

$$1. \quad R_0 = 0 \quad \underline{Q_1(1)} = Q_1(2) = Q_1(3) = Q_1(4) = 0$$

$t_1 \quad A_1=1$ Since $Q_1(a)$ was same, ϵ -greedy may occurred

$$R_1 = -1 \quad \underline{Q_2(1)} = Q_2(3) = Q_2(4) = 0$$

$t_2 \quad A_2=2$ Since $Q_2(2) = Q_2(3) = Q_2(4)$, ϵ -greedy may occurred

$$R_2 = 1 \quad \underline{Q_3(1)} = 1, Q_3(2) = 0, Q_3(3) = 0$$

$t_3 \quad A_3=2$ Since Action 1 was not selected, ϵ -greedy may occurred.

$$R_3 = -2 \quad \underline{Q_4(1)} = -1 + \frac{1}{2}(-2-1) = -\frac{1}{2}, Q_4(2) = 0, Q_4(3) = 0$$

$t_4 \quad A_4=2$ Since $Q_4(2) < Q_4(3) = Q_4(4)$, ϵ -greedy definitely occurred.

$$R_4 = 2 \quad Q_5(1) = -1, Q_5(2) = -\frac{1}{2} + \frac{1}{3}(2+2) = \frac{1}{3}, \underline{Q_5(3)} = 0, Q_5(4) = 0$$

$t_5 \quad A_5=3$ Since $Q_5(3) < Q_5(2)$, ϵ -greedy definitely occurred.

In sum, the ϵ -greedy may occurred on time step 1, 2, 3

the ϵ -greedy definitely occurred on time step 4, 5

$$2. \quad Q_{n+1} = Q_n + d_n [R_n - Q_n]$$

$$= d_n R_n + (1-d_n) Q_n$$

$$= d_n R_n + (1-d_n) [Q_{n-1} + d_{n-1} [R_{n-1} - Q_{n-1}]]$$

$$= d_n R_n + (1-d_n) [d_{n-1} R_{n-1} + (1-d_{n-1}) Q_{n-1}]$$

$$= d_n R_n + (1-d_n) d_{n-1} R_{n-1} + (1-d_n)(1-d_{n-1}) Q_{n-1}$$

$$= d_n R_n + (1-d_n) d_{n-1} R_{n-1} + (1-d_n)(1-d_{n-1}) [Q_{n-2} + d_{n-2} [R_{n-2} - Q_{n-2}]]$$

$$= d_n R_n + (1-d_n) d_{n-1} R_{n-1} + (1-d_n)(1-d_{n-1}) d_{n-2} R_{n-2} + (1-d_n)(1-d_{n-1})(1-d_{n-2}) Q_{n-2}$$

$$= d_n R_n + \sum_{i=1}^{n-1} (d_i R_i \prod_{j=i+1}^n (1-d_j)) + Q_1 \prod_{i=1}^n (1-d_i)$$

$$3. (a) Q_n = \frac{R_1 + \dots + R_{n-1}}{n-1} = \bar{g}^* \text{ when } n \text{ is large enough}$$

for E.S. 2.1

$$\begin{aligned} E(Q_n) &= E\left(\frac{R_1 + \dots + R_{n-1}}{n-1}\right) \\ &= \frac{1}{n-1} E(R_1 + \dots + R_{n-1}) \\ &= \frac{1}{n-1} E(R_1) + \dots + E(R_{n-1}) \end{aligned}$$

$$\therefore E(R_n) = \bar{g}^*$$

$$\therefore E(Q_n) = \frac{1}{n-1} (n-1) \bar{g}^* = \bar{g}^*$$

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

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

$$\therefore Q_1 = 0$$

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

$$\begin{aligned} \therefore E(Q_{n+1}) &= \sum_{i=1}^n d(1-d)^{n-i} E(R_i) \\ &= \sum_{i=1}^n d(1-d)^{n-i} \bar{g}^* \end{aligned}$$

$$\text{If } \sum_{i=1}^n d(1-d)^{n-i} = 1 \quad E(Q_{n+1}) = \bar{g}^* \text{ (unbiased)}$$

$$\text{else } E(Q_{n+1}) \neq \bar{g}^* \text{ (biased)}$$

(c)

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

$$\boxed{E(Q_{n+1}) = g^*}$$

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

If the $E(Q_{n+1})$ is unbiased,

$$E(Q_{n+1}) = g^*$$

$$\therefore (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} g^* = g^*$$

\therefore Conditions:

$$(1-\alpha)^n Q_1 = 0 \rightarrow \boxed{Q_1 = 0}$$

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

$$1 - (1-\alpha)^{n-1}$$

$$(d) \text{ Due to } (1-\alpha)^n + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} = 1$$

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

when $n \rightarrow \infty$

$$\underline{(1-\alpha)^n \rightarrow 0}$$

$$\underline{1 - (1-\alpha)^n \rightarrow 1}$$

$$\therefore (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} g^* = g^*$$

that is:

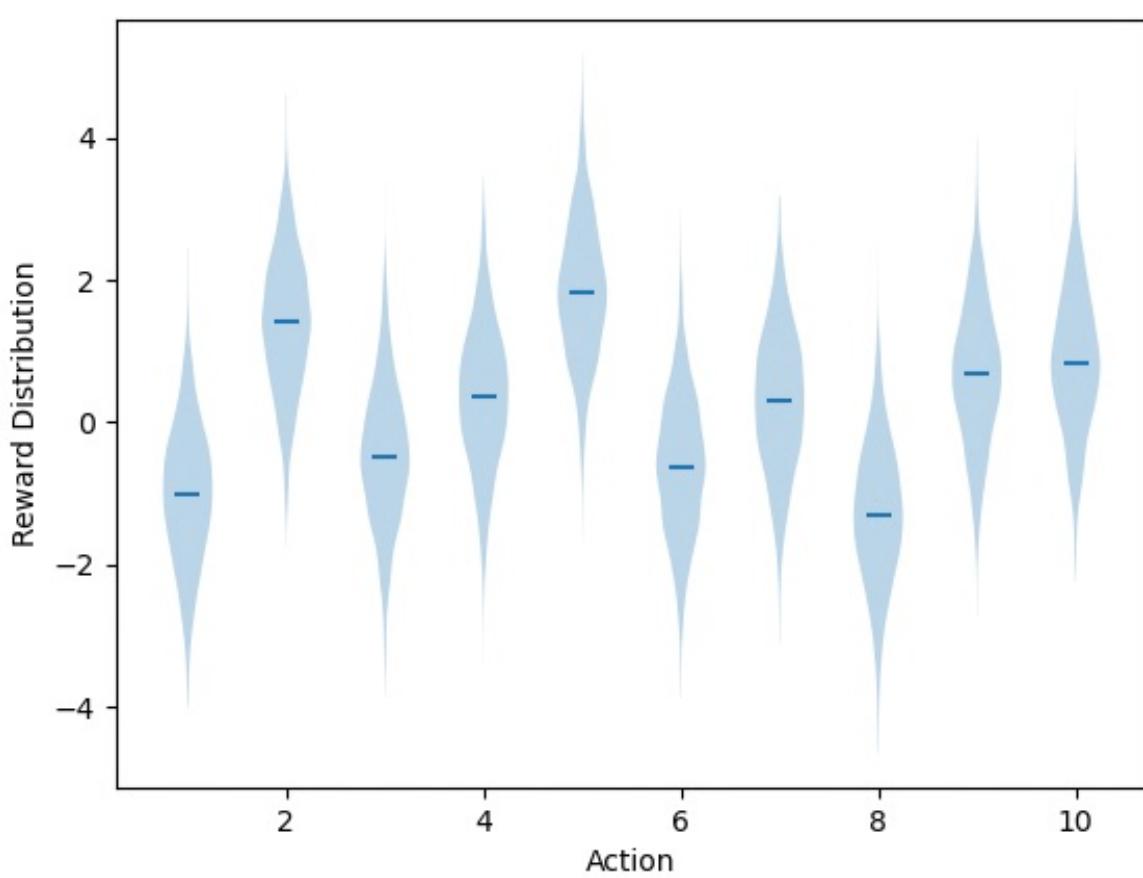
$$0 \cdot Q_1 + 1 \cdot g^* = g^*$$

$$\text{as: } \underline{g^* = g^*}$$

Hence, Q_n is asymptotically

(e) First, We cannot acquired infinite samples, which means that n is not equal to ∞ . As a result, $\sum_{t=1}^n \alpha(t\alpha)^{n-t} \neq 1$, so it must have some bias in general. Besides, $Q_n = \frac{R_{t+1} + R_{t+1}}{n-1} = g^*$ when n is large enough (finite). But if n is infinite, $Q_n \neq g^*$, which also denote some bias in the equation. All in all, we should expect that the exponential recency-weighted average will be biased in general.

4. Plot



5.

$\epsilon = 0.01$ will perform best in the long run.

When $\epsilon = 0.01$,

$$\text{Optimal Action} = |-0.0| + 0.01 \times \frac{1}{10} = 0.991 = 99.1\%$$

When $\epsilon = 0.1$,

$$\text{Optimal Action} = |-0.1| + 0.1 \times \frac{1}{10} = 0.91 = 91\%$$

When $\epsilon = 0$,

It will always select the highest reward of current state without any exploration. In the long run, it will lead to low probability of selecting optimal action.

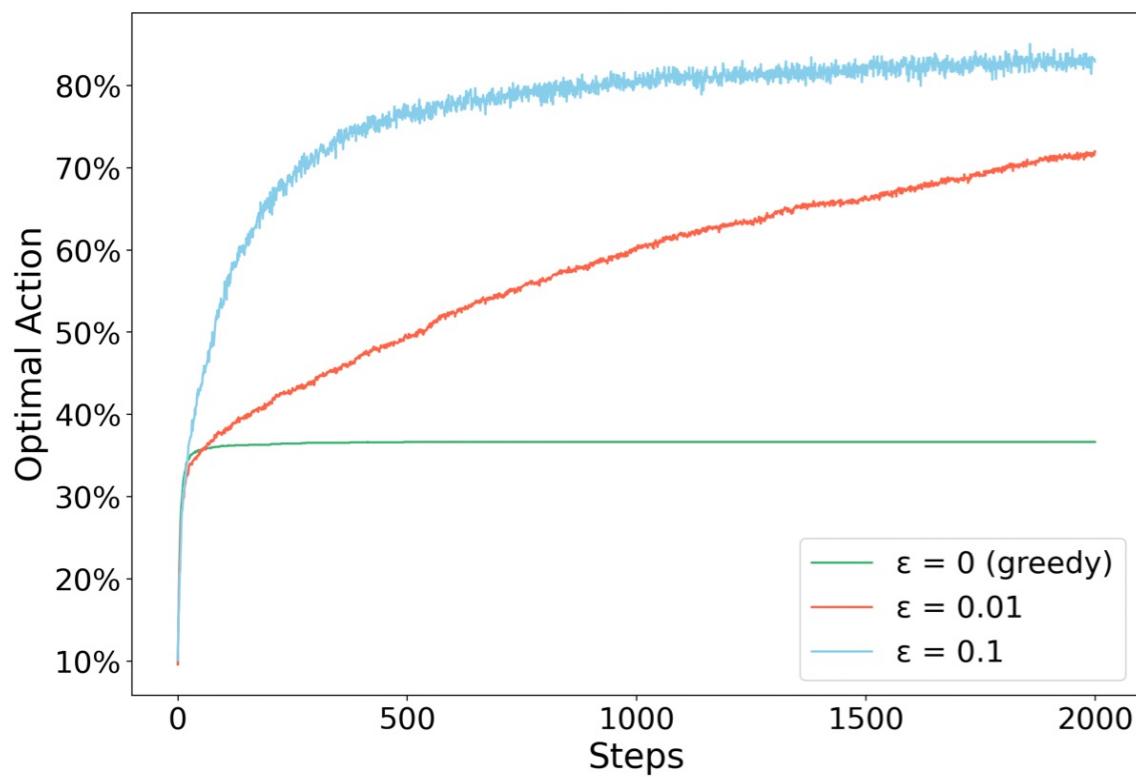
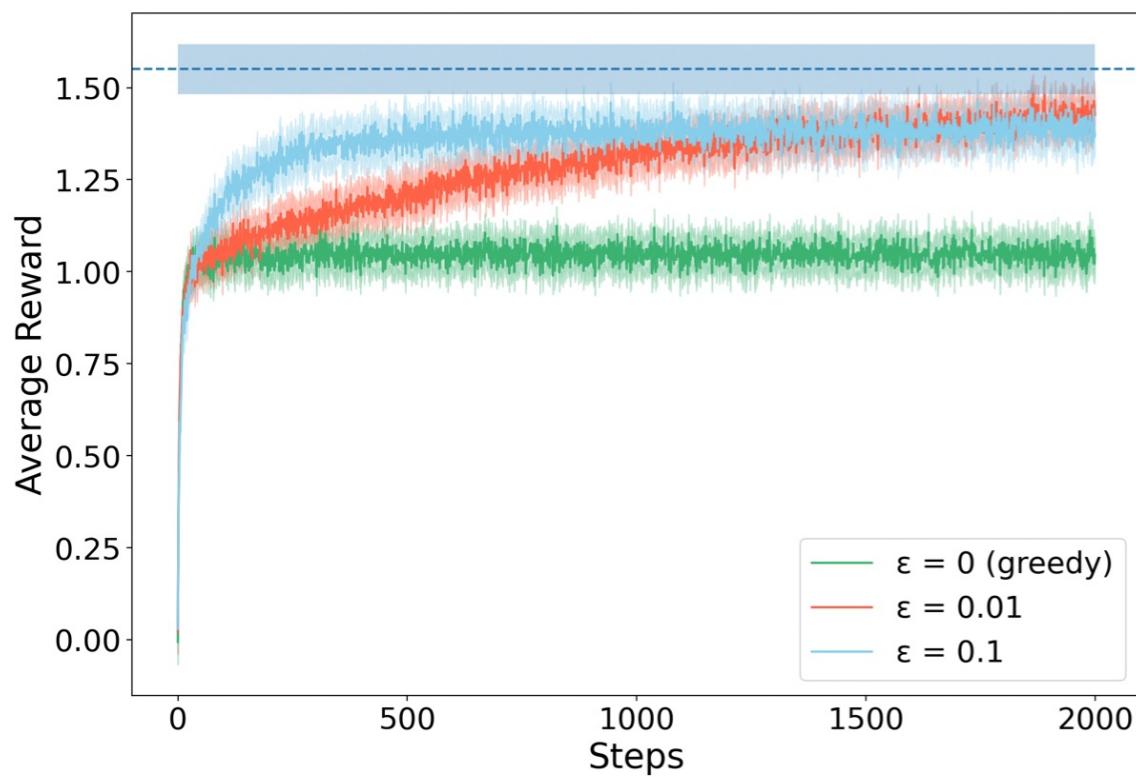
Quantitatively:

$\epsilon = 0.01$ is better than $\epsilon = 0.1$ about 8.1 %

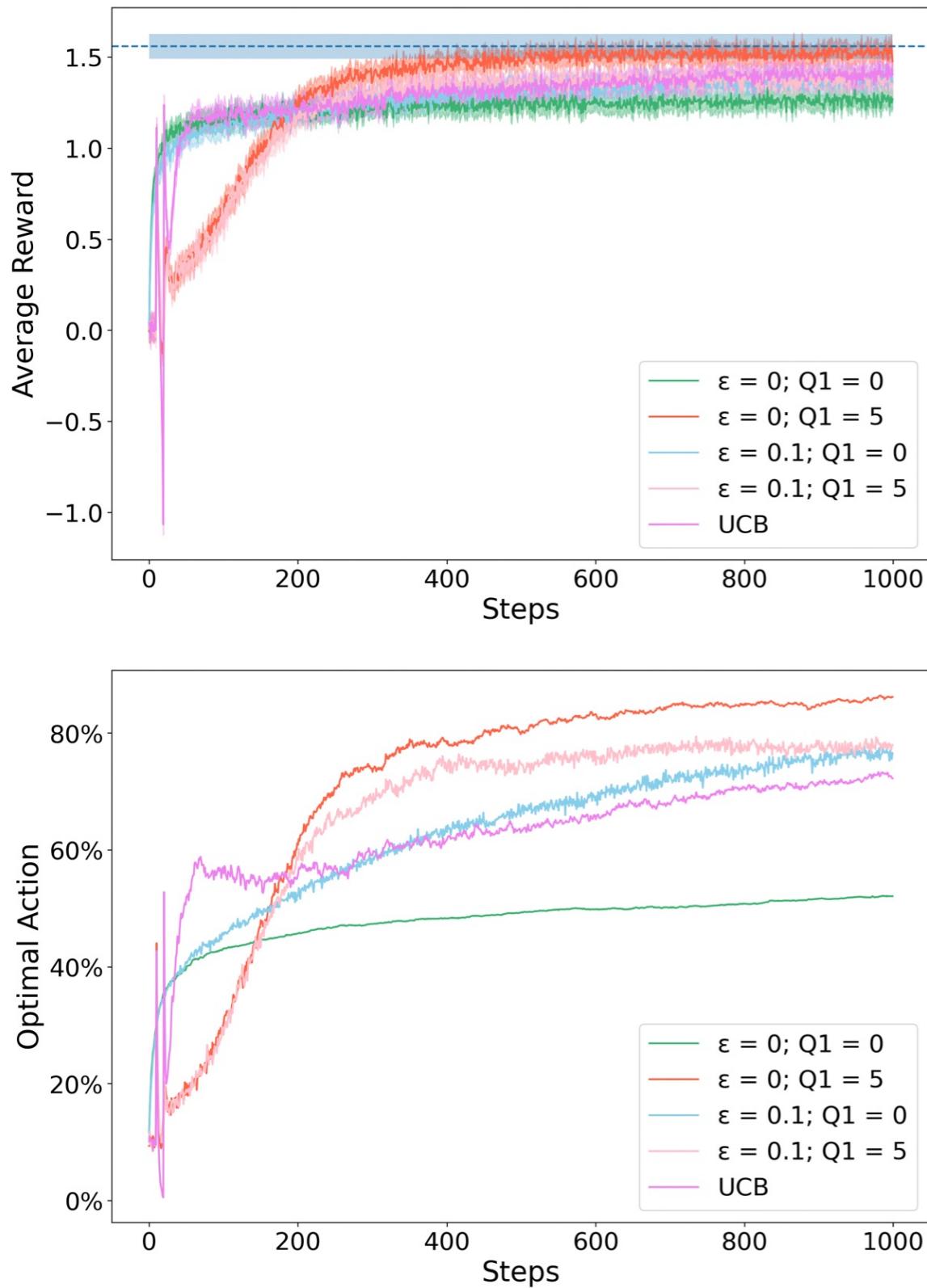
b. Written

From my observation (as shown below), the averages don't reach the asymptotic level. I think the reason is that the steps are still not enough. If we have much more steps, even approaching infinity, the performance and the average reward will reach the asymptotic level.

6. Plot



7. Plot



7. Written

I think the reason of spikes is because the exploration is not enough. Also, the past experience cannot support the agent to select the optimal actions. With the continuous exploration process, the agent can gradually select the optimal actions. During this process, the average rewards are increased/ rapidly. That's why the spikes appear.
decreased

so the agent will randomly select actions at begining.

From the VCB in my experimental data, we could find the spikes (both the sharp increase and decrease). The reason is that the VCB first needs to try all actions. Some actions have high reward, some are not, that's why the spikes appear. The inner property of the VCB (explore all actions at the begining) leads to the spikes.

Since the rewards are initialized randomly according to the normal distribution.