1-\alpha)^{n-j} R_j \end{aligned}$$

~~The~~ We can see that Q_{n+1} still depends on Q_1 irrespective of the action ~~which~~ picked by the agent.

To get an independent estimate using a ~~constant~~ step size, we can use

$$\beta_n = \frac{\alpha}{\overline{Q}_n} \quad (\alpha \rightarrow \text{constant step size})$$

$$\overline{Q}_n = \overline{Q}_n + \alpha (1 - \overline{Q}_n) \quad \text{for } n \geq 0$$

$$\overline{Q}_0 = 0$$

$$Q_2 = Q_1 + \beta_1 [R_1 - Q_1]$$

$$\beta_1 = \frac{\alpha}{(1-\alpha) \times Q_0 + \alpha} = \frac{\alpha}{\alpha}$$

$$\therefore Q_2 = Q_1 + R_1 - Q_1$$

\therefore We can see that Q_2 is independent of our initial estimate Q_1 and only depends on R_1 [reward]

We can extend it further,

$$Q_3 = Q_2 + \beta_2 [R_2 - Q_2]$$

$$\beta_2 = \frac{\alpha}{(1-\alpha) \times Q_1 + \alpha} = \frac{1}{2-\alpha}$$

$$\therefore Q_3 = Q_2 + \frac{1}{2-\alpha} [R_2 - Q_2]$$

$$= \frac{R_2}{2-\alpha} + Q_2 \left[\frac{1 - \frac{1}{2-\alpha}}{2-\alpha} \right]$$

$$= \frac{R_2}{2-\alpha} + Q_2 \left[\frac{1-\alpha}{2-\alpha} \right]$$

We know that $Q_2 = R_1$,

$$\therefore Q_3 = \frac{R_3}{2-\alpha} + R_1 \left[\frac{1-\alpha}{2-\alpha} \right]$$

which is independent of Q_1 .

Similarly we can know for
 $Q_4, Q_5 \dots Q_n$.

Hence $\beta = \frac{\alpha}{\bar{Q}_n}$, $\bar{Q}_n = Q_{n-1}(1-\alpha) + \alpha$

is a parameter that removes bias and provides constant step size.