

Q3

Ex 2.3 : Which method will perform best in the long run in terms of cumulative reward & probability of selecting the best selecting the best action. Do this for  $\epsilon = 0, 0.1, 0.01$

Solution: Let  $Q_1(a) = 0 \quad \forall a$

Initially, it will choose an action with equal probability.



~~in fact, by our greedy algorithm, we get the best estimate.~~ After ~~exploring~~ choosing an action which gives the reward, it won't explore further as with the probability  $(1-\epsilon)$  i.e. 1, it will do  $\arg \max_a Q_\epsilon(a)$ !

~~When~~ When  $\epsilon = 0.01$ , algorithm will choose greedy action  $\left(1 - \epsilon + \frac{\epsilon}{K}\right) * 100\% = 99.1\%$  times

$$\text{i.e. } \left(1 - 0.01 + \frac{0.01}{10}\right) * 100\% = 99.1\%$$

while  $\epsilon = 0.1$ , algorithm will choose greedy action  $\left(1 - 0.1 + \frac{0.1}{10}\right) * 100\% = 91\%$  the time.

Hence,  $\epsilon = 0.01$  will perform better in the long run by 8.1%.



In case of  $\epsilon(t) = \frac{1}{t}$ ,

Probability of choosing a greedy action  $= 1 - \frac{1}{t} + \frac{1}{tk}$

$$= 1 - \frac{1}{t} + \frac{0.10}{t} = 1 - \frac{0.9}{t}$$

as  $t \rightarrow \infty$ ,  $P(\text{choosing greedy action}) = 1$