

Bandits

Exercise 1

- (a) 50% is the probability for choosing the greedy action and from the other 50%, 50% is also the probability for the greedy, because we have two actions selected at random.

$$(1 - 0.5) + 0.5 * 0.5 = 0.75$$

- (b) At the beginning, every choice is non-epsilon, but $A_1 = 1$ was made.

$$Q_2(1) = 1, Q_2(2) = Q_2(3) = Q_2(4) = 0 \text{ but chosen action was } A_2 = 2$$

$$Q_3(1) = Q_3(2) = 1, Q_3(3) = Q_3(4) = 0 \text{ both } A_3 = 1 \text{ and } A_3 = 2 \text{ can be chosen}$$

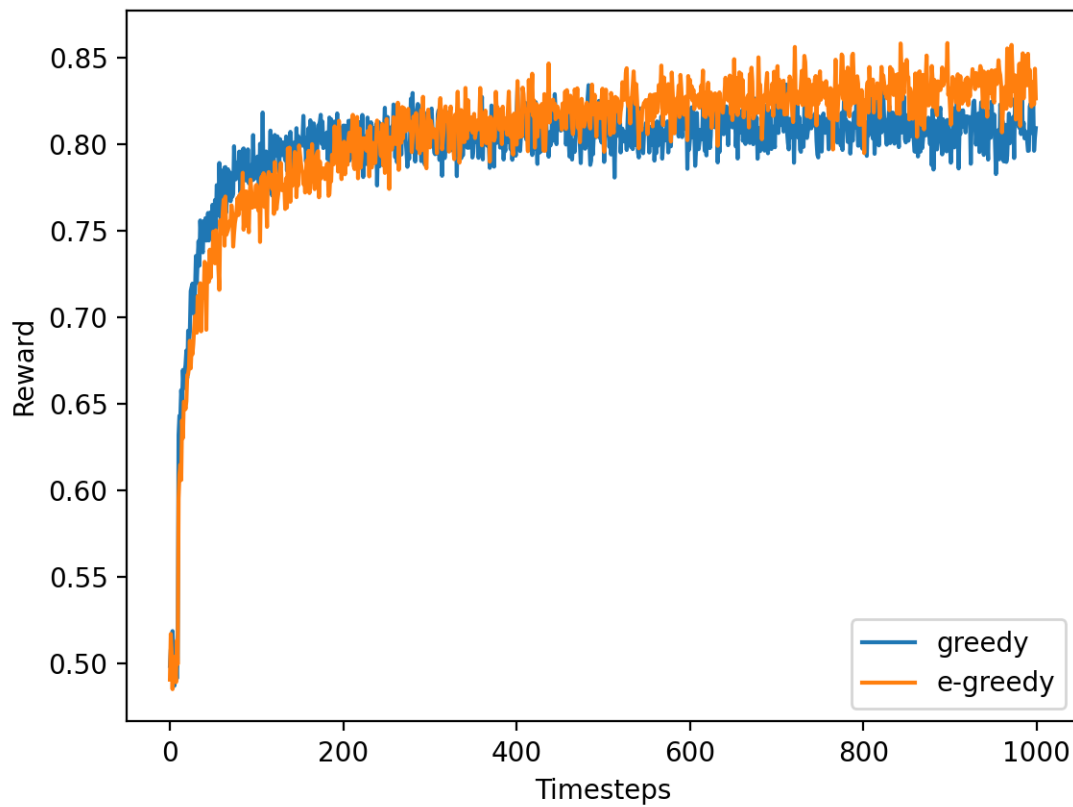
$$Q_5(1) = 1, Q_5(2) = 5/3 = 1.7 \text{ and } Q_5(3) = Q_5(4) = 0 \text{ but chosen action was } A_5 = 3$$

- (1) Definitely, choices 2 and 5 were ϵ choices

- (2) Choices 1 and 3 could have been ϵ choices

Exercise 2

- (c) The ϵ -greedy method improves slower, but to a better average reward:



- (d) We can use initialization with optimistic values to force exploration on the initial stages. Also Upper-Confidence-Bound action selection is shown to perform better than ϵ -greedy search, because it takes into account the uncertainty of the value of the chosen action - more frequently chosen actions have lower uncertainty.