

# HW1

Theo McArn

January 2026

## Question 1: 10-Armed Bandit Testbed

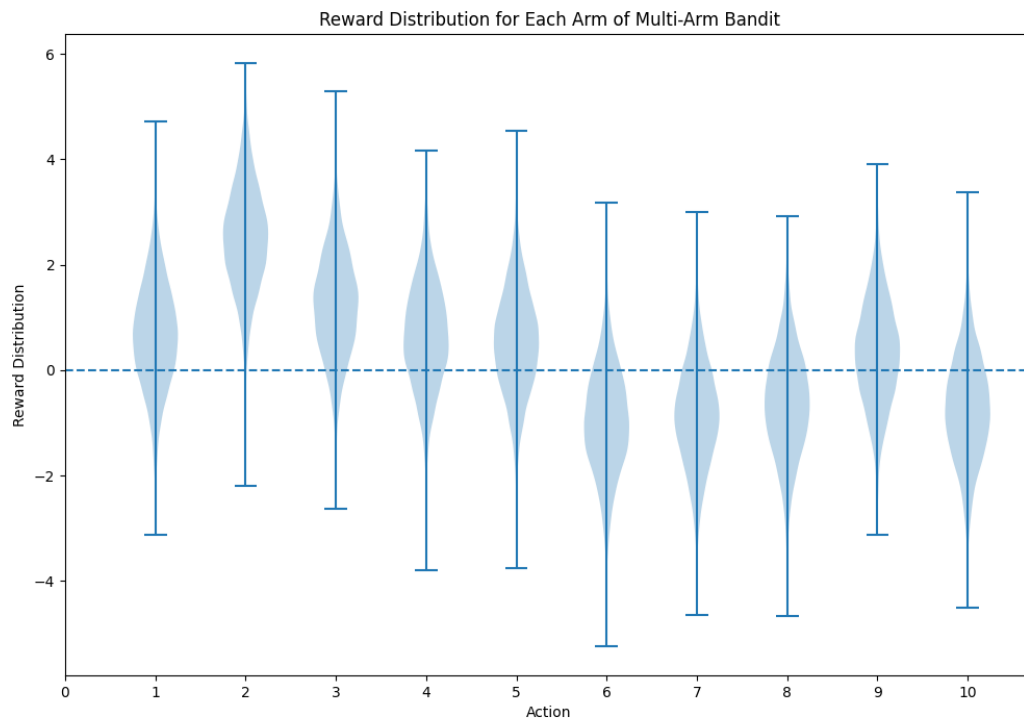


Figure 1: Reward Distribution for each Bandit

## Question 2: $\epsilon$ -Greedy

When  $\epsilon = 0$ , there is no exploration. As a result the agent quickly converges on a bandit that is suboptimal. In this case, this results in an average reward of about 1.

When  $\epsilon = 0.01$ , there is a slight amount of exploration. 1% of actions are random. As a result, the average reward seems to steadily increase as more exploration happens and more optimal bandits are discovered. However, because the exploration is so slow, the average rewards per step does not seem to converge by the end of 1000 time steps. However, at the end of 1000 steps, the average reward is about 1.3.

When  $\epsilon = 0.1$ , there is a lot of exploration. 10% of actions are random. As a result, the average reward seems to quickly increase as more exploration happens and more optimal bandits are discovered. And it seems to converge with an average reward of about 1.4 per step once it has found the best bandit.

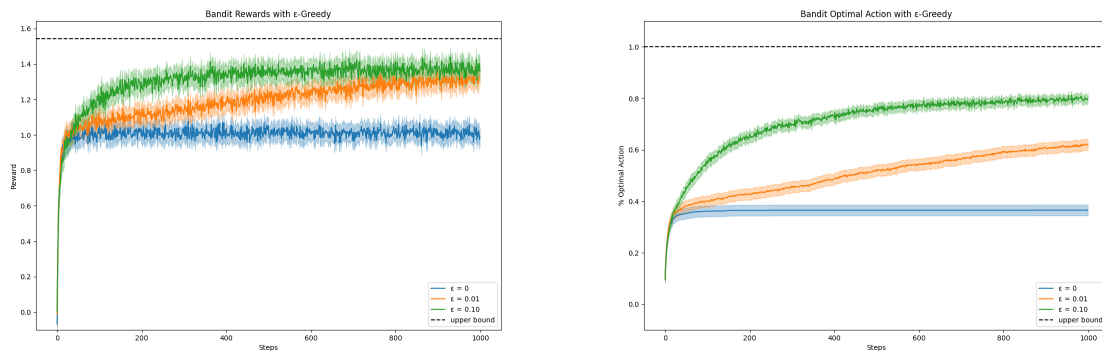


Figure 2: Average Reward and Percentage of Optimal Action for  $\epsilon$ -Greedy Agents

### Question 3: $\epsilon$ -Greedy vs. UCB

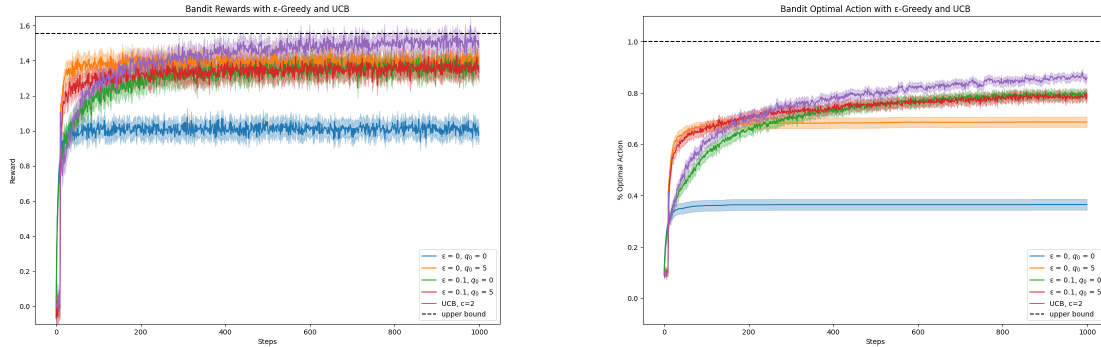


Figure 3: Average Reward and Percentage of Optimal Action for  $\epsilon$ -Greedy Agents and UCB Agents

#### Why spikes appear in UCB and Optimistically Initialized Agents:

For both UCB agents and optimistically initialized agents, every bandit is tried once before returning to any duplicate actions. As a result, for the first 10 timesteps, every bandit is tested, regardless of what the previous rewards were. After 10 steps, the "best" action is then chosen based on the bandit that had the highest reward after the first 10 actions. This chosen action represents the "up" part of the spike. However, after this selection, the UCB score for this action is reduced because it has been visited twice and the scores for every other action is increased uniformly since they have only been visited once. As a result, this shift in UCB scores outweighs the Q-Values and as a result, the agent selects a new action that has a lower expected reward than the initial choice at step 10. Taking this action with a lower expected reward results in the "down" component of the spike.

This hypothesis is further supported by figure 3 because you can see from the rewards plot that this spike (in purple) happens roughly at timestep 10.

## Varying Step-Size Weights

→ Non-constant  $\alpha$

$$\begin{aligned}
 Q_{n+1} &= Q_n + \alpha_n [R_n - Q_n] \\
 &= \alpha_n R_n + (1 - \alpha_n) Q_n \\
 &= \alpha_n R_n + (1 - \alpha_n) [\alpha_{n-1} R_{n-1} + (1 - \alpha_{n-1}) Q_{n-1}] \\
 &= (1 - \alpha_n)(1 - \alpha_{n-1}) Q_{n-1} + (1 - \alpha_n) \alpha_{n-1} R_{n-1} + \alpha_n R_n \\
 &= (1 - \alpha_n)(1 - \alpha_{n-1}) [\alpha_{n-2} R_{n-2} + (1 - \alpha_{n-2}) Q_{n-2}] \\
 &\quad + (1 - \alpha_n) \alpha_{n-1} R_{n-1} + \alpha_n R_n \\
 &= (1 - \alpha_n)(1 - \alpha_{n-1})(1 - \alpha_{n-2}) Q_{n-3} + \\
 &\quad (1 - \alpha_n)(1 - \alpha_{n-1}) \alpha_{n-2} R_{n-2} + \\
 &\quad (1 - \alpha_n) \alpha_{n-1} R_{n-1} + \alpha_n R_n \\
 &= \left[ \prod_{i=0}^n (1 - \alpha_{n-i}) \right] Q_1 + \sum_{i=1}^n \left[ \alpha_i \prod_{j=i+1}^n (1 - \alpha_j) \right] R_i
 \end{aligned}$$

$$w_i = \alpha_i \prod_{j=i+1}^n (1 - \alpha_j)$$

Figure 4: Derivation of Reward Weighting for Non-Constant Step Sizes

## Bias in Q-Value Estimates

1. The sample-average estimate of  $Q_n$  is unbiased.

Sample Average Estimate

$$Q_n = \frac{R_1 + R_2 + \dots + R_{n-1}}{n-1}$$

$$\begin{aligned} E[Q_n] &= \frac{1}{n-1} E[R_1 + R_2 + \dots + R_{n-1}] \\ &= \frac{1}{n-1} \left[ E[R_1] + E[R_2] + \dots + E[R_{n-1}] \right] \end{aligned}$$

$$E[R_i] = q_*$$

$$E[Q_n] = \frac{1}{n-1} \cdot [n \cdot q_*]$$

$$E[Q_n] = \frac{n-1}{n-1} q_*$$

$$E[Q_n] = q_* \text{ (if } n \geq 2)$$

∴ unbiased

2. For an exponential recency-weighted average estimate, if  $Q_1 = 0$ ,  $Q_n$  is biased when  $n > 1$ .

Exponentially Recency-Weighted Average

$$Q_n = \cancel{(1-\alpha)^n} Q_1 + \sum_{i=0}^n \alpha (1-\alpha)^{n-i} R_i \quad Q_1 = 0$$

$$E[Q_n] = q_* \cdot \sum_{i=0}^n \alpha (1-\alpha)^{n-i}$$

$$E[Q_n] \neq q_* \quad \therefore \text{biased estimate}$$

3. The exponential recency-weighted average estimate will be unbiased only when  $Q_1 = q_*$

Conditions for Unbiased Estimator  $E[Q_n] = q_*$

$$Q_n = (1-\alpha)^n Q_1 + \sum_{i=0}^{n-1} \alpha (1-\alpha)^{n-i-1} R_{i+1} \quad Q_1 = 0$$

$$E[Q_n] = E[(1-\alpha)^n Q_1] + \alpha q_* \sum_{i=0}^{n-1} (1-\alpha)^{n-i-1} = q_*$$

using geometric series

$$q_* = (1-\alpha)^n Q_1 + q_* [1 - (1-\alpha)^n]$$

$$q_* = (1-\alpha)^n Q_1 + q_* - q_* (1-\alpha)^n$$

$$0 = (1-\alpha)^n (Q_1 - q_*)$$

$$Q_1 = q_* \quad \text{or} \quad \alpha = 1 \quad \text{violates initial conditions}$$

$$E[Q_n] = q_* \quad (\text{unbiased})$$

iff  $Q_1 = q_*$

4. For an exponential recency-weighted average estimate,  $Q_n$  is asymptotically unbiased as  $n \rightarrow \infty$

Show  $Q_n$  is asymptotically unbiased:

$$Q_n = (1-\alpha)^n Q_1 + \sum_{i=0}^{n-1} \alpha (1-\alpha)^{n-i} R_i$$

$$\lim_{n \rightarrow \infty} E[Q_n] = \lim_{n \rightarrow \infty} \left[ \cancel{(1-\alpha)^n} Q_1 + q_* \alpha \sum_{i=0}^{n-1} (1-\alpha)^{n-i} \right]$$

since  $\alpha < 1$   $\lim_{n \rightarrow \infty} (1-\alpha)^n = 0$  geometric series

$$= \lim_{n \rightarrow \infty} \left[ q_* \left[ 1 - \cancel{(1-\alpha)^n} \right] \right]$$

$$\lim_{n \rightarrow \infty} E[Q_n] = q_* \rightarrow \therefore \text{asymptotically unbiased}$$

5. **Why should we expect that the exponential recency-weighted average will be biased in general?** In general, we can expect that the exponential recency-weighted average will be biased because it is always placing more weight on the current and recent samples and less weight on the older observations. Inherently, this difference in weighting the rewards from the same distribution results in bias.