

RL Theory Assignment - 1

Ans 2.6

Initial action-value estimates are equal & very large. So, there is equal probability to pick ~~the first~~ any bandit. The reward it produces would most likely be lower than the initial estimate. Hence action-value estimate would reduce. This will happen for 10 times for each bandit. on 11th step we would choose the highest value estimate bandit which would likely be the optimal action. Hence, we would get a high reward & observe a spike.

Ans 2.7

$$\beta_n = \frac{\alpha}{\bar{Q}_n}, \quad \bar{Q}_n = \bar{Q}_{n-1} + \alpha(1 - \bar{Q}_{n-1}), \quad \bar{Q}_0 = 0$$

|
step size

$$\begin{aligned} \text{We know that - } Q_{n+1} &= Q_n + \beta_n(R_n - Q_n) \\ &= Q_n(1 - \beta_n) + \beta_n R_n \end{aligned}$$

$$\begin{aligned} Q_{n+1} &= [Q_{n-1} + \beta_{n-1}(R_{n-1} - Q_{n-1})](1 - \beta_n) + \beta_n R_n \\ \Rightarrow Q_{n+1} &= Q_{n-1}(1 - \beta_n)(1 - \beta_{n-1}) + (1 - \beta_n)\beta_{n-1}R_{n-1} + \beta_n R_n \end{aligned}$$

As we want to prove that Q_{n+1} is independent of Q_1 , we will look only at coefficient of Q_1 :

$$\Rightarrow Q_{n+1} \text{ depends on } (1 - \beta_1)(1 - \beta_2) \dots (1 - \beta_n) Q_1 \quad \text{--- (1)}$$

$$\text{now, } \beta_1 = \frac{\alpha}{\bar{Q}_1} = \frac{\alpha}{\bar{Q}_0 + \alpha(1 - \bar{Q}_0)} = 1$$

Put in eq. (1) we see that coefficient of Q_1 is 0
 $\Rightarrow Q_{n+1}$ is independent of Q_1 & only depends on R_1, R_2, \dots, R_n

$$\begin{aligned} Q_{n+1} &= \beta_1(1 - \beta_2)(1 - \beta_3) \dots (1 - \beta_n) R_1 + \beta_2(1 - \beta_3) \dots (1 - \beta_n) R_2 + \dots + \beta_n R_n \\ &= \sum_{i=1}^n \beta_i R_i \prod_{j=i+1}^n (1 - \beta_j) \end{aligned}$$