

Homework-1

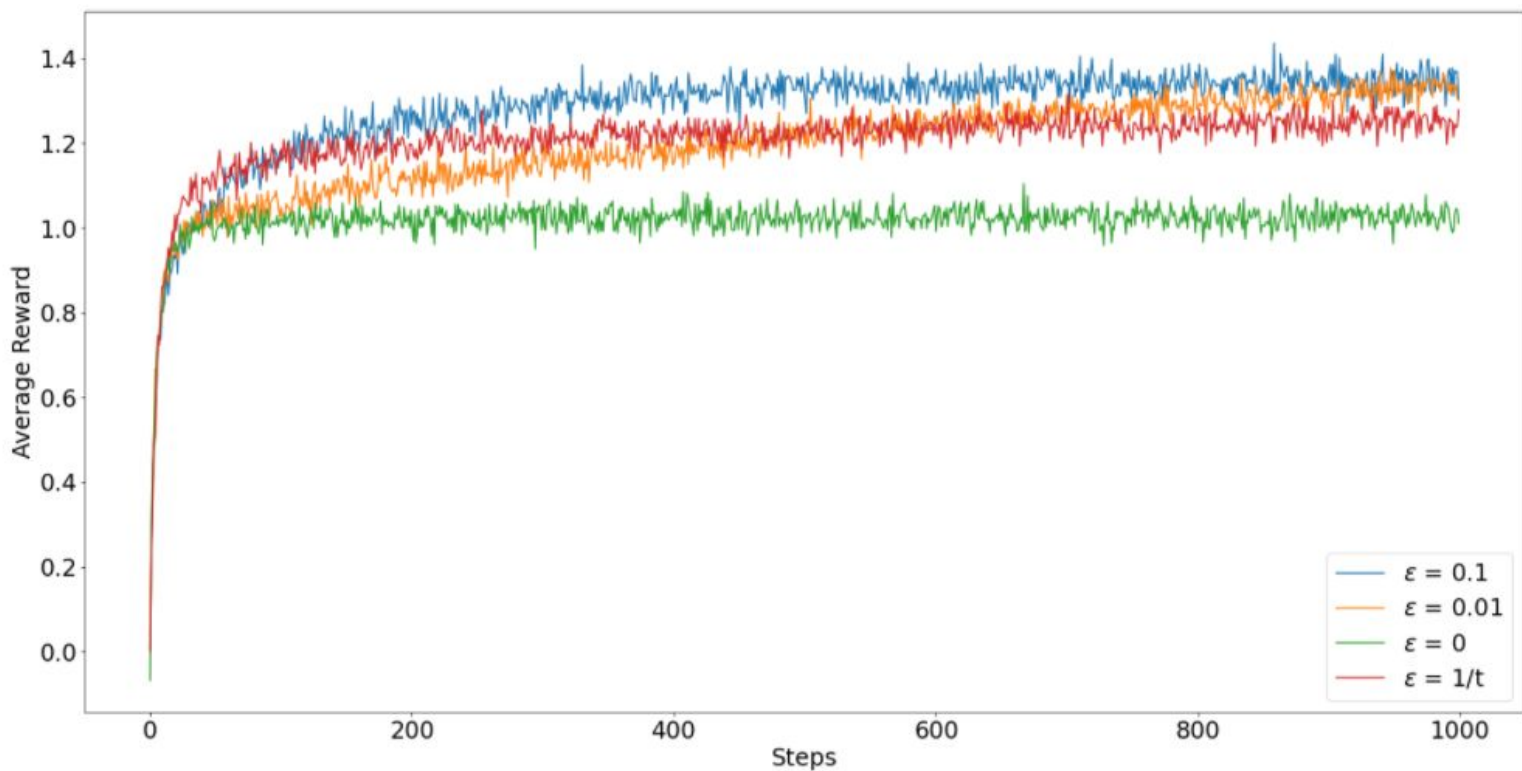
Question 1:

The 10-armed testbed with $q^*(a)$, $a=1,2,\dots,10$ selected according to a normal distribution with mean 0 and variance 1.

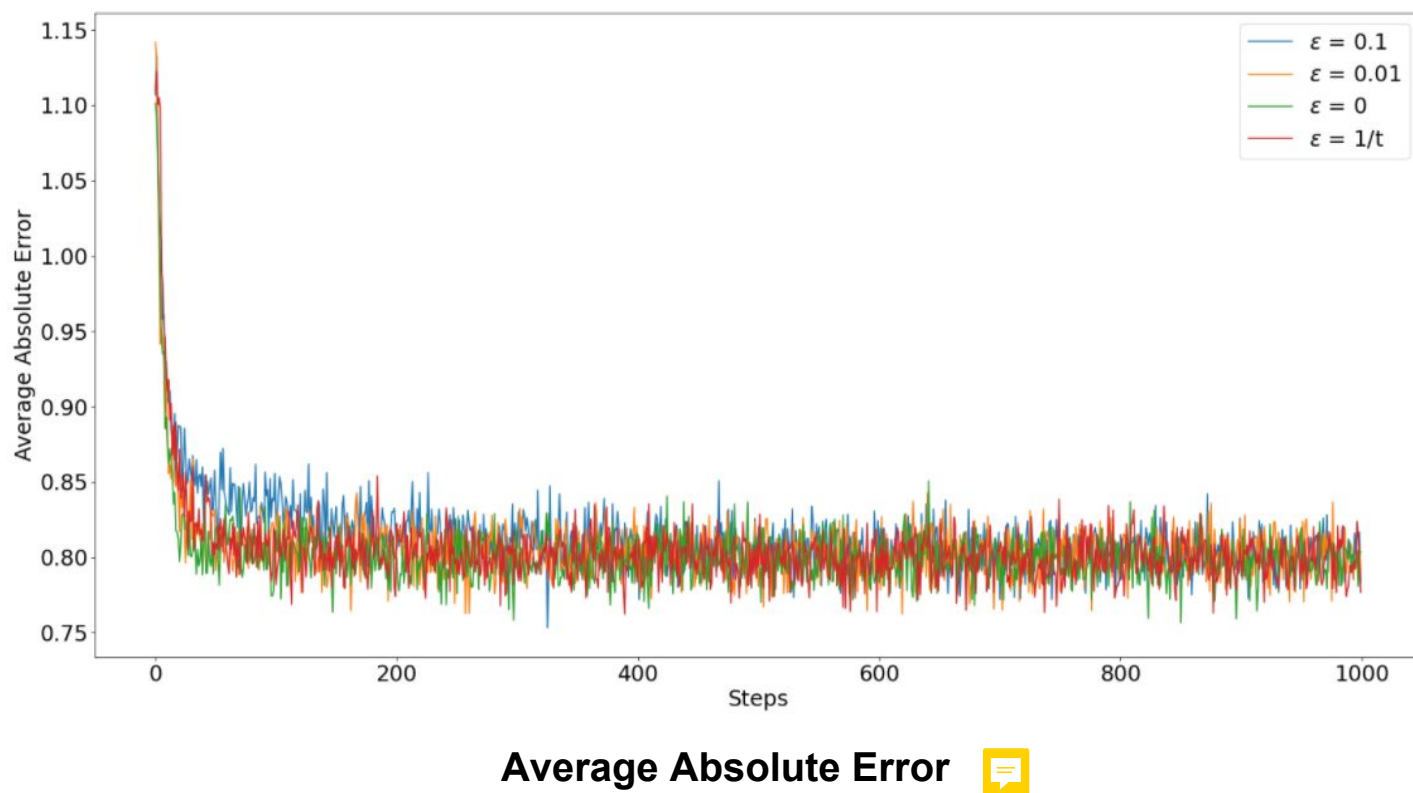
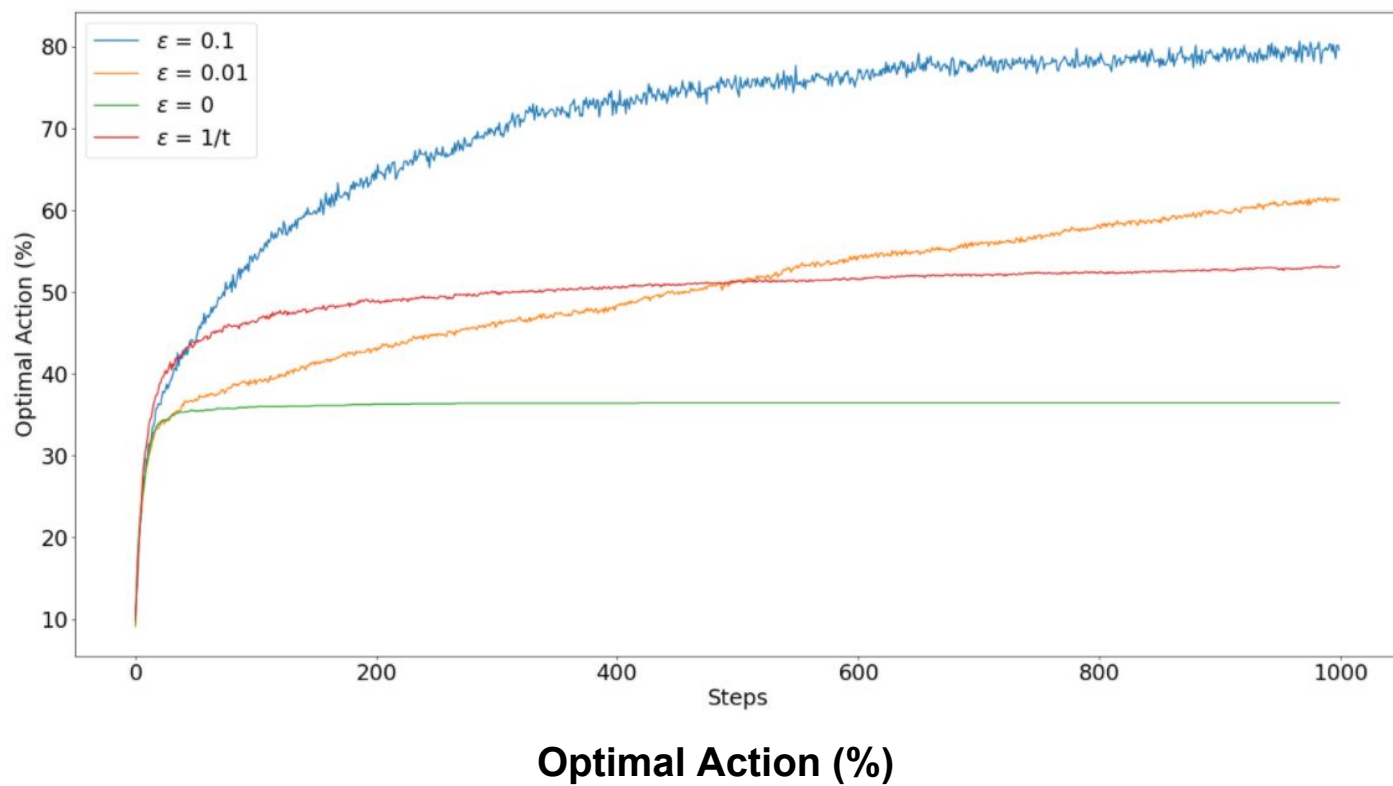
Rewards for each arm are selected using normal distribution with mean $q^*(A_t)$ and variance 1.

Graph drawn for $\varepsilon = 0.1, 0.01, 0$ and $1/t$.

Time steps=1000, Number of simulations=2000



Average Reward

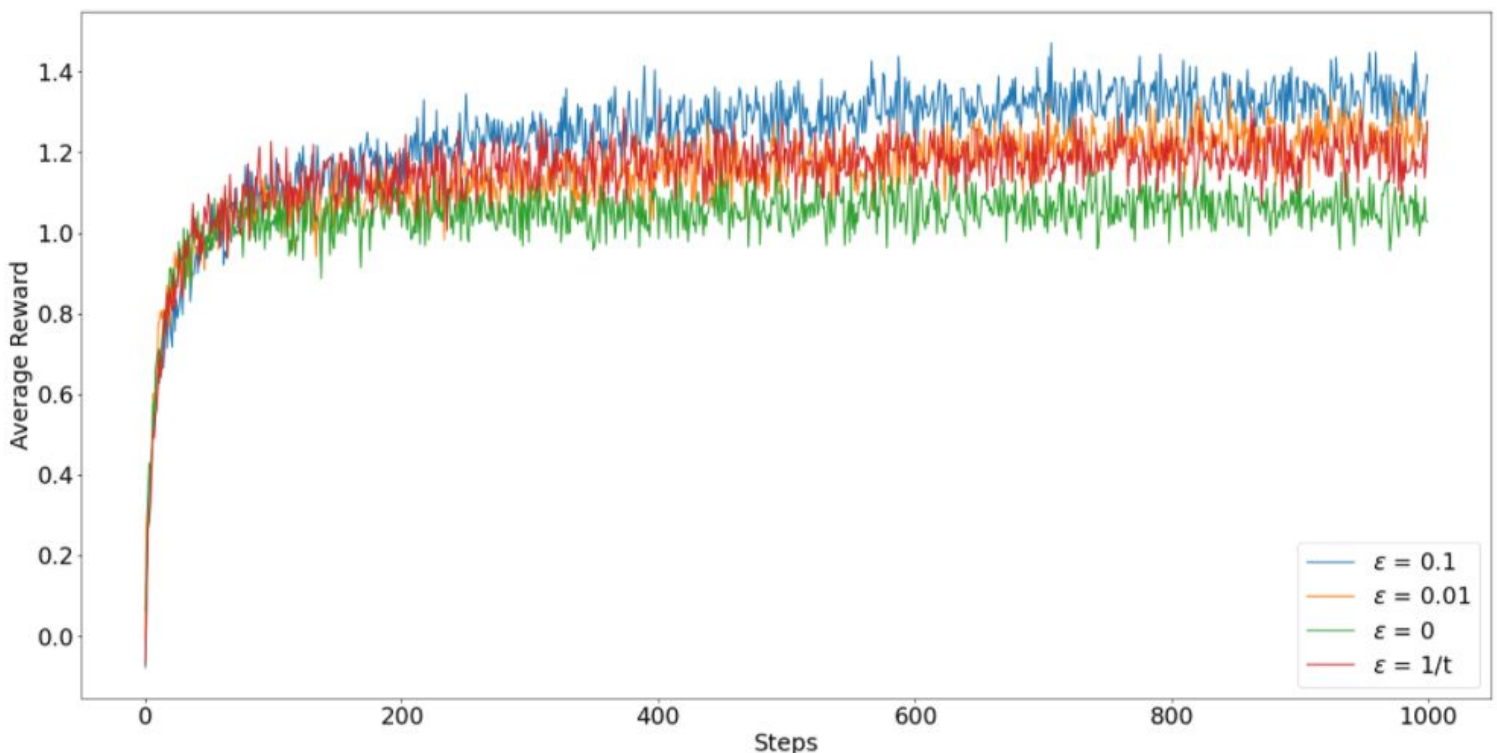


Result: For time steps=1000, $\epsilon = 0.1$ gives the better result but this is not true for long runs as then $\epsilon = 0.01$ gives the better result as it picks up optimal action throughout the simulation with greater probability as compared to $\epsilon = 0.1$.

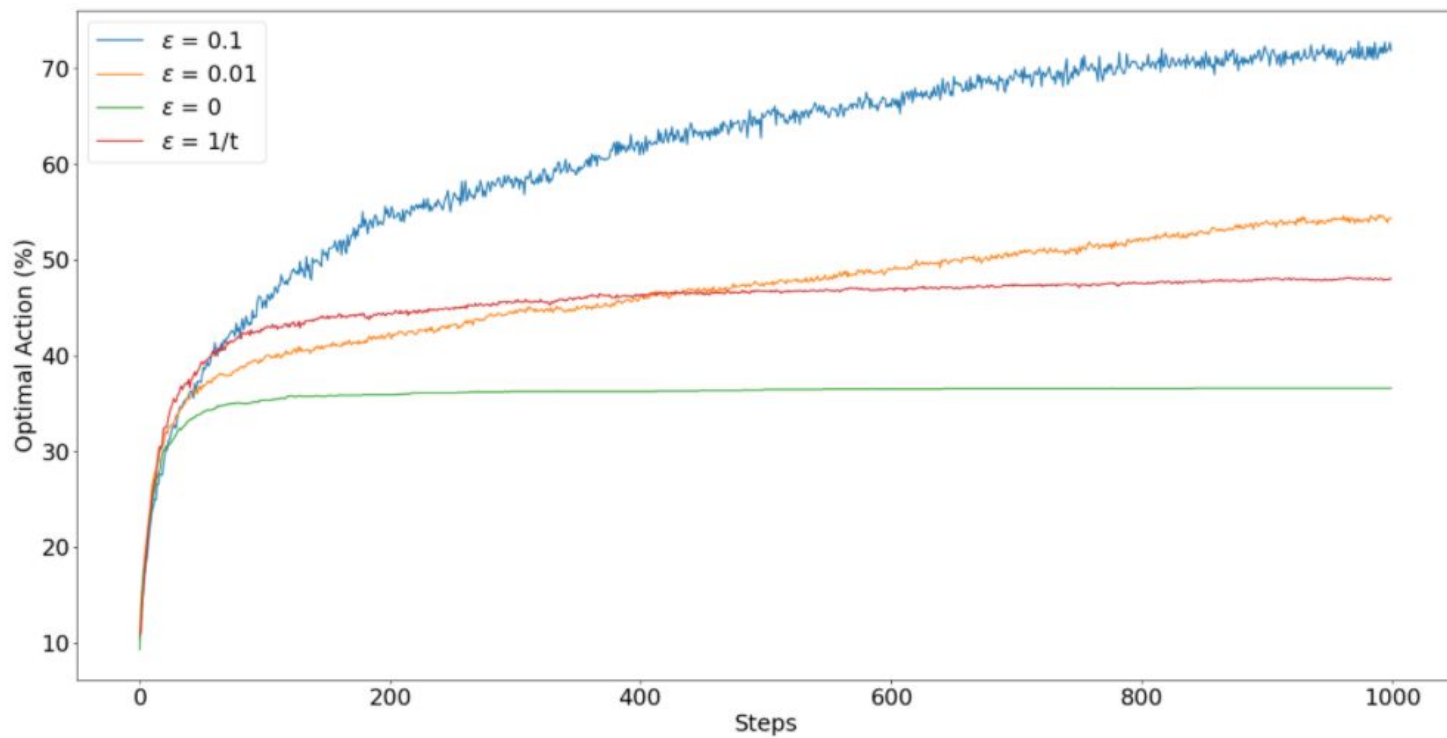
Also, for long runs, $\epsilon = 1/t$ gives us the best result as it increases its probability to exploit as t increases thus motivating exploration early on to find optimal estimate and then exploiting once it has been found (as $t \rightarrow \infty, \epsilon \rightarrow 0$)

Question 2 :

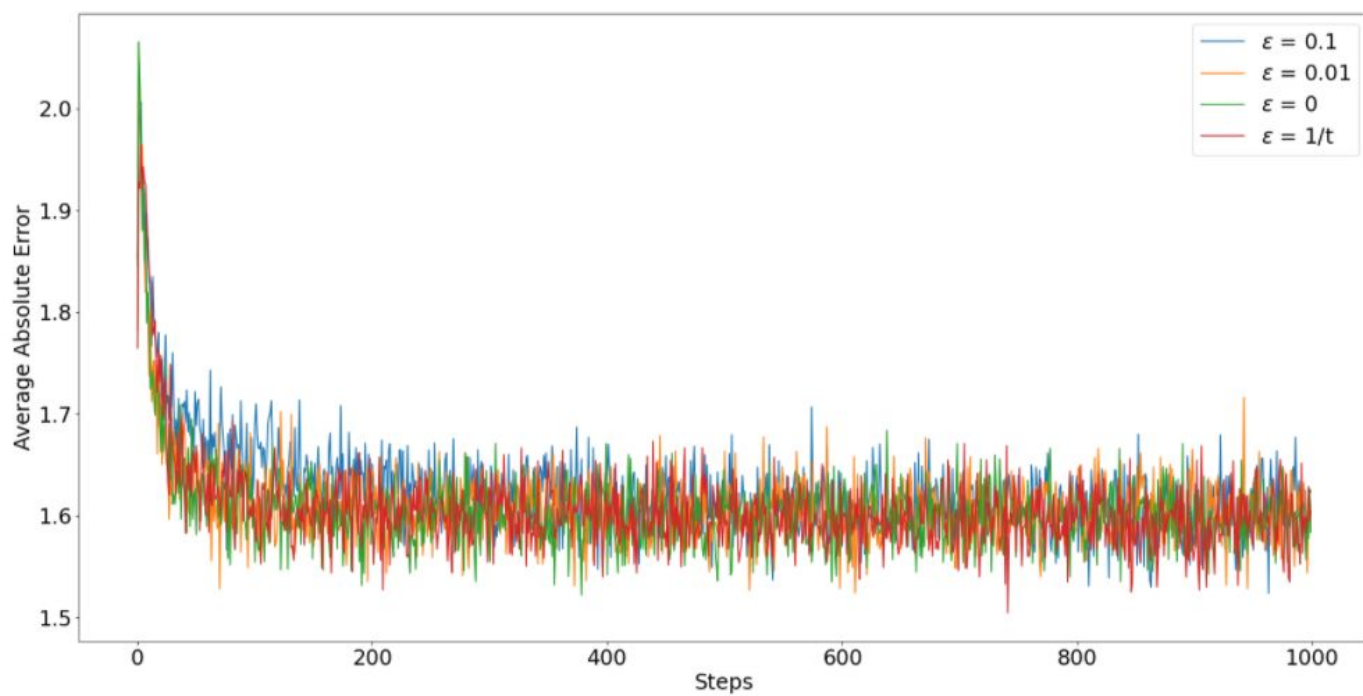
Same as Question1 but now rewards for each arm are selected according to a normal distribution with mean $q^*(A_t)$ and variance 4.



Average Reward



Optimal Action (%)



Average Absolute Error

Result: The same trend follows as in Question1 for different values of ϵ . But, now the rewards and estimates are much noisier and variance= 4 also negatively effects the optimal action for time steps=1000 (can be seen through the range of values on the y-axis)

Question 3

In the long run, $Q_t(a)$ for $\epsilon = 0.01, 0.1$ and $\frac{1}{t}$ would converge to become $Q^*(a)$ (following sample-average method).

- This is not the case for $\epsilon = 0$ as it always picks the greedy estimate hence it performs the poorest.
- $\epsilon = 0.01$ in the long run performs better than $\epsilon = 0.1$ as :
- for $\epsilon = 0.01$

$$P[\text{picking greedy action}] = 1 - \epsilon + \frac{\epsilon}{|A|}$$

$$= 0.9991$$

- thus it exploits and picks up the optimal arm with more probability when compared to $\epsilon = 0.1$ as $P[\text{picking greedy action} | \epsilon = 0.1] = 0.91$

- ~~this is~~ It may be possible that initially $\epsilon = 0.01$ takes more time to explore when compared to $\epsilon = 0.1$ case but since $t \rightarrow \infty$, ~~then it can be said~~ so now $\epsilon = 0.01$ would come up with a better estimate which it now chooses with greater probability.

★ For $\epsilon = \frac{1}{t}$,

Our probability for exploring is initially high and reduces over time. This helps in knowing a better estimate of Q_t and as t increases, we shift our focus to exploitation as t then a pretty good estimate of Q_t that is closest to Q_t^* would be found.

Also as $t \rightarrow \infty$, $\epsilon = 0$ and thus when Q_t^* is approached then we always exploit.

★ Thus, overall $\epsilon = \frac{1}{t}$ will perform best in the long run. t (may not be the case for small time steps)

Question-4

→ Sample mean

- It is not influenced by initial choice of $Q_1(a)$ as soon as all the arms get selected atleast once. From then onwards, $Q_1(a)$ has no effect on Q_n .

$$\begin{aligned}
 \Rightarrow Q_n &= \frac{1}{n-1} \sum_{i=1}^{n-1} R_i^0 \\
 &= \frac{1}{n-1} \left(R_{n-1} + \sum_{i=1}^{n-2} R_i^0 \right) \\
 &= \frac{1}{n-1} \left(R_{n-1} + \frac{(n-2)}{(n-2)} \sum_{i=1}^{n-2} R_i^0 \right) \\
 &\Rightarrow \frac{1}{n-1} \left(R_{n-1} + (n-2) Q_{n-1} \right) \\
 &\quad \text{as } Q_{n-1} = \frac{\sum_{i=1}^{n-2} R_i^0}{n-2}
 \end{aligned}$$

$$\Rightarrow \frac{1}{n-1} \left(R_{n-1} + \overset{(n-1)}{Q_{n-1}} - Q_{n-1} \right)$$

$$Q_n = Q_{n-1} + \frac{1}{(n-1)} (R_{n-1} - Q_{n-1})$$

- No dependence on Q_1 can be seen for $n \neq 1$ apart from $n=2$.

→ Constant step-size α

This method aims at weighing past rewards with decreasing weight as t increases. Thus even though α becomes very large, there is some dependence seen ~~seen~~ on $Q_1(a)$.

This can be shown as:

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

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

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

$$= \alpha R_n + (1-\alpha)\alpha R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$= \alpha R_n + (1-\alpha)\alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2}$$

$$+ \dots + (1-\alpha)^{n-1} \alpha R_1 + (1-\alpha)^n Q_1$$

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

- Here, clear dependence on Q_1 can be seen with weight $(1-\alpha)^n$.

Also, if α is ~~more~~ smaller than ~~for~~

though $(1-\alpha)^n$ continues to decrease over

time ~~less~~, when compared to case where α is larger, it is ~~larger~~ still larger.

$$\bullet \quad (1-\alpha_1)^n > (1-\alpha_2)^n$$

where if $\alpha_1 < \alpha_2$

→ we can propose our step size to be:

$$\alpha_{\text{new}} = \frac{\alpha_{\text{old}}}{\bar{\theta}_n}$$

wherein $\bar{\theta}_n = \bar{\theta}_{n-1} + \alpha_{\text{old}} (1 - \bar{\theta}_{n-1})$

for all $n \geq 0$

and $\bar{\theta}_0 = 0$

* Now to show independence from $Q_1(a)$:

Now ~~for~~ for Q_1 as

$$\alpha_{\text{new}} = \frac{\alpha_{\text{old}}}{\bar{\theta}_1}$$

and $\bar{\theta}_1 = 0 + \alpha_{\text{old}} (1 - 0)$

$$\bar{\theta}_1 = \alpha_{\text{old}}$$

$$\boxed{\alpha_{\text{new}} = 1}$$

$$\bar{\theta}_0 = 0$$

and.

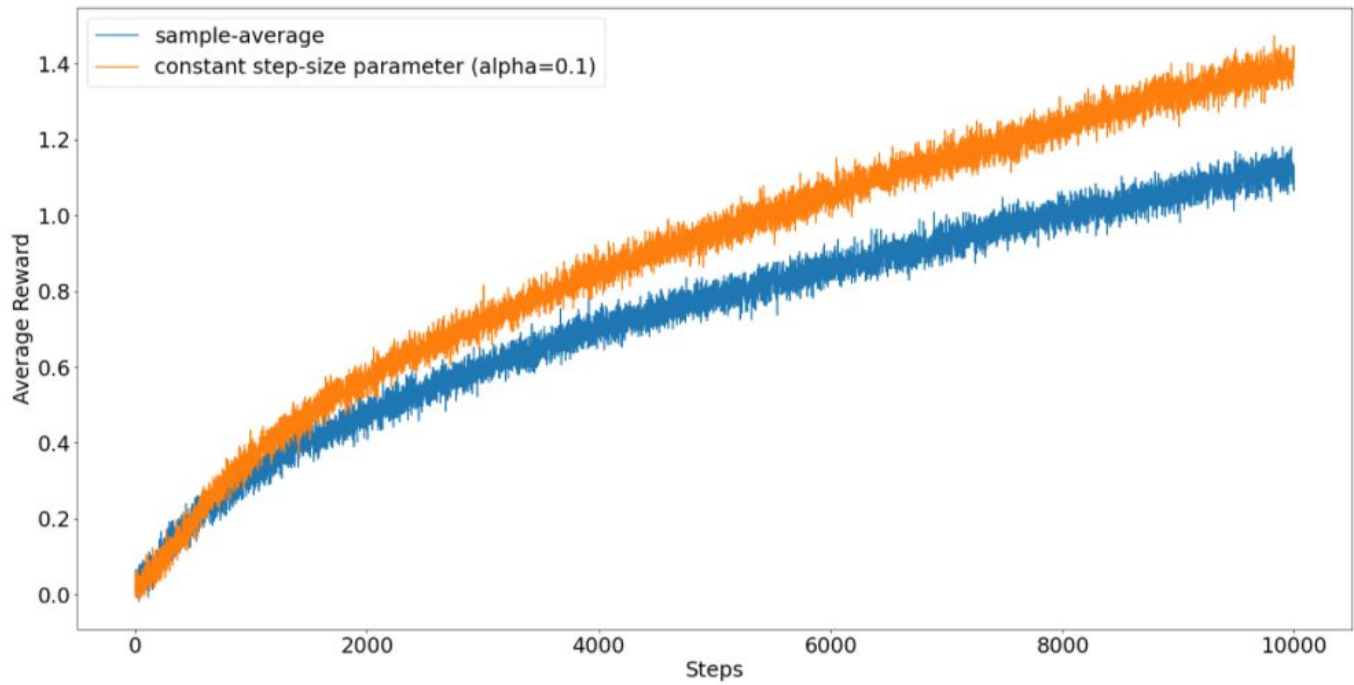
$$Q_{n+1} = (1-\alpha) Q_1 + \sum_{i=1}^n \alpha \left(\frac{\alpha_{old}}{\bar{\alpha}_n} \right) \left(\frac{1-\alpha_{old}}{\bar{\alpha}_n} \right)^{n-i} R_i$$

Thus it can be clearly seen that

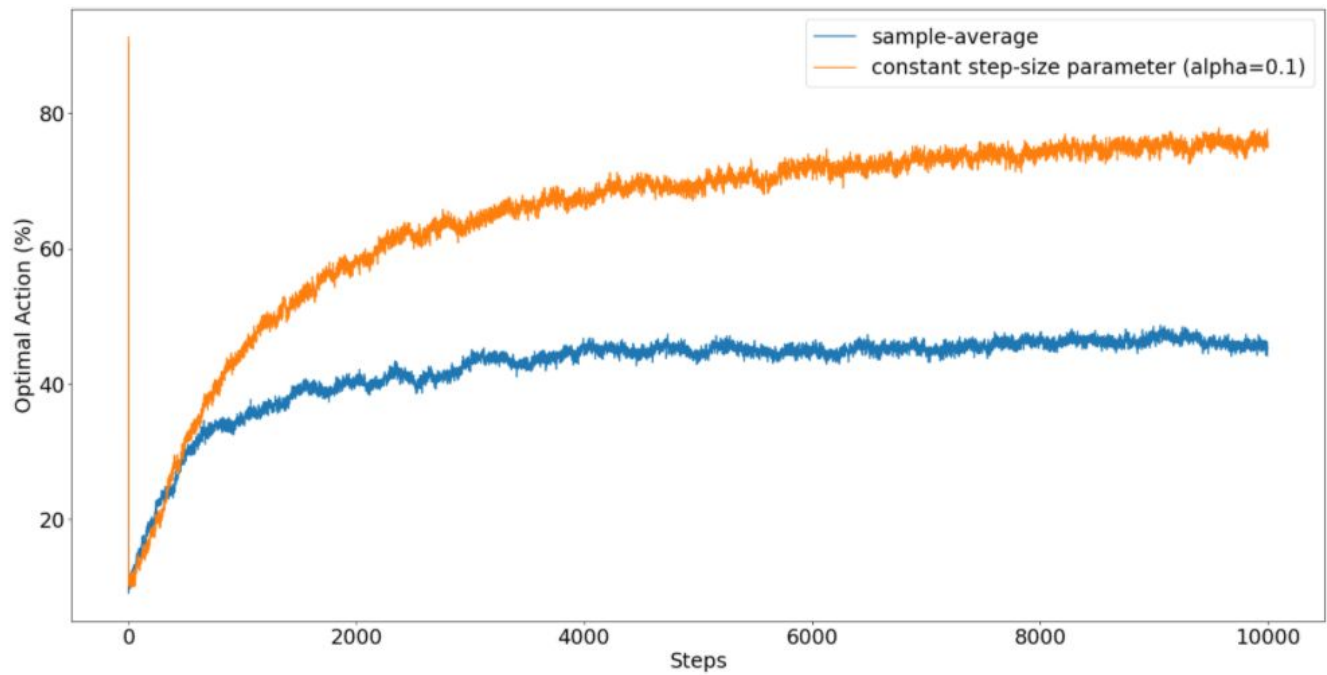
for $\alpha_{new} = \frac{\alpha_{old}}{\bar{\alpha}_n}$, Q_{n+1} does not

depend on $Q_1(\alpha)$.

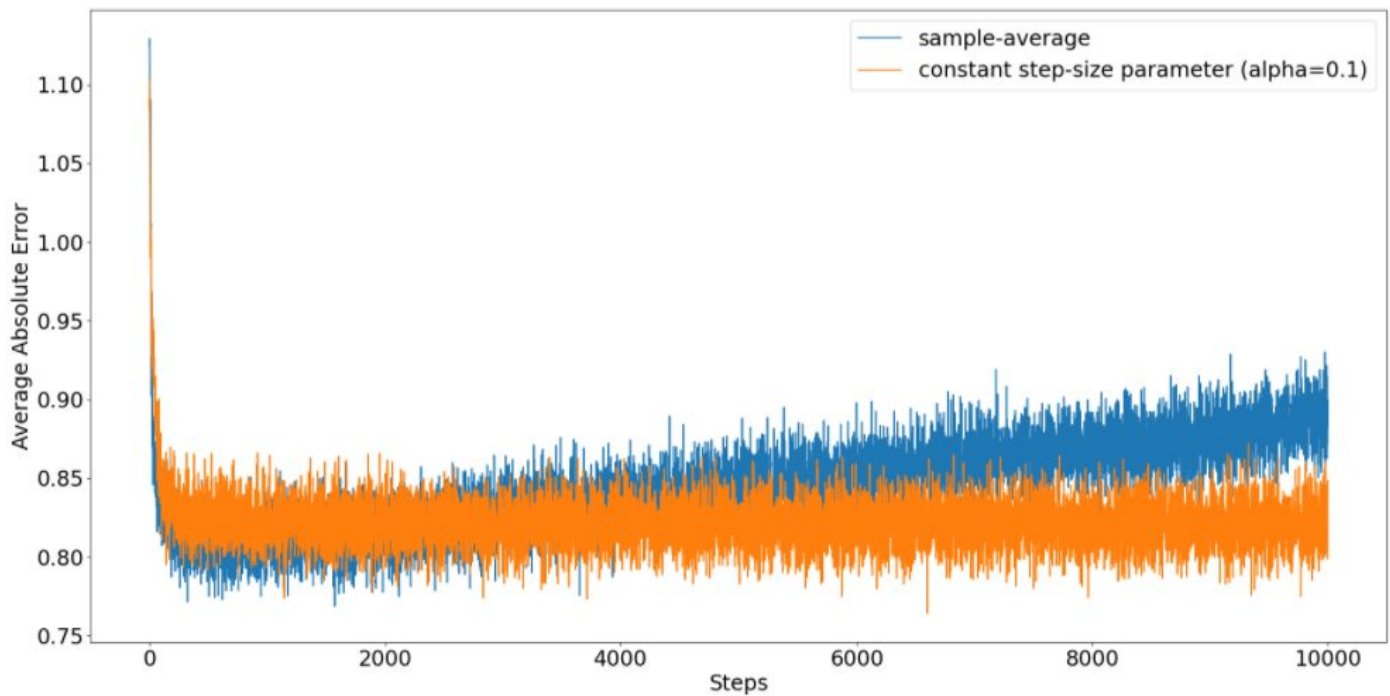
Question 5



Average Reward



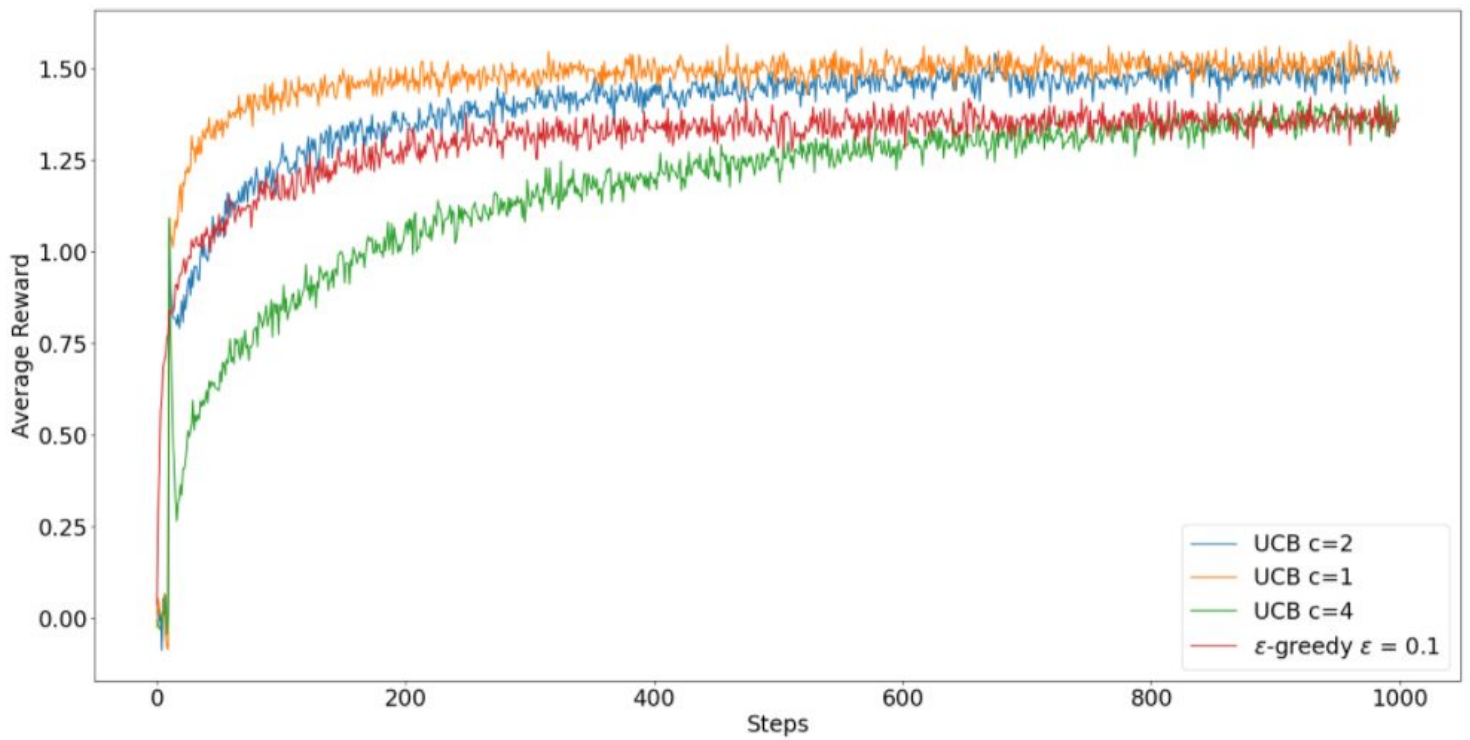
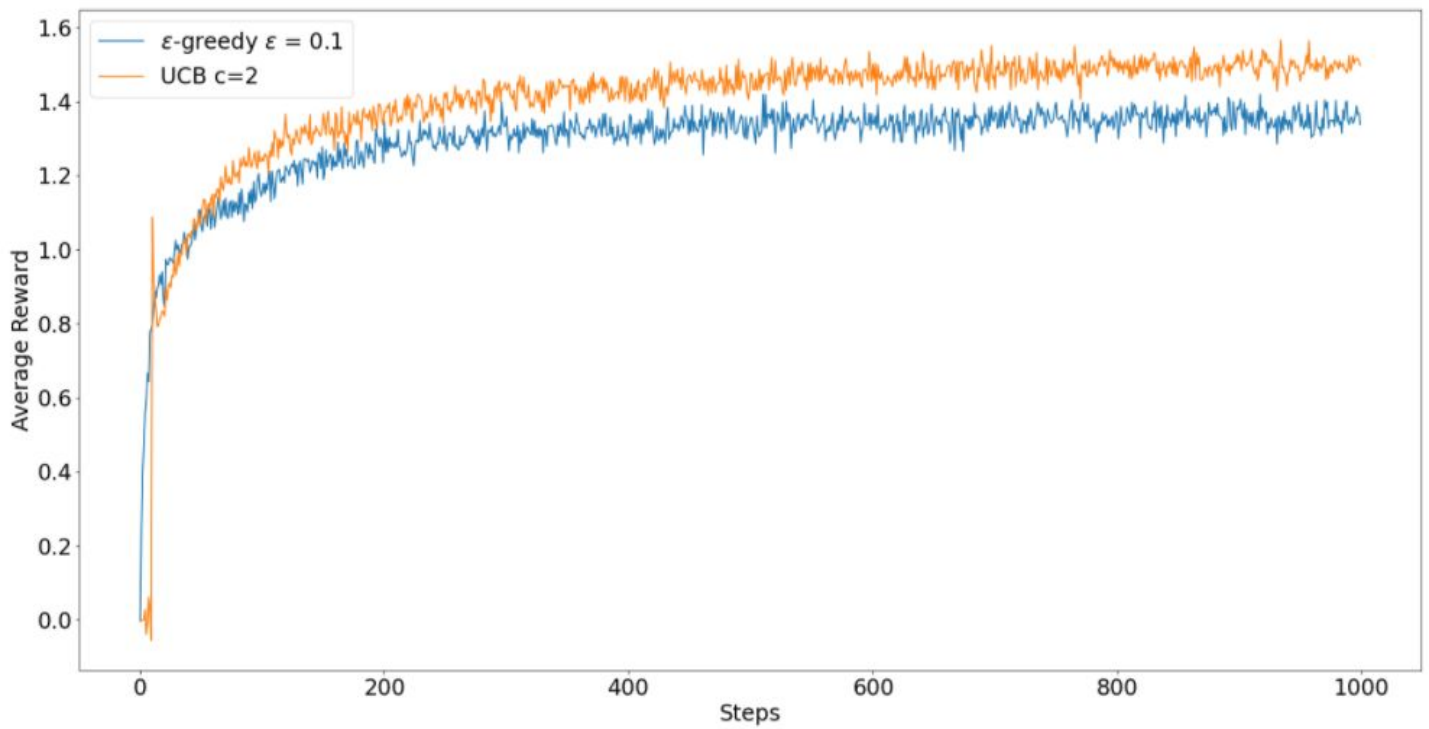
Optimal Action (%)



Average Absolute Error

Result: For non-stationary setting, constant step-size parameter tends to give better results as they, unlike sample-average, do not converge and thus readjust themselves according to any change in the distribution of the arms.

Question 6:

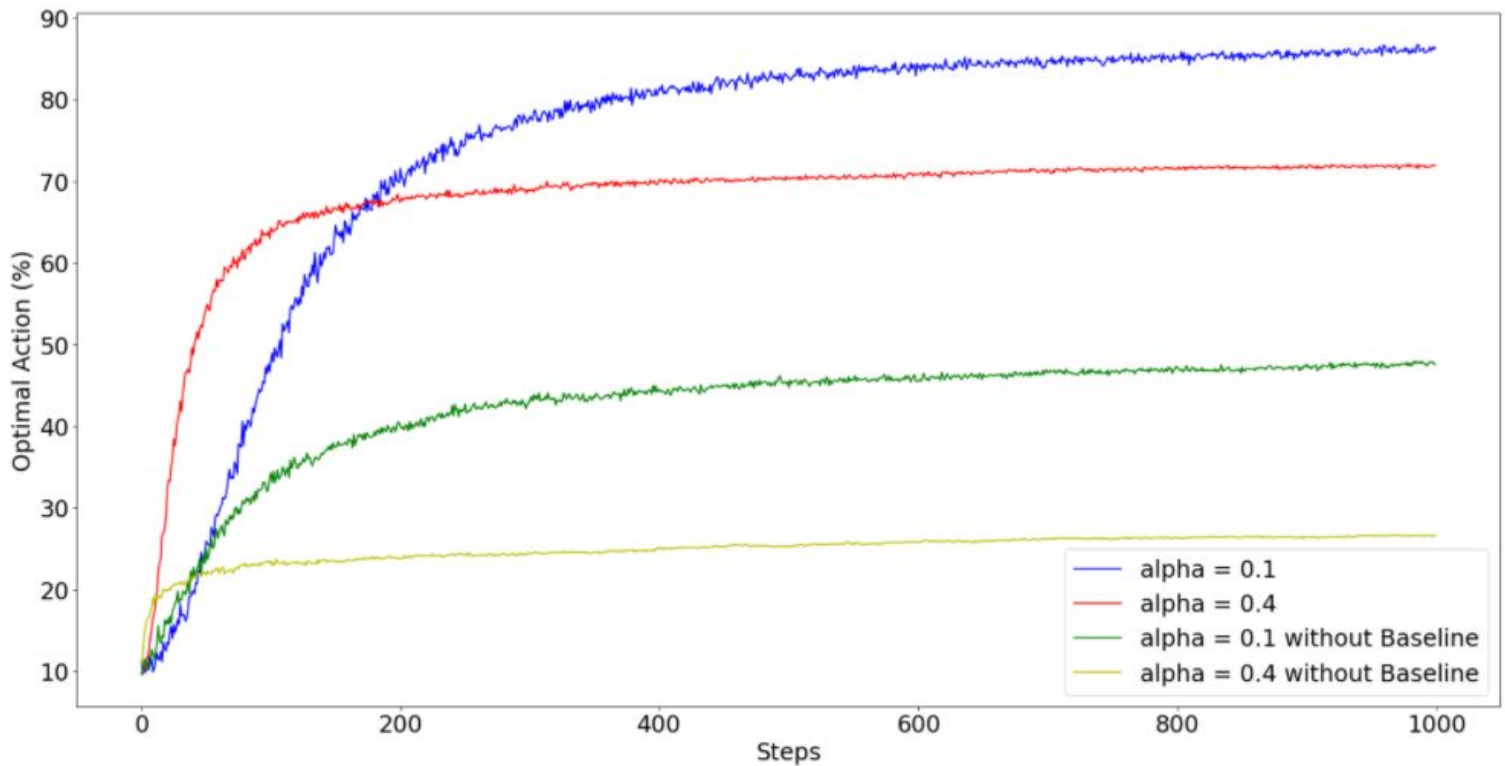


Result: According to UCB, when $N_t(a)=0$ for any arms A_t , then that arm is considered to be the maximizing action.

Now, this occurs for all the 10 arms initially in random order. Thus, we see the low average reward for the first 10 steps as then we are exploring and picking a different arm each step.

Now for the 11th step, the term $c * \sqrt{\ln t / N_t(a)}$ is the same irrespective of the picked arm. Thus, now the discriminatory factor will only be the value of Q_t which is dependent on the first round of exploration (affected by the reward of each arm and Q_1). Now, that action is picked whose Q_t is greater and this would be seen across all 2000 simulations, thus a spike in average reward can be seen for the 11th step. Further onwards $c * \sqrt{\ln t / N_t(a)}$ is different for each arm and thus average rewards progress steadily as the model acquires more certainty about the system.

Question 7:



Result: The optimal action is better for $\alpha = 0.1$ and with Baseline rather than for without baseline. It is because the Baseline is used to capture the deviation from the mean rather than the absolute value.

Both lower α and Baseline lead to lower variance changes thus performing better.