

Q2 Initially  $Q_0$  dominates over any result  
given by any of the bandits. So the  
algorithm runs in batches of 10, where  
~~but~~ at every  $(10n+1)$ th step we  
each arm is picked once in each  
batch. But within these batches the  
order of picking the bandits is entirely  
dependant on the result shown so far.  
And as ~~bandits~~ best bandits, have the highest  
chance of having the best reward on the  
batches so far. Thus, the initial  
few picks <sup>in each batch</sup> have a much higher chance to be  
optimum. Thus when going from the last i.e.  $10n$   
step to  $10n+1$  there is a sharp increase in  
~~probability~~ % chance optimum, leading to the  
spike.

$$Q_{n+1} = Q_n + \beta_n [e_n - Q_n]$$

$$= Q_n + \alpha$$

Q3  $\bar{Q}_n = \bar{Q}_{n-1} + \alpha(1 - \bar{Q}_{n-1})$

$$= \alpha + (1-\alpha)\bar{Q}_{n-1}$$

$$= \alpha + (1-\alpha)\left(\alpha + \frac{(1-\alpha)}{\alpha}\bar{Q}_{n-2}\right)$$

$$= \alpha + (1-\alpha)\alpha + (1-\alpha)^2\alpha \dots (1-\alpha)^{n-1}\alpha + 0$$

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

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

$$= 1 - (1-\alpha)^n \quad n \geq 1$$

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

$$= Q_n + \frac{\alpha}{\bar{Q}_n} [R_n - Q_n]$$

$$= Q_n \left[ 1 - \frac{\alpha}{\bar{Q}_n} \right] + \beta_n R_n$$

$$= \cancel{Q_n} \beta_n R_n + \left(1 - \frac{\alpha}{\bar{Q}_n}\right) (\beta_{n-1} R_{n-1} + (1 - \frac{\alpha}{\bar{Q}_{n-1}}) \dots)$$

$$\begin{aligned}
 & \Rightarrow \alpha R_n + (1-\alpha) \alpha R_{n-1} \\
 & = \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)(1-\alpha) R_{n-2} \\
 & \quad + (1-\alpha)(1-\alpha)(1-\alpha) \dots (1-\alpha) Q_1
 \end{aligned}$$

$$\begin{aligned}
 & \Rightarrow \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2} + \dots \\
 & \quad \left( \text{using } (1-\alpha)^n = \frac{(1-\alpha)(1-(1-\alpha)^{n-1})}{1-(1-\alpha)} \right) \\
 & \quad \text{\& then telescoping} \\
 & \alpha R_n + (1-\alpha) \alpha R_{n-1} + \dots \\
 & = \alpha R_n + \frac{(1-\alpha)(1-(1-\alpha)^{n-1})}{1-(1-\alpha)} \alpha R_{n-1} + \frac{(1-\alpha)^2(1-(1-\alpha)^{n-2})}{1-(1-\alpha)} \alpha R_{n-2} \\
 & \quad \dots + \frac{(1-\alpha)^n(1-(1-\alpha)^0)}{1-(1-\alpha)} Q_1
 \end{aligned}$$

$$\Rightarrow \left( \text{using } 1-(1-\alpha)^0 = 1-1=0 \right)$$

$$\begin{aligned}
 & = \alpha R_n + (1-\alpha) \left( \alpha R_{n-1} + (1-\alpha) \alpha R_{n-2} + \dots \right) \\
 & = \alpha R_n + \sum_{i=1}^{n-1} \frac{(1-\alpha)^{n-i} (1-(1-\alpha)^i)}{1-(1-\alpha)} \alpha R_i \\
 & = \sum_{i=1}^n \frac{(1-\alpha)^{n-i} (1-(1-\alpha)^i)}{1-(1-\alpha)} \alpha R_i = \sum_{i=1}^n \frac{R_i (1-\alpha)^{n-i}}{1-(1-\alpha)}
 \end{aligned}$$



$$= \frac{d}{1-d(1-d)^n} \left( R_n + \cancel{R_{n-1}} (1-d) R_{n-1} + (1-d)^{n-1} R_1 \right)$$

(This has no factor of  $d$ , &  $R_1, R_2, \dots$  are clearly ~~identical~~ <sup>factorial</sup> with geometric  $\sum_{i=1}^{\infty} (1-d)^{i-1}$ )

To show weighted sum of factors to be 1

$$\text{sum of factors } \frac{1}{(1-d)^n} = \frac{1 + (1-d) + \dots + (1-d)^{n-1}}{(1-d)^n}$$

$$= \frac{d}{1-(1-d)^n} \left( \frac{1 - (1-d)^n}{1-(1-d)} \right)$$

$$= \frac{d}{1-(1-d)^n} \frac{(1-(1-d)^n)}{d} = 1$$

hence  $Q_{n+1}$  is weighted sum of  $R_i$ 's

Q4

Initially optimistic explores a lot, but allowing it make better & greedy decisions but this effect quickly wears off, & then the ~~actual~~  $\Phi^*$  change significantly over time but optimistic due to its  $\epsilon = 0$ , doesn't explore enough to keep up while  $\epsilon$ -greedy doesn't suffer from this as it continues to explore  $\epsilon = 0$  & isn't weighed down by past explorations too much due to its constant step propagation.

In stationary both UCB improve over time as ~~the~~ their estimate of  $\Phi^*$  improves and they increasingly waste less time exploring the UCB with higher  $\epsilon$  lags behind as spends more time exploring. Optimistic however spends very little time exploring with most of its exploration in the early part, by which time it's



In non-stationary ~~UCB~~ ~~initially~~ optimistic  
 starts logging, since with ( $E=0$ ) it  
~~stop~~ severely reduces its exploration after  
 initial few terms, so as  $\phi^*$  changes, ~~optimal~~  
 fails to adjust.

UCB explores a little more but not as  
 much  $\epsilon$ -greedy since  $\ln t$  ~~reduces~~  
 $\epsilon$  grows far more slowly than  $\ln t$ ,  
 thus eventually UCB also ~~has~~ becomes  
 outdated, if  $c$  is higher there is more  
 exploration hence ~~UCB~~ UCB with higher  
 $c$  perform better overtime.