

Initialize, for $a = 1$ to k :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Loop forever:

$$A \leftarrow \begin{cases} \arg \max_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$$

$$R \leftarrow \text{bandit}(A)$$

$$N(A) \leftarrow N(A) + 1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$$

$$Q_{n+1} = (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i$$

$$\sum_{n=1}^{\infty} \alpha_n(a) = \infty \quad \text{and} \quad \sum_{n=1}^{\infty} \alpha_n^2(a) < \infty$$

$$A_t \doteq \arg \max_a \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Dynamics of the MDP

$$p(s', r | s, a) \doteq \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

State-transition probabilities

$$p(s' | s, a) \doteq \Pr\{S_t = s' \mid S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a).$$

Expected rewards state-action

$$r(s, a) \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a)$$

Expected rewards state-action-state

$$r(s, a, s') \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)}.$$

State-value function for policy π

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

$$= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) \left[r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s'] \right]$$

$$= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_{\pi}(s') \right], \quad \text{for all } s \in \mathcal{S},$$

Action-value function for policy π

$$q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

Bellman optimality equation

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a)$$

$$= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$

$$q_*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

$$= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right].$$

Some equations related to Bellman

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a) + v_{\pi}(s')]$$

$$v_*(s) = \max_a \left\{ \sum_{s'} p(s' | s, a) [r(s, a) + v_*(s')] \right\}$$

$$q_{\pi}(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a) + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right]$$

$$q_*(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a) + \gamma \max_{a'} q_*(s', a') \right]$$

the step-size parameters, α_n , are not constant, then the estimate Q_n is a weighted average of previously received rewards with a weighting different from that given by Equation 2.6. What is the weighting on each prior reward for the general case, analogous to Equation 2.6, in terms of the sequence of step-size parameters?

$$\begin{aligned} Q_{n+1} &= Q_n + \alpha_n (R_n - Q_n) \\ &= \alpha_n R_n + (1 - \alpha_n) [Q_{n-1} + \alpha_{n-1} (R_{n-1} - Q_{n-1})] \\ &= \dots \\ &= \sum_{i=n}^1 \alpha_i R_i \prod_{j=i+1}^n (1 - \alpha_j) + Q_1 \prod_{j=1}^n (1 - \alpha_j) \end{aligned}$$

- If $Q_1 = 0$, is Q_n (for $n > 1$) biased? Justify your answer with brief words and equations

(b) equation 2.6

$$\begin{aligned} Q_{n+1} &= (1-\alpha) Q_n + \frac{\alpha}{n+1} (R_n - Q_n) \\ \text{when } Q_1 &= 0 \\ Q_{n+1} &= \frac{\alpha}{n+1} (1-\alpha)^n R_n \\ \Rightarrow E(Q_{n+1}) &= E\left(\frac{\alpha}{n+1} (1-\alpha)^n R_n\right) \\ &= \frac{\alpha}{n+1} (1-\alpha)^n E(R_n) \\ \Rightarrow \text{when } \frac{\alpha}{n+1} (1-\alpha)^n &= 1 \\ E(Q_{n+1}) &= E(R_n) \\ \Rightarrow \text{equation is unbiased} \end{aligned}$$

- Derive condition(s) for Q_1 for when Q_n will be unbiased.

(c) Q_n will be unbiased if $\begin{cases} Q_1 = 0 \\ \sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} = 1 \end{cases}$

- Show that Q_n is an unbiased estimator as n , (which is often referred to as asymptotically unbiased).

(d) $Q_{n+1} = (1-\alpha) Q_n + \frac{\alpha}{n+1} (R_n - Q_n)$

$$\begin{aligned} \lim_{n \rightarrow \infty} Q_{n+1} &= \lim_{n \rightarrow \infty} (1-\alpha) Q_n + \lim_{n \rightarrow \infty} \frac{\alpha}{n+1} (R_n - Q_n) \\ &= \lim_{n \rightarrow \infty} \alpha \cdot \frac{1 - (1-\alpha)^{n+1}}{1 - (1-\alpha)} \cdot R_i \\ &= \lim_{n \rightarrow \infty} \alpha \cdot \frac{1 - (1-\alpha)^{n+1}}{\alpha} \cdot R_i \\ &= \lim_{n \rightarrow \infty} 1 - (1-\alpha)^{n+1} \cdot R_i \\ &= R_i \\ \Rightarrow \lim_{n \rightarrow \infty} E(Q_{n+1}) &= E(R_i) \end{aligned}$$

- Why should we expect that the exponential recency-weighted average will be biased in practice?

For α , we are able to control the learning rate by modifying α , the learning process could be speed up by increasing α and vice versa. For Q_1 , the agent could start explore more in the beginning rather than always choosing the best choice by increase Q_1 .

- Show that in the case of two actions, the soft-max distribution is the same as that given by the logistic, or sigmoid, function often used in statistics and artificial neural networks

$$Q_i = \frac{\exp(H_i(a))}{\sum_{b=1}^2 \exp(H_i(b))}$$

case of two action:

$$\begin{aligned} \int_0^1 \frac{\exp(H_i(a))}{\exp(H_i(a)) + \exp(H_i(b))} &= \frac{1}{1 + \exp(H_i(b) - H_i(a))} \\ \text{if } Q &= - (H_i(b) - H_i(a)) \\ \text{then } \int_0^1 \frac{1}{1 + \exp(-Q)} &= \sigma(Q) \rightarrow \text{sigmoid} \end{aligned}$$

- UCB also produce spikes in the very beginning in both the two reproduced figures. Explain in your own words why the spikes appear (both the sharp increase and sharp decrease)

- In the optimistic initialization, the first few steps are not really random. Instead, we loop through all the actions multiple times, picking our large optimistic value at random. In the first round, all actions have equal optimization percentages (10%). The next action will have a spike, because on average the best action will have the largest value. This will be repeated with decreasing values until the effect of the optimal initial value fades away.