Reinforcement Learning
HW1
By Madhur Tandon
(2016053)

1) Comparison between Constant Step Size and Sample Averages

Analysis:

Constant Step Size gives more weight to recent rewards and thus easily adapts and surpasses the method using Sample Averages

2) Mysterious Spikes

Optimistic Initial Values encourage exploration in the initial phase as every action when explored is disappointing in some sense.

as Q(n) = Q(n-1) + alpha*(R(n) - Q(n-1))

Since Q(n-1) is large (because optimistic initial value), Q(n) decreases.

Thus, the agent chooses the action with max Q(n) which is essentially the action which hasn't been tried before.

After K turns, all possible K actions would have been chosen (on an average). Thus, in the (K+1)th turn, the action which gives the maximum reward will be chosen by the agent.

On an average, this action is the optimal action i.e. the one that gives max(q*) thereby resulting in a spike.

_Non-stationary case_:

The Mysterious Spike disappears and Optimistic Initial Value method performs equally (if not better) to Epsilon Greedy. This is because on the (K+1)th turn, the optimal action aka the one that gives max(q*) has changed since q* has changed itself.

3) Making constant step size independent of Q(0)

Q. $\beta_n = \dfrac{\alpha}{\bar{o}_n}$ and $\bar{o}_n = \bar{o}_{n-1} + \alpha(1-\bar{o}_{n-1})$ —— (A)

for $n \geq 0$ with $\bar{o}_0 = 0$

$\rule{2cm}{0.4pt}$ ①

Also, $o_n = o_{n-1} + \beta_n(R_n - o_{n-1})$ —— ②

Put ① in ②

$o_n = o_{n-1} + \dfrac{\alpha}{\bar{o}_n}(R_n - o_{n-1})$

Substituting $\bar{o}_n$ from (A)

$o_n = o_{n-1} + \dfrac{\alpha(R_n - o_{n-1})}{\bar{o}_{n-1} + \alpha(1-\bar{o}_{n-1})}$

Put $n = 1$

$\therefore \ o_1 = o_0 + \dfrac{\alpha(R_1 - o_0)}{\bar{o}_0 + \alpha(1-\bar{o}_0)}$

From (A), we know that $\bar{o}_0 = 0$

$\therefore \ o_1 = o_0 + \dfrac{\alpha(R_1 - o_0)}{0 + \alpha(1-0)}$

$o_1 = o_0 + \dfrac{\cancel{\alpha}(R_1 - o_0)}{\cancel{\alpha}}$

$o_1 = \cancel{o_0} + R_1 - \cancel{o_0}$

$\therefore \ o_1 = R_1$

Now, since $o_1 = R_1$

$\therefore \ o_1$ does not depend on $o_0$

$\Rightarrow \ o_1$ does not have initial bias

$\therefore \ o_2$ does not have initial bias (since $o_2$ depends on $o_1$)

In a similar way,

$o_n$ does not have initial bias

$\therefore \ o_n$ is an exponential recency-weighted average without initial bias.