

Q1:

$A_0=0, R_0=0, \boxed{Q_1=0}, Q_2=0, Q_3=0, Q_4=0.$   
 $\downarrow$  Choose 1. Maybe random. Cuz  $Q_1=Q_2=Q_3=Q_4=0$  } Consistent with  $\epsilon$ -greedy action selection.

$A_1=1, R_1=-1, Q_1=\frac{-1}{1}=-1, \boxed{Q_2=0}, Q_3=0, Q_4=0.$   
 $\downarrow$  Choose 2. Maybe random. Cuz  $Q_1=-1 < Q_2=Q_3=Q_4=0$  } Consistent with  $\epsilon$ -greedy action selection.

$A_2=2, R_2=-1, Q_1=-1, \boxed{Q_2=\frac{1}{1}=1}, Q_3=0, Q_4=0.$   
 $\downarrow$  Choose 2. Maybe random. Cuz  $Q_2=1 > Q_3=Q_4=0 > Q_1=-1$  } Consistent with  $\epsilon$ -greedy action selection.

$A_3=2, R_3=-2, Q_1=-1, \boxed{Q_2=\frac{-1}{2}=-0.5}, Q_3=0, Q_4=0.$   
 $\downarrow$  Choose 2. Must be random. Cuz  $Q_3=Q_4=0 > Q_2=-0.5 > Q_1=-1$ .  
 According to  $\epsilon$ -greedy action selection, it should choose 3 or 4.

$A_4=2, R_4=2, Q_1=-1, Q_2=\frac{1}{3}=0.33, \boxed{Q_3=0}, Q_4=0.$   
 $\downarrow$  Choose 3. Must be random. Cuz  $Q_2=0.33 > Q_3=Q_4=0 > Q_1=-1$ .  
 According to  $\epsilon$ -greedy action selection, it should choose 2.

$A_5=3, R_5=0, Q_1=-1, Q_2=0.33, Q_3=\frac{0}{1}=0, Q_4=0.$

Thus, it may occurred on step 1, 2 and 3

It definitely occurred on step 4 and 5

Q2:

When  $\alpha_n$  is non-stationary.

$$Q_{n+1} = Q_n + \alpha_n [R_n - Q_n]$$

$$= \alpha_n R_n + (1 - \alpha_n) Q_n$$

$$= \alpha_n R_n + (1 - \alpha_n) [\alpha_{n-1} R_{n-1} + (1 - \alpha_{n-1}) Q_{n-1}]$$

$$= \alpha_n R_n + (1 - \alpha_n) \alpha_{n-1} R_{n-1} + (1 - \alpha_n) (1 - \alpha_{n-1}) [\alpha_{n-2} R_{n-2} + (1 - \alpha_{n-2}) Q_{n-2}]$$

...

$$= \alpha_n R_n + (1 - \alpha_n) \alpha_{n-1} R_{n-1} + \dots + Q_1 \prod_{i=1}^n (1 - \alpha_i)$$

$$= \alpha_n R_n + \sum_{i=1}^{n-1} \left[ \prod_{j=i+1}^n (1 - \alpha_j) \alpha_i R_i \right] + Q_1 \prod_{i=1}^n (1 - \alpha_i).$$

Q3:

(a) Equation 2.1:  $Q_n = \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$ ,  $E(R_n) = q^*$

$$E(Q_n) = E\left[\frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}\right]$$

$$= \frac{1}{n-1} \times E(R_1 + R_2 + \dots + R_{n-1})$$

$$= \frac{1}{n-1} \times (E(R_1) + E(R_2) + \dots + E(R_{n-1}))$$

$$= \frac{1}{n-1} \times (n-1)q^*$$

$$= q^*.$$

Thus, it is unbiased.

(b) If  $Q_1 = 0$  for  $n > 1$ :

$$Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i = \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i$$

$$E(Q_{n+1}) = E\left[\sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i\right] = \sum_{i=1}^n \alpha(1-\alpha)^{n-i} E(R_i)$$

$$= \sum_{i=1}^n \alpha(1-\alpha)^{n-i} q^*$$

Thus, only when  $\sum_{i=1}^n \alpha(1-\alpha)^{n-i} = 1$ , it is unbiased.



Q3:

$$(c) \quad Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

When  $E|Q_{n+1}| = q^*$ , it is unbiased.

When

$$\text{Thus, } E|Q_{n+1}| = E\left|(1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i\right|$$

$$= (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} E|R_i|$$

$$= (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} q^* = q^*, \text{ it is unbiased.}$$

The condition:

$$\text{Therefore, } (1-\alpha)^n Q_1 = 0, \rightarrow \underbrace{Q_1 = 0}, \underbrace{\sum_{i=1}^n \alpha (1-\alpha)^{n-i} = 1}.$$

$$(d) \quad Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

When  $n \rightarrow \infty$ . since  $0 < \alpha < 1$ .  $(1-\alpha)^n Q_1 \rightarrow 0$ .

$$\sum_{i=1}^n \alpha (1-\alpha)^{n-i} = 1 - (1-\alpha)^{n+1} \text{ and } (1-\alpha)^{n+1} \rightarrow 0. \text{ thus}$$

$$\text{thus, } \sum_{i=1}^n \alpha (1-\alpha)^{n-i} \rightarrow 1. \text{ then it satisfied the condition } \begin{cases} Q_1 = 0 \\ \sum_{i=1}^n \alpha (1-\alpha)^{n-i} = 1 \end{cases}$$

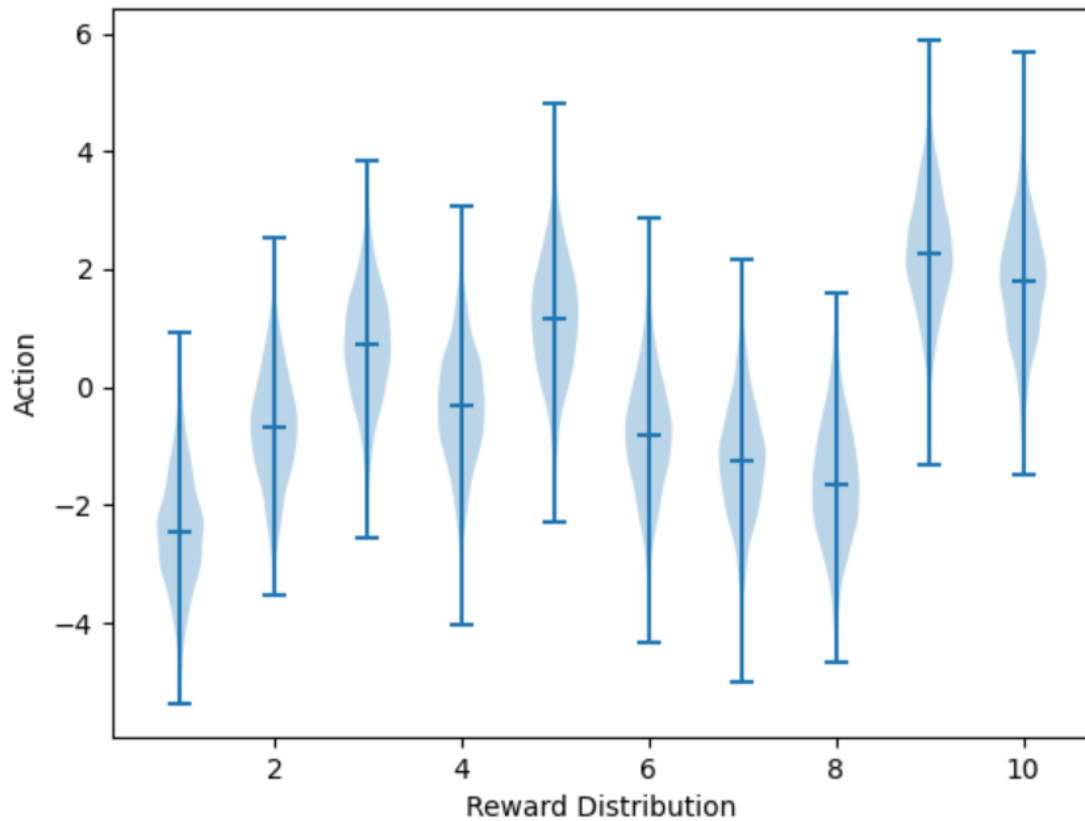
so it is unbiased when  $n \rightarrow \infty$ .

Q3:

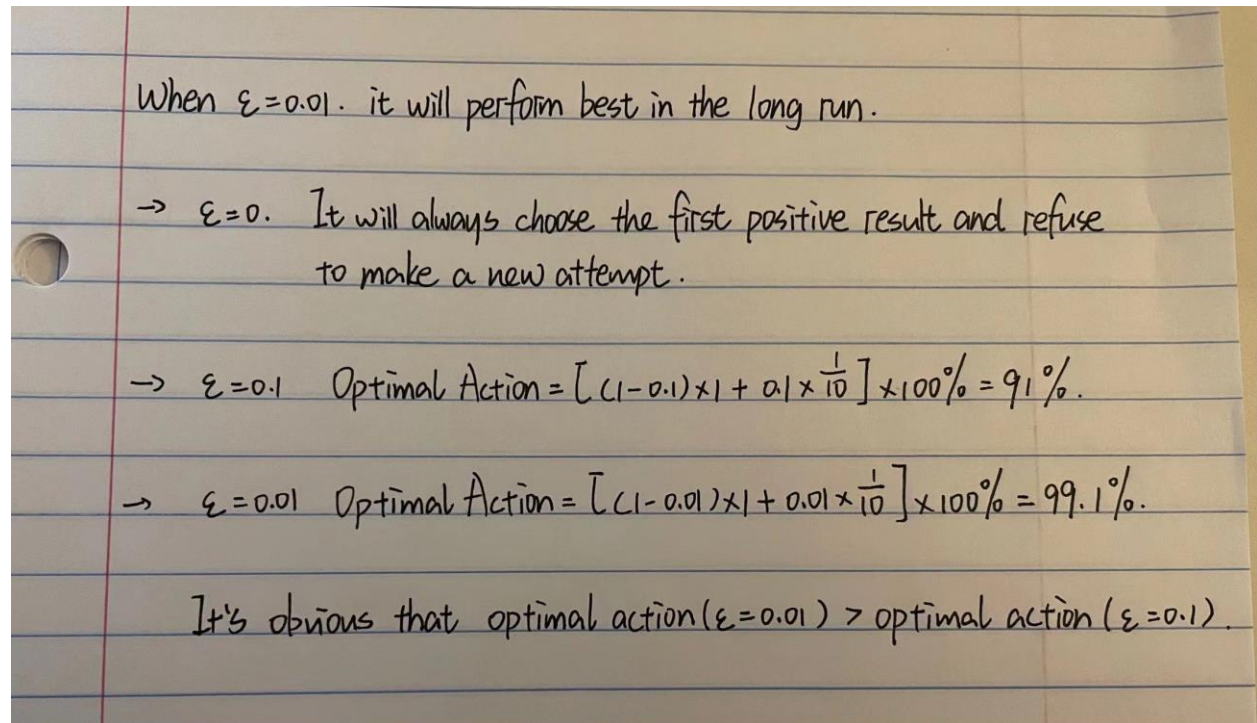
(e) **Why should we expect that the exponential recency-weighted average will be biased in general?**

1. *the actual situation will never be perfectly.*
2. *we can never obtain infinite values, but we can obtain values large enough that this still leads to the possibility of some bias*
3. *The occurrence of bias is more common than the absence of bias*

Q4:



Q5:



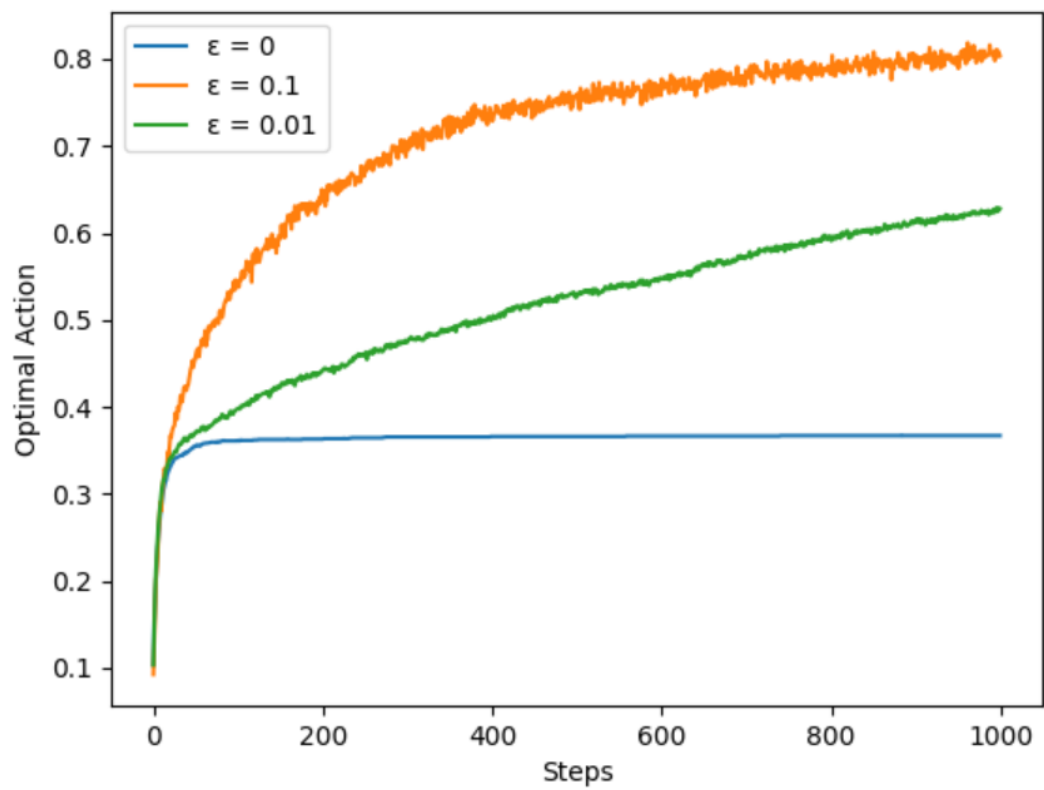
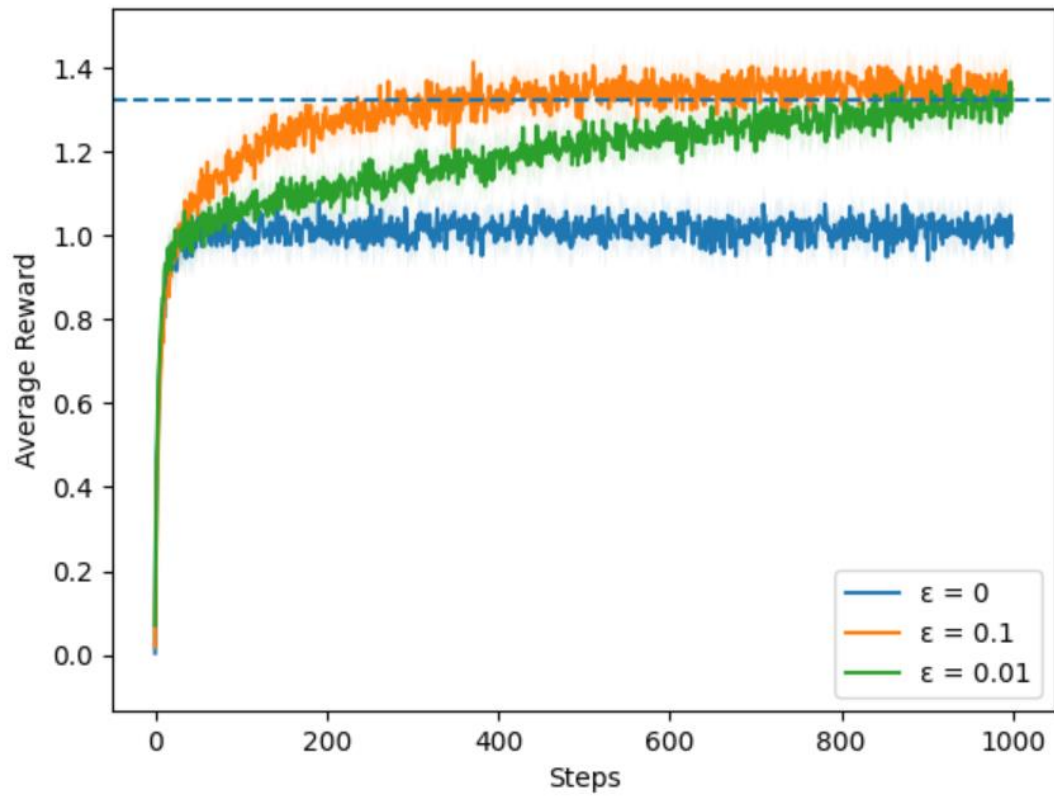
The optimal action ( $\epsilon = 0.01$ ) is 8.1% better than the optimal action ( $\epsilon = 0.1$ )

Q6:

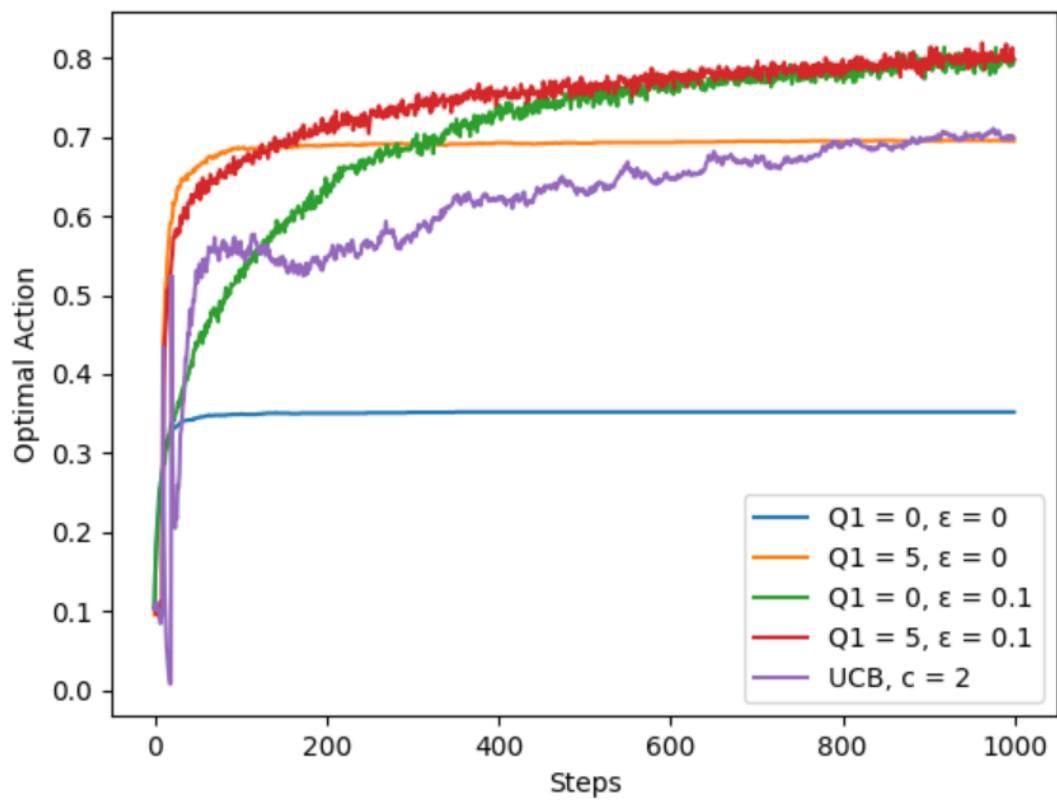
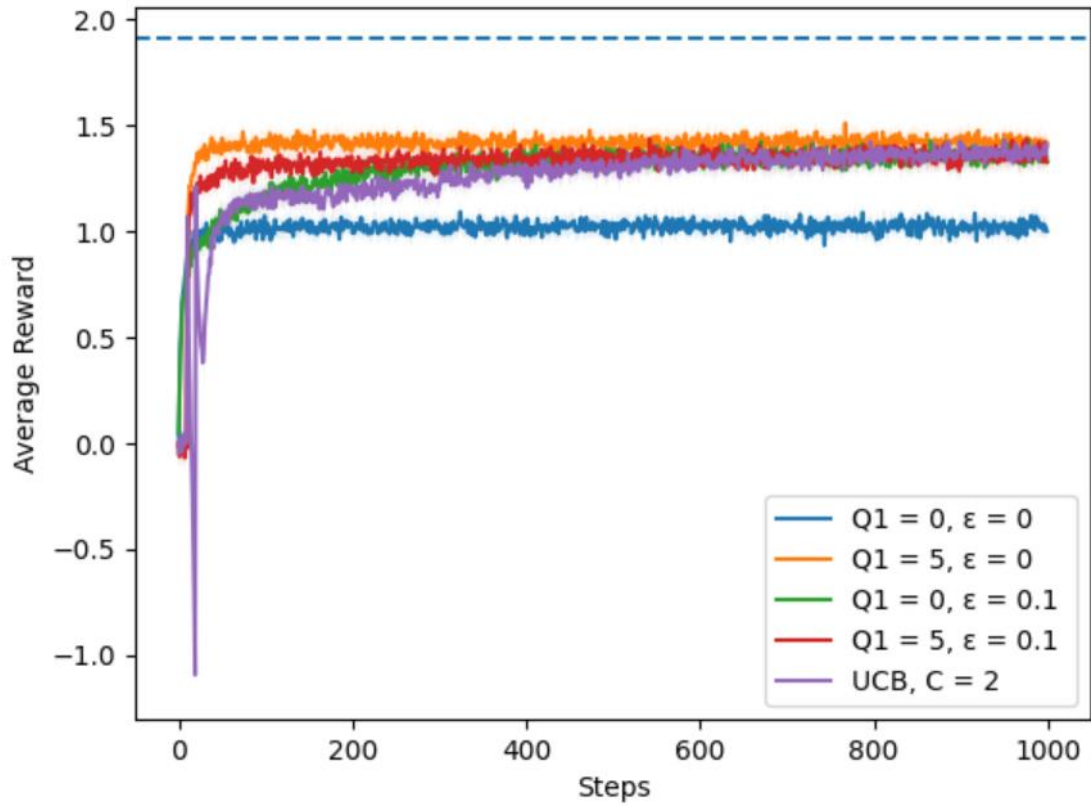
**Do the averages reach the asymptotic levels predicted in the previous question?**

*I don't think we have enough steps to get to the asymptotic levels predicted in the previous question. The more steps we have (close to infinity), the better the performance.*

Q6:



Q7:





Q7:

**Observe that both optimistic initialization and UCB produce spikes in the very beginning. In lecture, we made a conjecture about the reason these spikes appear. Explain in your own words why the spikes appear (both the sharp increase and sharp decrease). Analyze your experimental data to provide further empirical evidence for your reasoning.**

1. I believe that the optimistic initialization products spike is due to the fact that a high Q value was given at the beginning and that a relatively stable result could not be obtained immediately after further exploration.
2. For the UCB, I think it is due to the fact that it first needs to try all the actions to find the one that is likely to obtain a high reward, thus generating the spike, and after 10 actions he will adjust in order to obtain a more stable performance, and that products another spike.