

Reinforcement learning
Assignment – 1
Gopisainath Mamindlapalli
Study partner: Harin Kumar

- 1) Given $Q_1 = [0 \ 0 \ 0 \ 0]$, actions : $A_1 = 1, A_2 = 2, A_3 = 2, A_4 = 2, A_5 = 3$ and the rewards are: $R_1 = -1, R_2 = 1, R_3 = -2, R_4 = 2, R_5 = 0$.

Using the rewards the following Q values are:

$$\begin{aligned} Q_1(A_1) &= -1/1 = -1 & : Q_1 &= [-1 \ 0 \ 0 \ 0] \\ Q_2(A_2) &= 1/1 = 1 & : Q_2 &= [-1 \ 1 \ 0 \ 0] \\ Q_3(A_3) &= (1-2)/2 = -1/2 & : Q_2 &= [-1 \ 1/2 \ 0 \ 0] \\ Q_4(A_4) &= (1-2+2)/3 = 1/3 & : Q_2 &= [-1 \ 1/3 \ 0 \ 0] \\ Q_5(A_5) &= 0/1 = 0 & : Q_2 &= [-1 \ 1/3 \ 0 \ 0] \end{aligned}$$

From the above Q values we can say that actions 4,5 are occurred randomly also when there is no action is taken then Q values will be same for all actions (zero in this case) then this event may occur.

2)

$$\begin{aligned} Q_{n+1} &= Q_n + \alpha_n [R_n - Q_n] \\ Q_{n+1} &= \alpha_n R_n + (1 - \alpha_n) Q_n \end{aligned}$$

where,

$$\begin{aligned} Q_n &= \alpha_{n-1} R_{n-1} + (1 - \alpha_{n-1}) Q_{n-1} \\ Q_{n-1} &= \alpha_{n-2} R_{n-2} + (1 - \alpha_{n-2}) Q_{n-2} \end{aligned}$$

$$Q_{n+1} = Q_1 \prod_{i=1}^n (1 - \alpha_i) + \sum_{i=1}^n \alpha_i \prod_{j=1+1}^n (1 - \alpha_j) R_i$$

- 3) a) The sample-average estimate in Equation 2.1, which is used to calculate the action value Q_n , is demonstrated to be unbiased. In this estimate, Q_n is determined by dividing the total rewards from time step 1 to time step $n-1$ by the number of time steps $(n-1)$. Comparing the predicted value of Q_n , designated as $E[Q_n]$, with the actual expected reward q of a particular arm is the main focus of bias assessment. The estimate is thought to be unbiased when $E[Q_n]$ equals q , and in this case the law of large numbers ensures that it is. $E[Q_n]$ converges to q as the number of samples (n) approaches infinity, ensuring that the estimate improves with time and eventually converges to the true value.

b) The estimate Q_n for $n > 1$ is biased in the exponential recency-weighted average estimate specified in Equation 2.5 with $0 < \alpha < 1$ and a starting estimate $Q_1 = 0$. Because α is strictly between 0 and 1, it gives the most recent reward R_n less weight than the original estimate $Q_1 = 0$, which leads to this bias. As a result, there is a persistent bias in the estimate as n grows since the estimate Q_n fails to converge to the genuine expected reward q and is still strongly influenced by the initial estimate.

c) We must make sure that Q_n converges to the genuine expected reward q over time in order to determine the circumstances under which the exponential recency-weighted average estimate Q_n will be unbiased. In other words, we're looking for circumstances in which: $E[Q_n] = q^*$. Remember that $Q_n = (1 - \alpha) * Q_{n-1} + \alpha * R_n$ is the update equation for Q_n in the exponential recency-weighted average estimate.

To make Q_n unbiased,

1. $E[R_n] = q$: The anticipated benefit at each time step must be equal to the actual anticipated benefit, q . To put it another way, the rewards ought to be a fair estimation of the real rewards.

2. $E[Q_{n-1}] = q$: The anticipated worth of the earlier estimate Q_{n-1} ought to be equivalent to the actual anticipated reward q . This requirement makes sure that as n increases, the initial estimate loses its impact. Both of these requirements must be met for Q_n to be impartial and to eventually converge to the genuine expected reward q . These requirements guarantee that the rewards and the earlier estimate both offer accurate information about the actual expected reward.

d) The expected value of Q_n must converge to the actual expected reward q in order to show that Q_n is asymptotically unbiased (i.e., $n \rightarrow \infty$).

$$Q_n = (1 - \alpha) * Q_{n-1} + \alpha * R_n$$

$$E[Q_n] = E[(1 - \alpha) * Q_{n-1} + \alpha * R_n]$$

$$E[Q_n] = (1 - \alpha) * E[Q_{n-1}] + \alpha * E[R_n]$$

$$\lim_{n \rightarrow \infty} E[Q_n] = \lim_{n \rightarrow \infty} [(1 - \alpha) * E[Q_{n-1}] + \alpha * E[R_n]]$$

$$\lim_{n \rightarrow \infty} E[Q_n] = (1 - \alpha) * \lim_{n \rightarrow \infty} E[Q_{n-1}] + \alpha * q^*$$

Now, consider the limit of $E[Q_{n-1}]$ as n approaches infinity:

$$\lim_{n \rightarrow \infty} E[Q_{n-1}] = \lim_{n \rightarrow \infty} E[Q_n]$$

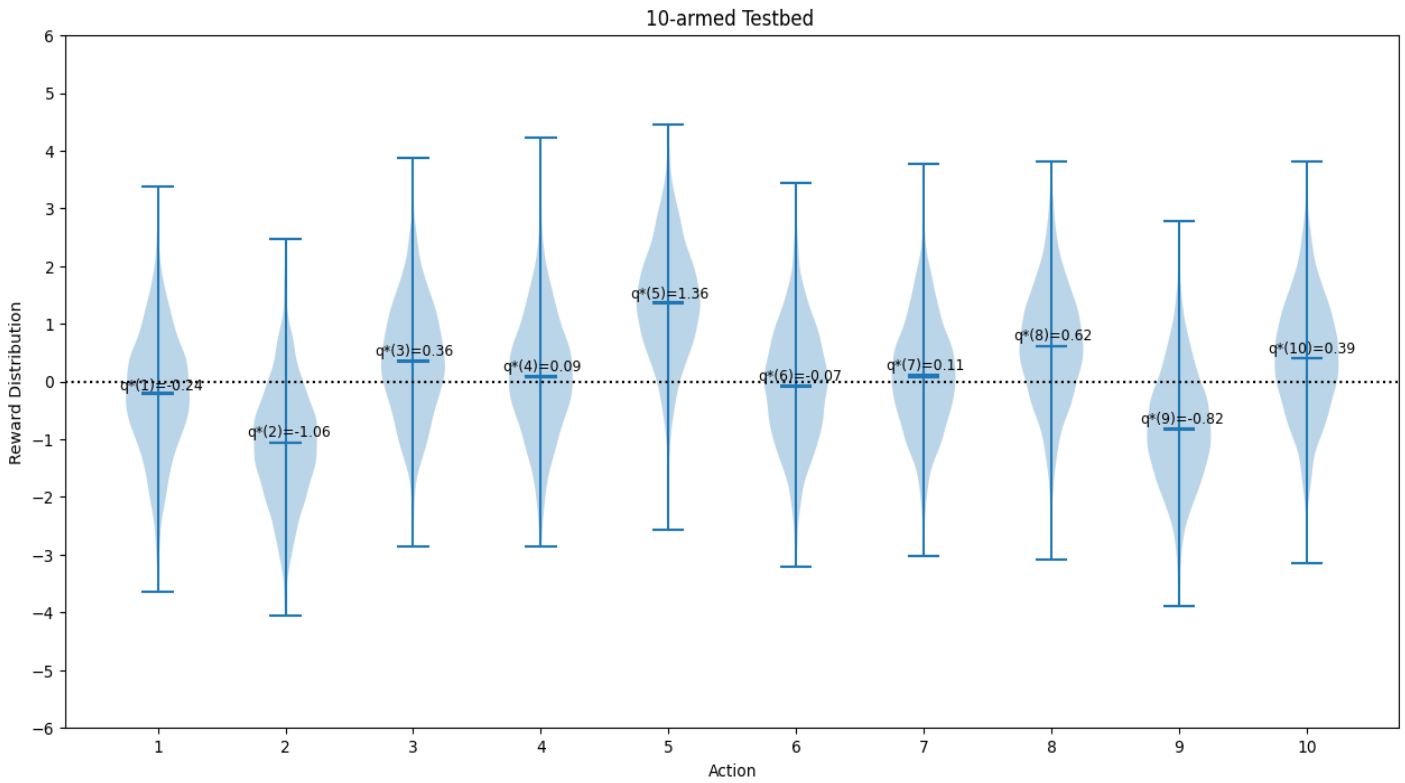
$$\text{Then, } \lim_{n \rightarrow \infty} E[Q_n] = (1 - \alpha) * \lim_{n \rightarrow \infty} E[Q_n] + \alpha * q^*$$

$$\lim_{n \rightarrow \infty} E[Q_n] = q^*$$

e) When employing the exponential recency-weighted average estimate, the selection of the weighting parameter is crucial. Although the exponential recency-weighted average is not by nature biased, the choice of α is critical to how well it performs. If α is set too closely to 1, the algorithm may respond to environmental changes more slowly since it would place a heavy emphasis on prior experiences, thus missing out on important information from more recent rewards.

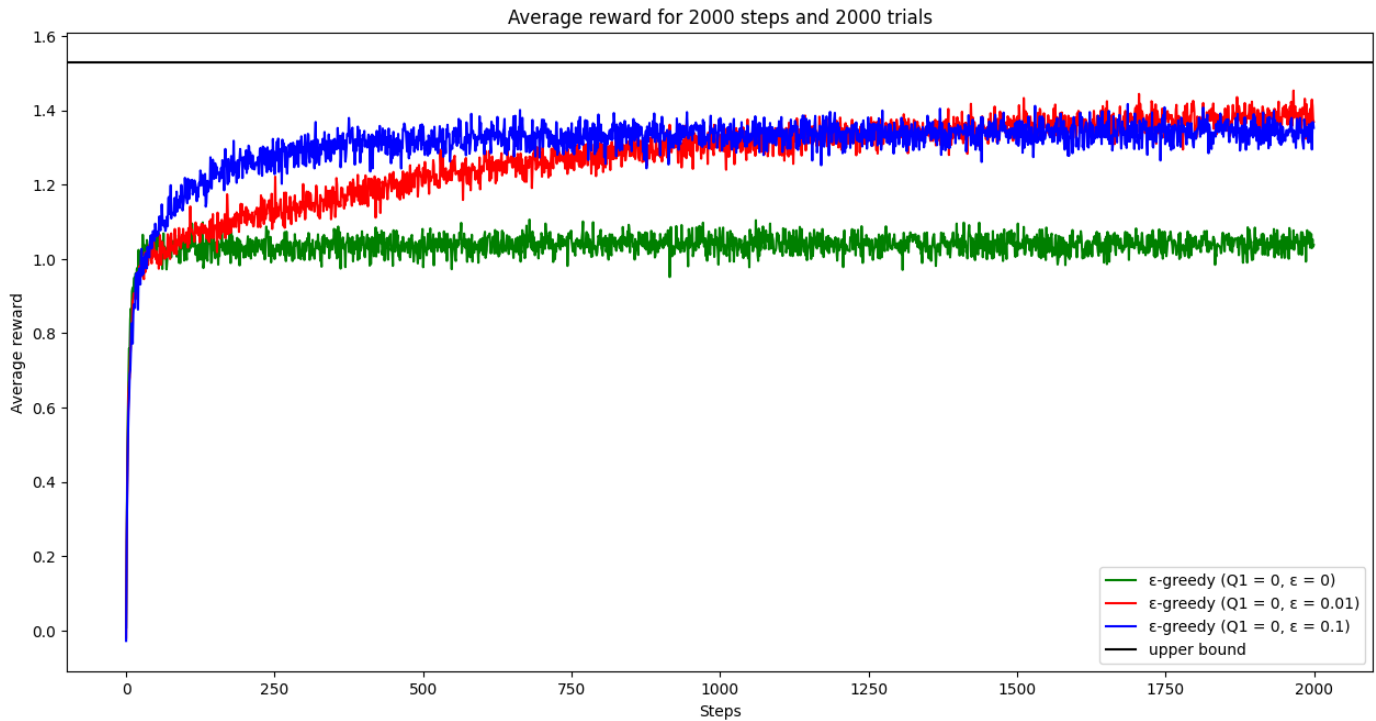
Alternatively, if α is set to be too small, the algorithm might disproportionately emphasise recent rewards, producing estimates with more variance. For the algorithm to successfully navigate different learning contexts, providing both adaptability and accuracy in forecasting the expected rewards, it is crucial to strike the proper balance while choosing a suitable α . This point emphasises the significance of tailoring solutions to the particulars of the current issue.

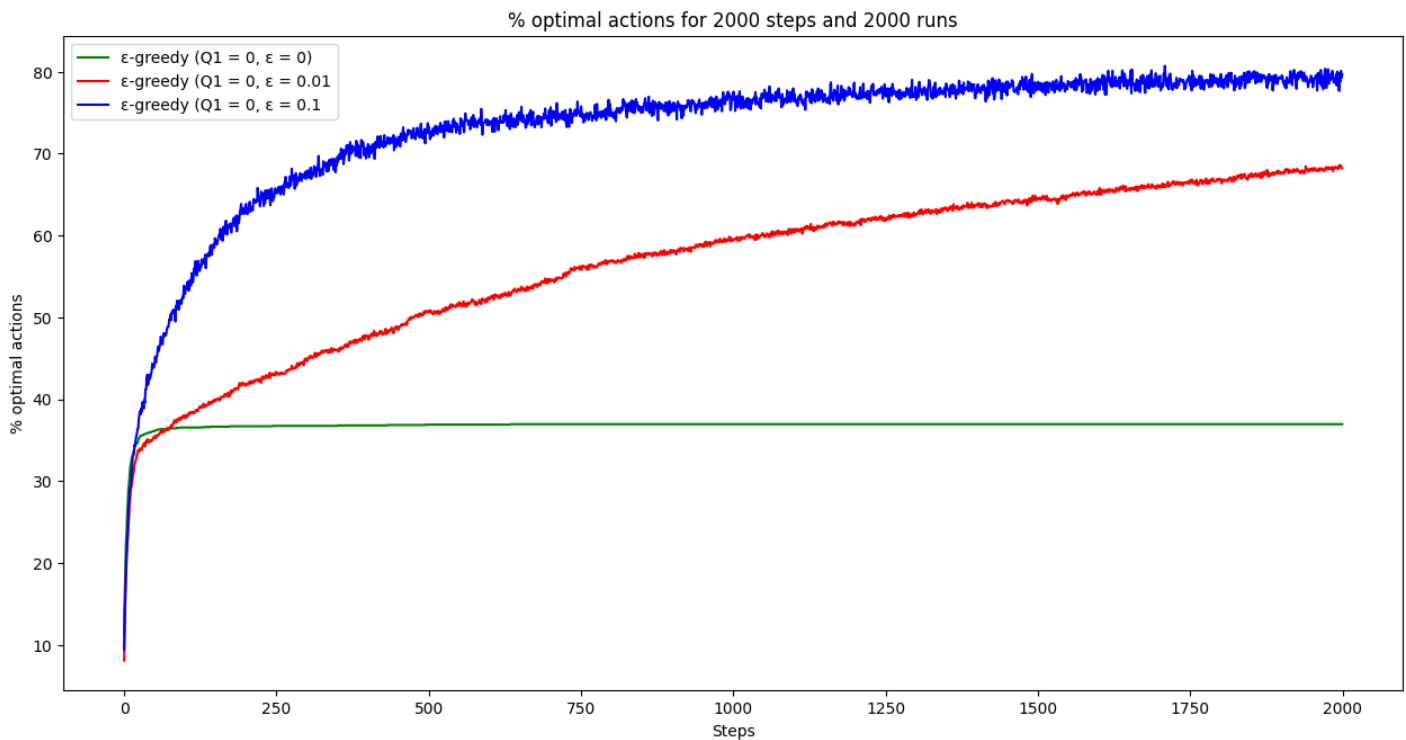
4)



- 5) The epsilon greedy method with epsilon value of 0.01 will perform in the long run. The probability of selecting the optimal action is $1 - \epsilon + (\epsilon/10)$.
 Therefore, for $\epsilon = 0.01$ method this value is 0.991
 for $\epsilon = 0.1$ method this value is 0.91
 so, it will perform by $(0.991/0.91) = 1.089$ times better

6)

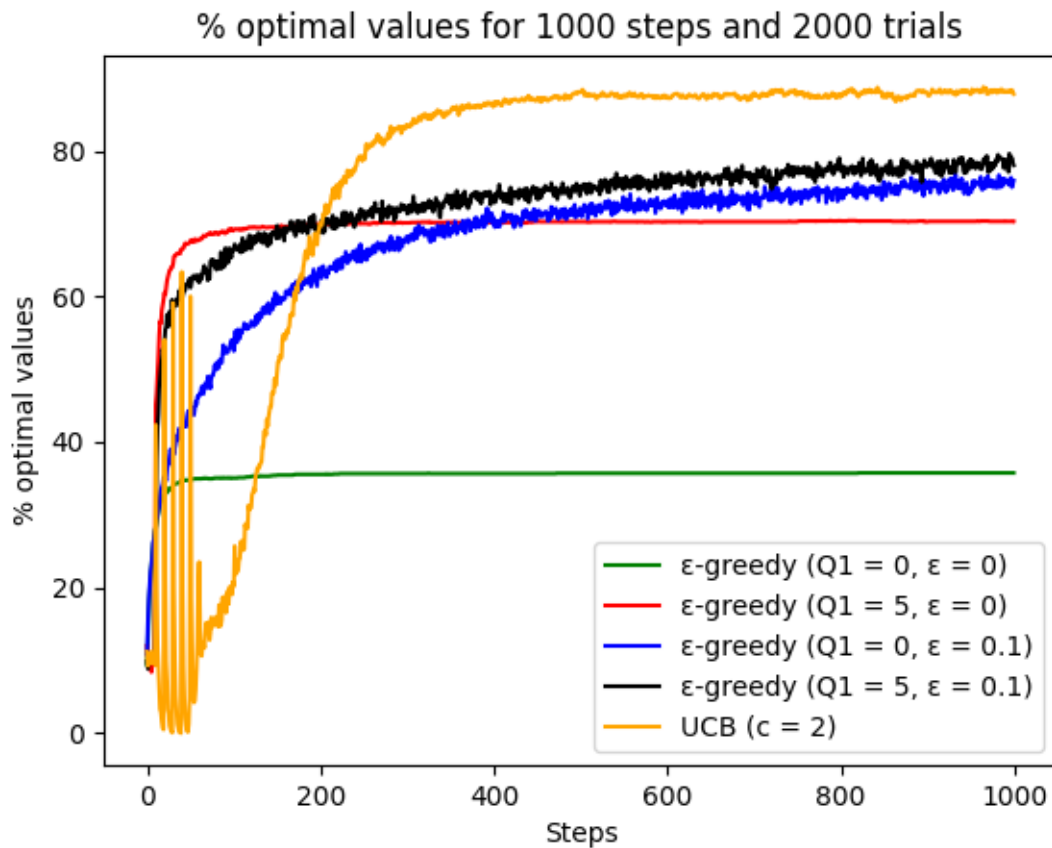




The average reward over steps graph shows that the epsilon=0.01 agent attempts to cross the Upper Bound line. So, as the number of steps tends to infinity, we can claim that it reaches the asymptotic levels.

7)





The initial spikes observed in both optimistic initialization and UCB techniques can be attributable to the deliberate use of high starting estimations for action values (Q-values). As can be observed from the average reward vs steps graph from steps 0 to 100, optimistic initialization with inflated Q-values promotes early exploration because the agent initially favours actions with greater anticipated rewards. As a result, rewards increase quickly at first as different activities are investigated. As can be seen in the average reward vs. steps graph from steps 0 to 150, UCB, on the other hand, promotes exploration by balancing the selection of actions with high uncertainty. This causes an early spike in rewards as uncertain actions with potential high rewards are taken.

The initial spikes observed in both optimistic initialization and UCB techniques can be attributable to the deliberate use of high starting estimations for action values (Q-values). As can be observed from the average reward vs steps graph from steps 0 to 100, optimistic initialization with inflated Q-values promotes early exploration because the agent initially favours actions with greater anticipated rewards. As a result, rewards increase quickly at first as different activities are investigated. As can be seen in the average reward vs. steps graph from steps 0 to 150, UCB, on the other hand, promotes exploration by balancing the selection of actions with high uncertainty. This causes an early spike in rewards as uncertain actions with potential high rewards are taken.