

# RL Assignment-1

A2) (a) At the starting of optimistic greedy, we constantly get disappointed and explore for some initial turns. This happens for all 2000 individual 10-armed bandits. This causes accumulation of ~~the~~ majority taking good or bad step together. This causes spikes in early stages.

(b) As we can infer from the graph optimistic <sup>greedy</sup> is not good method for non-stationary methods. Because after some time this method is the same as greedy which causes no exploration even when the optimal might have changed in non-stationary problem.

A3) As mentioned in the question, taking step size

$$\beta_n = \frac{\alpha}{\bar{Q}_n} \quad \text{where} \quad \bar{Q}_n = \bar{Q}_{n-1} + \alpha(1 - \bar{Q}_{n-1}) \quad \text{for } n \geq 0 \text{ with } \bar{Q}_0 = 0$$

— ①

has no initial bias & is exponentially wated

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

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

$$Q_{n+1} = \frac{1}{\bar{Q}_n} \left[ (\bar{Q}_n - \alpha) Q_n + \alpha R_n \right]$$

— ②

simplifying eq ①

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

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

putting this in eq ②.

$$\theta_{n+1} = \frac{1}{\bar{\theta}_n} \left[ \bar{\theta}_{n-1} (1-\alpha) \theta_n + \alpha R_n \right] \quad (3)$$

putting value of  $\theta_n$  from eq (3)

$$\theta_{n+1} = \frac{1}{\bar{\theta}_n} \left[ \bar{\theta}_{n-1} (1-\alpha) \cdot \frac{1}{\bar{\theta}_{n-1}} (\bar{\theta}_{n-2} (1-\alpha) \theta_{n-1} + \alpha R_{n-1}) + \alpha R_n \right]$$

$$\theta_{n+1} = \frac{1}{\bar{\theta}_n} \left[ \bar{\theta}_{n-2} (1-\alpha)^2 \theta_{n-1} + \alpha R_n + (1-\alpha) \alpha R_{n-1} \right]$$

$$\theta_{n+1} = \frac{1}{\bar{\theta}_n} \left[ \bar{\theta}_{n-3} (1-\alpha)^3 \theta_{n-2} + \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2} \right]$$

$$\theta_{n+1} = \frac{1}{\bar{\theta}_n} \left[ \bar{\theta}_0 (1-\alpha)^n \theta_1 + \alpha R_n + \sum_{i=1}^{n-1} (1-\alpha)^i \alpha R_{n-i} \right]$$

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

we get

$$\theta_{n+1} = \frac{1}{\bar{\theta}_n} \left[ \sum_{i=0}^{n-1} (1-\alpha)^i \alpha R_{n-i} \right]$$

which is an exponential recency weighted average & as  $\theta_n$  does not depend  $\theta_1$  it is without bias



A4) In stationary, for 1000 iterations UCB and optimistic perform almost equally. But UCB performs slightly better than ~~average~~ C-greedy because of better method of exploration.

In Non-stationary, ~~UCB~~ UCB perform much better than optimistic. Since after some time optimistic becomes normal greedy but ~~doesn't~~ information gained changes with time hence it is of no-use ~~is~~ after some iterations.

UCB performs slightly better than C-greedy because of better method of exploration.