# RL - Assignment -1
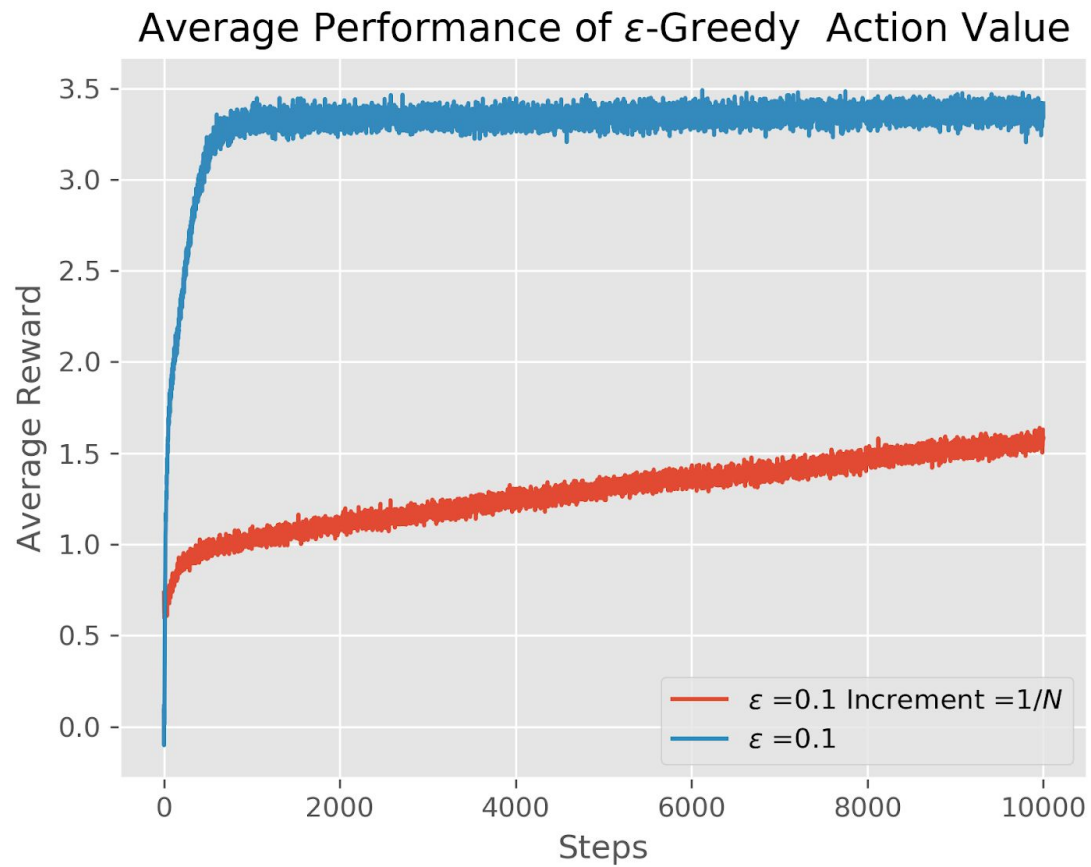
(Ashish Sethi-MT18024)

# Table of Contents

# Exercise - 2.5

Average Performance of $\varepsilon$-greedy action value on the ten-armed testbed

Optimal Actions of $\varepsilon$-greedy action value on the ten-armed testbed

# Figure 2.3

For Stationary



Effect of Optimistic Intial Action Value on 10 Armed Testbed

- Initially when our agent the exploring it selects all arms once and gets mostly disappointed until it tries out all the arms. But after trying out all the arms there is a high probability to select the optimal action that is why we see the spike at the $t = 10$.
- Another thing which we observed prior to the probability of selecting optimal actions is about $10\%$.

For Non-Stationary

## Effect of Optimistic Intial Action Value



In non-stationary case observe similar observation as compared to stationary case and some more observations which are more related to the non-stationary property of bandits.

- Initially when our agent the exploring it selects all arms once and gets mostly disappointed until it tries out all the arms. But after trying out all the arms there is a high probability to select the optimal action that is why we see the spike at the $t = 10$.
- Another thing which we observed prior to the probability of selecting optimal actions is about $10\%$.
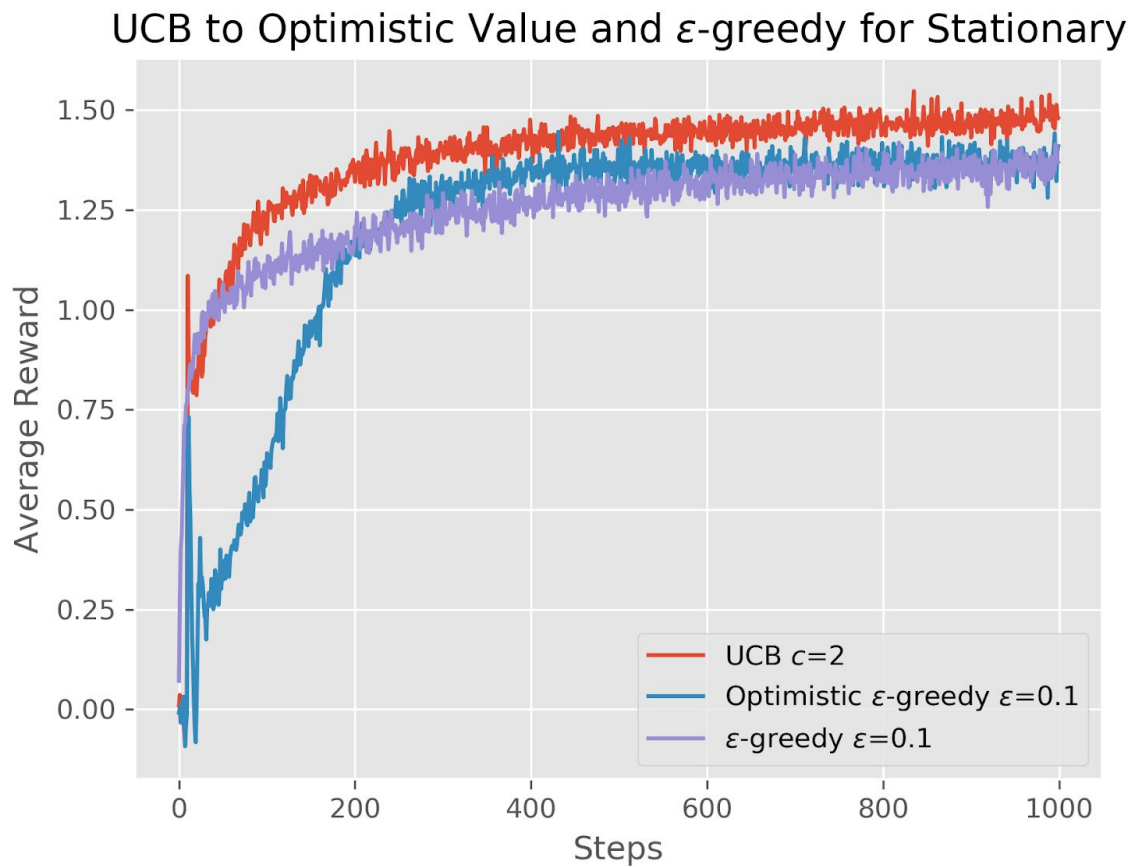- In the case of the non-stationary condition, we can see that $Q = 0$ is performing better than the $Q = 5$.

# Comparison of UCB to Optimistic Value and $\varepsilon$-greedy

## For Stationary

Average rewards Vs Number of total steps



Observations :
- In the case of stationary bandits, we can see that UCB is performing best as compared to $\varepsilon$-greedy method. After UCB optimistic $\varepsilon$- greedy method is performing which ideally should be.

% Optimal Action Vs Number of total steps



Observations :
Similar trends of average rewards we can see in the %optimal Actions graphs. Which theoretical also proved to be right.

# For Non-Stationary

Average rewards Vs Number of total steps



Observations:
In the case of non-stationary bandits $\varepsilon$- greedy is performing best among all the methods, Major reason for this is UCB often performs well, as shown here, but is more di!cult than "-greedy to extend beyond bandits to the more general reinforcement learning settings considered in the rest of this book.
One di!culty is in dealing with nonstationary problems.

% Optimal Action Vs Number of total steps

UCB to optimistic value and $\varepsilon$-greedy for non stationary

Q 3

As we know that

$$\mathcal{G}_{n+1} = \mathcal{G}_n + \alpha \left[ R_n - \mathcal{G}_n \right] \cdot --- \text{①}$$

$$= \alpha R_n + (1-\alpha) \mathcal{G}_n --- \text{②}$$

Now we change this to $\beta$ which is a step size constant.

and $\beta = \dfrac{\alpha}{\bar{O}_n}$

putting values of $\beta$ in eq$^n$ ②

$$\mathcal{G}_{n-1} = \frac{\alpha}{\bar{O}_n} R_n + \left( 1 - \frac{\alpha}{\bar{O}_n} \right) \mathcal{G}_n$$

$$\frac{\alpha}{\bar{O}_n} R_n + \left( \frac{\bar{O}_n - \alpha}{\bar{O}_n} \right) \mathcal{G}_n$$

By substituting $\bar{O}_n = \bar{O}_{n-1} + \alpha \left( 1 - \bar{O}_{n-1} \right)$

$$\frac{\alpha}{\bar{O}_n} R_n + \left( \frac{\bar{O}_{n-1} + \alpha (1 - \bar{O}_{n-1}) - \alpha}{\bar{O}_n} \right) \mathcal{G}_n$$

$$\frac{\alpha R_n}{\sigma_n} + \left[ \frac{\bar{\sigma}_{n-1} + \alpha - \alpha \bar{\sigma}_{n-1} - \alpha}{\sigma_n} \right] \phi_n.$$

$$\frac{\alpha R_n}{\sigma_n} + \frac{\bar{\sigma}_{n-1}}{\sigma_n} (1-\alpha) \phi_n \qquad — ③$$

Now from eq$^n$ ① we again put the value of $\phi_n$.

$$\phi_{n+1} = \frac{\alpha R_n}{\sigma_n} + \frac{\bar{\sigma}_{n-1}}{\sigma_n} (1-\alpha) \left[ \frac{\alpha}{\bar{\sigma}_{n-1}} R_{n-1} + \left( 1 - \frac{\alpha}{\sigma_{n-1}} \right) O_{n-1} \right]$$

$$= \frac{\alpha R_n}{\sigma_n} + \frac{\alpha(1-\alpha)}{\sigma_n} R_{n-1} + \frac{\bar{\sigma}_{n-1}}{\sigma_n} (1-\alpha)(1-\alpha) \frac{\bar{\sigma}_{n-2}}{\sigma_{n-1}} \phi_{n-1}$$

$$= \frac{\alpha . R_n}{\sigma_n} + \frac{\alpha(1-\alpha)}{\sigma_n} R_{n-1} + (1-\alpha)^2 \frac{\bar{\sigma}_{n-2}}{\sigma_n} \phi_{n-1}.$$

$$= (1-\alpha)^n \frac{\bar{\sigma}_0}{\sigma_n} \phi_1 + \sum_{i=1}^{n} \alpha (1-\alpha)^{n-1} \frac{R_i}{\bar{\sigma}_i} \qquad — ④$$

from equation ④

$$\phi_{n+1} = (1-\alpha)^n \underbrace{\bar{O}_0}_{\bar{O}_2} \cdot \phi_1 + \sum_{i=1}^{n} \alpha (1-\alpha)^{n-i} \underbrace{R_i}_{O_i}$$

since $\bar{O}_0 = 0$ (given)

equation ④ becomes.

$$\phi_{n+1} = \sum_{i=1}^{n} \alpha (1-\alpha)^{n-i} \frac{R_i}{O_i}$$

## Method II

In this method we solve our equation from starting point zero whreas in the last method we solved from the back side that is $n$.

$$O_1 = \bar{O}_0 + \alpha (1 - O_0) \qquad (given)$$

$$= \bar{O}_0 + \alpha - \alpha \bar{O}_0$$

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

$$\bar{O}_1 = \alpha$$

similarly for $\bar{O}_2$

$$\bar{O}_2 = \bar{O}_1 + \alpha(1 - \bar{O}_1)$$
$$= \bar{O}_1^{\alpha} + \alpha(1 - \alpha)$$
$$= \alpha + \alpha - \alpha^2$$
$$\bar{O}_2 = 2\alpha - \alpha^2$$

And from constant step size equation
We know

$$\mathcal{P}_{n+1} = \alpha \mathcal{P}_n + (1 - \alpha) \mathcal{P}_n$$

Substituting $\alpha$ with $\beta_n = \dfrac{\alpha}{\bar{O}_n}$

In this case $\mathcal{P}_1$ doesn not

because $\bar{O}_0$ is zero

so the Next term.

$$\phi_2 = \frac{\alpha}{\sigma_1} R_1 + \left(1 - \frac{\alpha}{\sigma_1}\right) \phi_1$$

$$= \frac{\alpha}{\alpha} R_1 \left(1 - \frac{\alpha}{\alpha}\right) \phi_1$$

$$= R_1 + 0.$$

$$\boxed{\phi_2 = R_1} \qquad - ①$$

for $\phi_3$.

$$\phi_3 = \frac{\alpha}{\sigma_2} R_2 + \left(1 - \frac{\alpha}{\sigma_2}\right) \phi_2$$

$$= \frac{\alpha}{\alpha(2-\alpha)} R_2 + \left(1 - \frac{\alpha}{2(1-\alpha)}\right) \phi_2$$

$$= \frac{R_2}{2-\alpha} + \left(\frac{2 - 2\alpha - \alpha}{2(1-\alpha)}\right) \phi_2$$

and $\phi_2 = R_2$

$$\frac{R_2}{2-\alpha} \left[ R_2 + (3-\alpha)R_1 \right] - ②$$

if we proceed further in similar fashion we can see that there will no term for $\theta_n$ involved. which was involved in the case of constant step size.