Divyam Anshumaan
2017147

**Q3.**

$$B_n = \frac{\alpha}{\theta_n}, \qquad \theta_n = \theta_{n-1} + \alpha(1 - \theta_{n-1})$$
$$\theta_0 = 0, \quad n \geq 0.$$

now, $\theta_n = \theta_{n-1} + \alpha(1 - \theta_{n-1})$

$$= \alpha + (1-\alpha)\theta_{n-1}$$
$$= \alpha + (1-\alpha)(\theta_{n-2} + \alpha(1 - \theta_{n-2}))$$
$$= \alpha + \alpha(1-\alpha) + (1-\alpha)^2 \theta_{n-2}$$
$$\vdots$$
$$= \sum_{i=1}^{n} \alpha(1-\alpha)^{n-i} + (1-\alpha)^n \theta_0 + \alpha$$

$$= \alpha \left( \sum 1-\alpha \right)$$

$$= \alpha(1-\alpha)^n \left( \sum_{i=1}^{n} \frac{1}{(1-\alpha)^i} \right) + (1-\alpha)^n \times 0$$

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

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

So, $$B_n = \frac{\alpha}{1 - (1-\alpha)^n} \qquad - (1)$$

Now, $\theta_{n+1} = R_n \cdot B_n + (1 - B_n)Q_n$

2) $\theta_{n+1} = R_n \cdot \frac{\alpha}{\theta_n} + \frac{\theta_n - \alpha}{\theta_n} \cdot Q_n$

Now

$$Q_n - \alpha = Q_{n-1}(1-\alpha) \quad --[11]$$

$$= \frac{R_n\,\alpha}{Q_n} + \frac{Q_n - \alpha}{Q_n} \cdot (R_{n-1}\,\beta_{n-1} + (1-\beta_{n-1})\,Q_{n-1})$$

$$= \frac{R_n\,\alpha}{Q_n} + \frac{Q_n - \alpha}{Q_n} \cdot \frac{\alpha}{Q_{n-1}} R_{n-1} + \frac{Q_n - \alpha}{Q_n} \cdot \frac{Q_{n-1} - \alpha\,Q_{n-1}}{Q_{n-1}}$$

$$\vdots$$

$$= R_n \cdot \frac{\alpha}{Q_n} + \frac{Q_{n-1}(1-\alpha)}{Q_n} \frac{\alpha}{Q_{n-1}} R_{n-1} + \frac{\beta(1-\alpha)^2 \alpha R_{n-2}}{Q_n} \cdots$$

$$\cdots + F \cdot \otimes \qquad \cancel{\frac{\alpha}{Q_{n-1}} \frac{(Q_{n-1})(1-\alpha)}{(Q_{n-1})} \cdots \frac{(1-\alpha)^n Q_0}{Q_1} Q_0}$$

$$= \frac{\alpha \sum R_i\,(1-\alpha)^{n-i}}{Q_n} + 0$$

→ As $R_i$'s are multiplied by $(1-\alpha)^{n-i}$

the weights are exponentially decreases for

<space sz=6 />$\underset{\times}{\longleftarrow} \quad \underset{\times}{\phantom{xx}} \quad \text{(from above)} \qquad \otimes$

$$\frac{Q_n - \alpha}{Q_n} \cdot \frac{Q_{n-1} - \alpha}{Q_{n-1}} \cdot \frac{Q_{n-2} \cdot \alpha}{Q_{n-2}} \cdots \frac{Q_1 - \alpha}{Q_1} \cdot Q_1$$

<space sz=30 />$= F$.

$$= \frac{(1-\alpha) \cdot Q_{n-1}}{Q_n} \cdot \frac{(1-\alpha)\,Q_{n-2}}{Q_{n-1}} \cdots \frac{Q_0}{Q_1} = F = 0$$

<space sz=16 />as $Q_0 = 0$.

* Now weights for $R$ is: $\dfrac{\alpha(1-\alpha)^{n-i}}{1-(1-\alpha)^n}$

now $\sum \cdot \dfrac{\alpha}{1-(1-\alpha)^n} \cdot (1-\alpha)^n \sum_{i=1}^{n} \dfrac{1}{(1-\alpha)^i}$

$= \dfrac{\alpha}{1-(1-\alpha)^n} \cdot (1-\alpha)^n \left( \dfrac{1}{(1-\alpha)} \cdot \dfrac{\dfrac{1-(1-\alpha)^n}{(1-\alpha)^n}}{\dfrac{1-1+\alpha}{(1-\alpha)}} \right)$

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

$= 1$. Hence weights sum to 1.

— x ———— x ———— x ———— 1

• Checking for convergence: $\sum_{i=1}^{n} \beta_i \quad \beta = \alpha$

Now, $\beta_n = \dfrac{\alpha}{1-(1-\alpha)^n}$, in general,

$\sum \dfrac{\alpha}{1-(1-\alpha)^n} = \alpha \quad \text{for } 1 > \alpha > 0,$

$\alpha \quad \to \quad 1$

for $1 \geq \alpha \geq 0$, $\sum \beta_n$

$= \alpha \sum_{n=1}^{\infty} \dfrac{1}{1-(1-\alpha)^n}$

$= \alpha \sum$

$$(1-\alpha) \geq (1-\alpha)^n \geq 0 \quad \text{for } n \geq$$

now

$$|z| = 1 - (1-\alpha)^n \geq 1 - (1-\alpha) \geq \quad \text{——(III)}$$

or

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

$$\alpha \leq \frac{\alpha}{1-(1-\alpha)^n} \leq n \quad \text{for } 0 \leq \alpha \leq n$$

now

$$\frac{\alpha}{1-(1-\alpha)^n} \geq \frac{\alpha}{1} \quad (\text{from III})$$

and

$$\sum_{i=1}^{\infty} \frac{\alpha}{1} \leq \sum_{n=1}^{\infty} \frac{\alpha}{1-(1-\alpha)^n}$$

so

$$\sum_{i=1}^{\infty} \beta_i = \infty$$

Now checking for $\sum \beta_i^2 < \infty$:

ie, $\sum_{n=1}^{\infty} \frac{\alpha^2}{\left(1-(1-\alpha)^n\right)^2} \qquad \sum_{n=1}^{\infty} \frac{\alpha^2}{} < \infty$

## Q2.

### Stationary Optimistic greedy:

1. Spikes are visible after the first pass of the greedy action search. This is:

   - Due to an optimistic initial value, which will decrease after it is selected for a given bandit

2. A large spike (40%) optimistic optimal action is selected, immediatly after the first pass since the actual optimal bandit (with the highest expected reward) will have the greatest estimate after the first pass.

3. It initially lags while exploring all bandits but eventually settles as exploration decreases.

### Non-Stationary Optimistic greedy:

1. Optimistic greedy initially performs better because $q^*(a)$ for bandits has not changed much. But after time passes the initial distribution has changed and if no longer has the correct estimate

2. Realistic greedy keeps exploring estimate

2. Realistic greedy keeps exploring other est actions and eventually gets a better estimate.

## Q4. (Optimal Action)

1. In stationary case, UCB explores all actions first. As $ln(t)$ is bounded the increments will eventually become negligible and will have explored all actions with lower estimates and more frequent examples explored lesser. this is better than a realistic greedy approach that explores. leafs exploring suboptimal actions even after a large no. of steps. Optimistic greedy explores actions and then immediately stops due to bias and then immediatly stops exploring due to sample averaging.

• The spike in UCB appears after the first pass over all bandits, as the actual optimal action will have a better estimate and will get selected more often.