

## Exercise 2.7

$$B_n = \frac{\alpha}{O_n}$$

$$O_0 = 0.$$

$$O_1 = 0 + \alpha \cdot (1 - 0) = \alpha.$$

$$\begin{aligned} O_2 &= \alpha + \alpha(1 - \alpha) \\ &= \alpha + \alpha - \alpha^2 = 2\alpha - \alpha^2 \\ &= 1 - (1 - \alpha)^2 \end{aligned}$$

$$\begin{aligned} O_3 &= 1 - (1 - \alpha)^2 + \alpha(1 - (1 - (1 - \alpha)^2)) \\ &= 1 - (1 - \alpha)^2 + \alpha(1 - \alpha)^2 \\ &= 1 - (1 - \alpha)^2 [1 - \alpha] \\ &= 1 - (1 - \alpha)^3 \end{aligned}$$

$$\therefore O_n = 1 - (1 - \alpha)^n.$$

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

$$= \beta_n R_n + (1 - \beta_n) [\beta_{n-1} R_{n-1} + (1 - \beta_{n-1}) \cdot Q_{n-1}]$$

$$= \beta_n R_n + (1 - \beta_n) \beta_{n-1} R_{n-1} + (1 - \beta_n) \cdot (1 - \beta_{n-1}) Q_{n-1}$$

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

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

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

$$= \sum_{i=0}^{n-1} \beta_{n-i} \prod_{j=0}^i (1 - \beta_{n-j}) R_{n-i}$$

$$+ \prod_{i=1}^n (1 - \beta_i) Q_1$$

$$= \lambda + \prod_{i=1}^n \left( 1 - \frac{d}{1 - (1 - \lambda)^i} \right) Q_1$$

as  $n \rightarrow \infty$

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

when  $n = 1$ .

$$\beta_1 = \frac{\alpha}{1 - (1 - \alpha)^1} = 1$$

$\Rightarrow$  As the product  $\prod_{i=1}^n (1 - \beta_i) = 0$

hence

the term  $\gamma$  will be 0, no bias

due to  $Q_1$  disappears.



### Ex. 2.6

The early part of the graph shows spike because it is exploring during that time and sometimes it will take the ~~most~~ optimal step and sometimes it won't. But once, it gets a stable estimate, the % of ~~est~~ optimal actions grows.

To avoid spikes, we can discourage exploring and encourage the system to do more of exploitation.