**Ans 2.** The oscillations and spikes in the early part of the curve for the optimistic method represents that the agent explores an ~~action~~ action and then gets disappointed by the reward for that action and explore a new action for the next time step. So, after trying out all the actions, the percentage of selecting an optimal action rises. (around 35 - 40%). The spikes in the graph are caused because many bandit agents selected an optimal action at that particular time step.

**Ans 4** ~~Ex. 2.7~~ As can be seen from the graph of Average Rewards in stationary case, the slope of graph of UCB is more ~~increasing~~ than the graphs of ε-greedy & optimal initial values. for time steps greater than 1000. ~~time~~ So, the optimal action taken by UCB is more than the rest.

For non-stationary case, the percentage of optimal action is reduced to 8-10% but UCB is still the best method among the 3. ~~because for the~~ ¶ Since UCB explores more, it has higher optimal action percentage.

**Ans 3** **Ex. 2.7**

$$Q_n = Q_{n-1} + \beta_n \left[ R_{n-1} - Q_{n-1} \right]$$

$$Q_n = Q_{n-1}(1-\beta_n) + \beta_n R_{n-1}$$

~~$Q_{n-1} = (1-\beta_n)\left[ Q_{n-2}(1-\beta_{n-1}) + \beta_{n-1} R_{n-2} \right]$~~

~~$Q_{n-1} = (1-\beta_n)\left[ Q_{n-2} \right]$~~

$$Q_{n-1} = Q_{n-2}(1-\beta_{n-1}) + \beta_{n-1} R_{n-2}$$

$$Q_n = (1-\beta_n)(1-\beta_{n-1}) Q_{m-2} + (1-\beta_n)\beta_{n-1} R_{n-2} + \beta_n R_{n-1}$$

$$\vdots$$

$$Q_n = \prod_{i=1}^{n}(1-\beta_i) Q_1 + \sum_{i=1}^{n-1} \beta_{i+1} R_i \prod_{j=i+1}^{n}(1-\beta_j)$$

~~$\prod^{n} \cdots (1-\beta_n)(1-\beta_{n-1})(1-\beta$~~

Given: $\beta_n = \dfrac{\alpha}{\bar{o}_n}$

$$1-\beta_i = 1-\frac{\alpha}{\bar{o}_i} = \frac{\bar{o}_i - \alpha}{\bar{o}_i} = \frac{\bar{o}_{i-1}(1-\alpha)+\alpha-\alpha}{\bar{o}_i} = \frac{(1-\alpha)\,\bar{o}_{i-1}}{\bar{o}_i}$$

So, $\displaystyle\prod_{i=1}^{n}(1-\beta_i) = \prod_{i=2}^{n}(1-\beta_i)\left[ (1-\alpha)\frac{\bar{o}_0}{\bar{o}_1} \right]$

& given that $\bar{o}_0 = 0$ (zero)

So, $\displaystyle\prod_{i=1}^{n}(1-\beta_i) = 0$

So, $Q_n = \sum\limits_{i=1}^{n-1} \beta_{i+1} R_i \prod\limits_{j=i+1}^{n} (1 - \beta_j)$

$= \sum\limits_{i=1}^{n-1} \beta_{i+1} R_i \prod\limits_{j=i+1}^{n} (1-\alpha)\left(\dfrac{\bar{\sigma}_{j-1}}{\bar{\sigma}_j}\right)$

$= \sum\limits_{i=1}^{n-1} (1-\alpha)\, \beta_{i+1} R_i \left[\dfrac{\bar{\sigma}_{j-1}}{\bar{\sigma}_j} \times \dfrac{\bar{\sigma}_j}{\bar{\sigma}_{j+r}} \times \dots \times \dfrac{\bar{\sigma}_{n-1}}{\bar{\sigma}_n}\right]_{j=i+1}$

$= \sum\limits_{j=1}^{n-1} (1-\alpha)\, \beta_{i+1} R_i \dfrac{\bar{\sigma}_i}{\bar{\sigma}_n}$

$= \dfrac{1-\alpha}{\bar{\sigma}_n} \sum\limits_{i=1}^{n-1} \bar{\sigma}_i\, \beta_{i+1} R_i$

$= \dfrac{(1-\alpha)\alpha}{\bar{\sigma}_r} \sum\limits_{i=1}^{n-1} \dfrac{\bar{\sigma}_i}{\bar{\sigma}_{i+1}} R_i$

$= \dfrac{(1-\alpha)\alpha}{\bar{\sigma}_n} \sum\limits_{i=1}^{n-1}$