

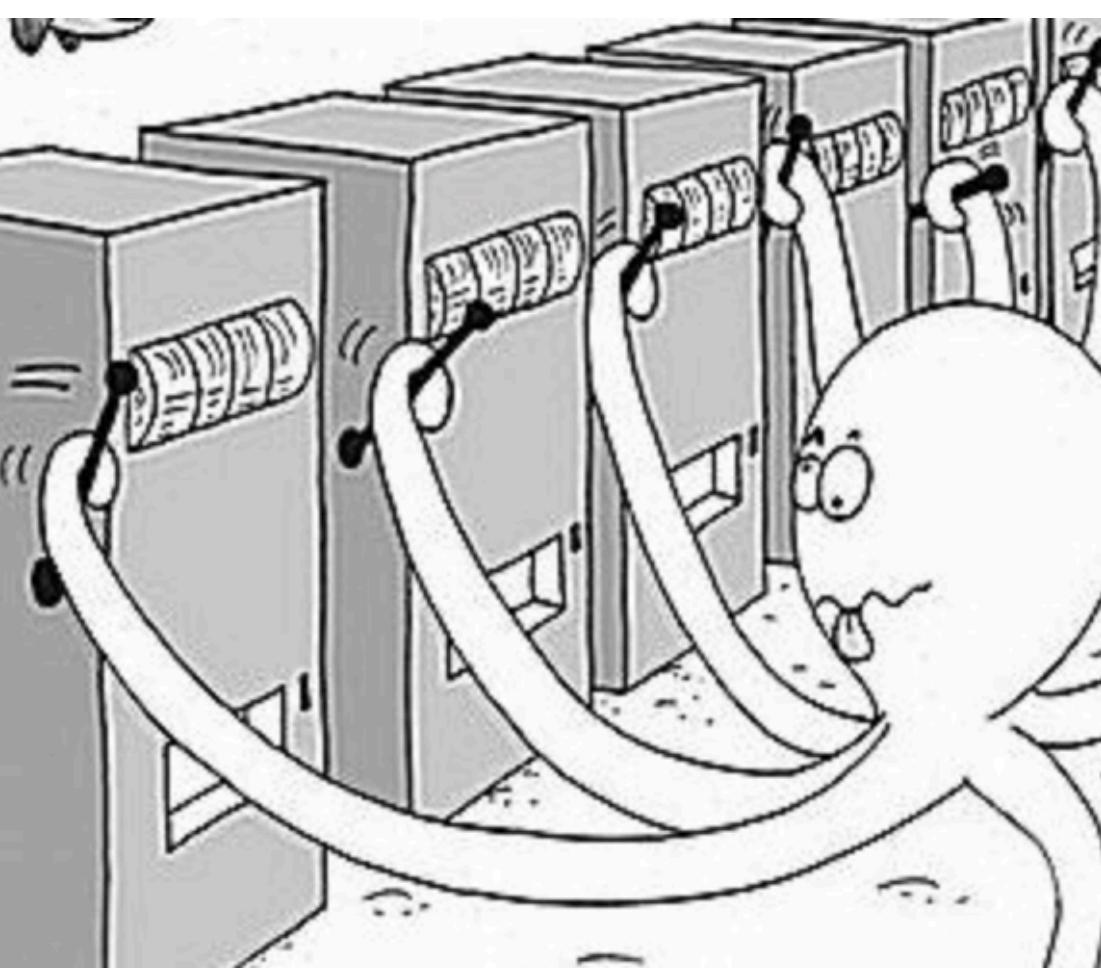
An introduction to Reinforcement Learning

28th of June 2022

Recap: Multi-armed bandits

Greedy action selection:

$$P(a_t = a) = \begin{cases} 1 & \text{if } a_t = \operatorname{argmax}_a V_t(a) \\ 0 & \text{otherwise} \end{cases}$$



Softmax action selection:

$$P(a_t = a) = \frac{e^{V_t(a) \cdot \beta}}{\sum_{i=1}^N e^{V_t(a_i) \cdot \beta}}$$

Action is governed by a **policy**:

$$\pi(a, s) = P(a_t = a | s_t = s)$$

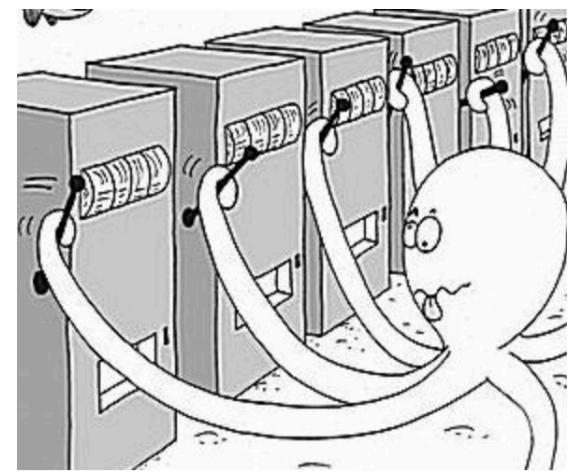
Epsilon-greedy action selection:

$$P(a_t = a) = \begin{cases} 1 - \epsilon & \text{if } a_t = \operatorname{argmax}_a V_t(a) \\ \epsilon/N & \text{otherwise} \end{cases}$$

Upper-confidence-bound (UCB) action selection:

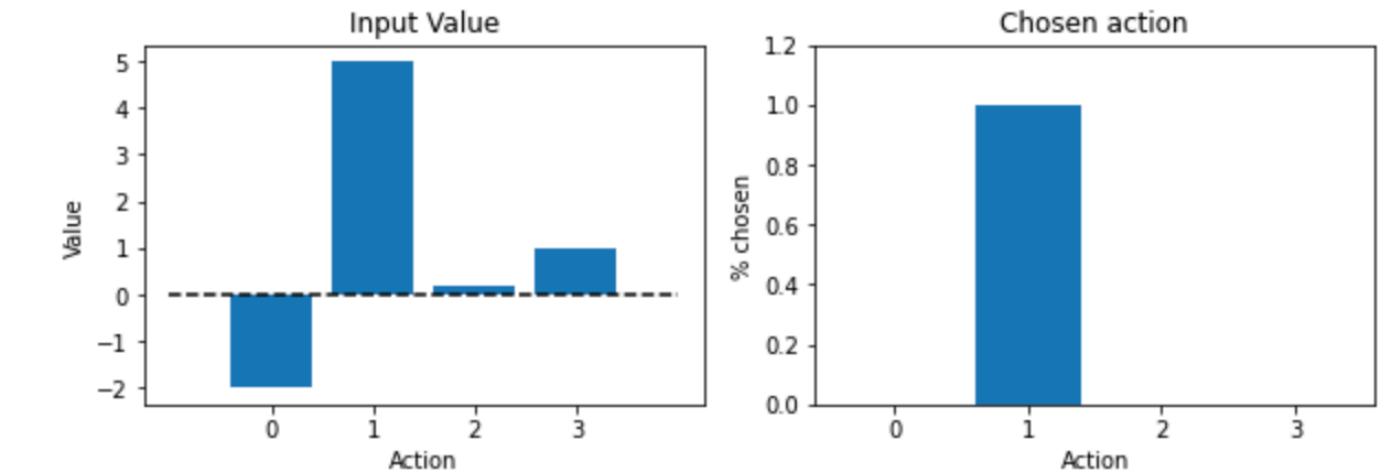
$$P(a_t = a) = \operatorname{argmax}_a [V_t(a) + c \cdot \sqrt{\frac{\ln t}{N_t(a)}}]$$

Performance comparison

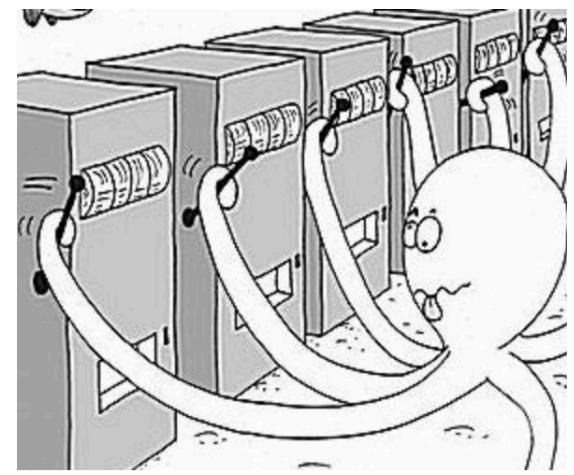


N arms = 10
N simulations = 1000

Greedy action selection:
$$P(a_t = a) = \begin{cases} 1 & \text{if } a_t = \operatorname{argmax}_a V_t(a) \\ 0 & \text{otherwise} \end{cases}$$



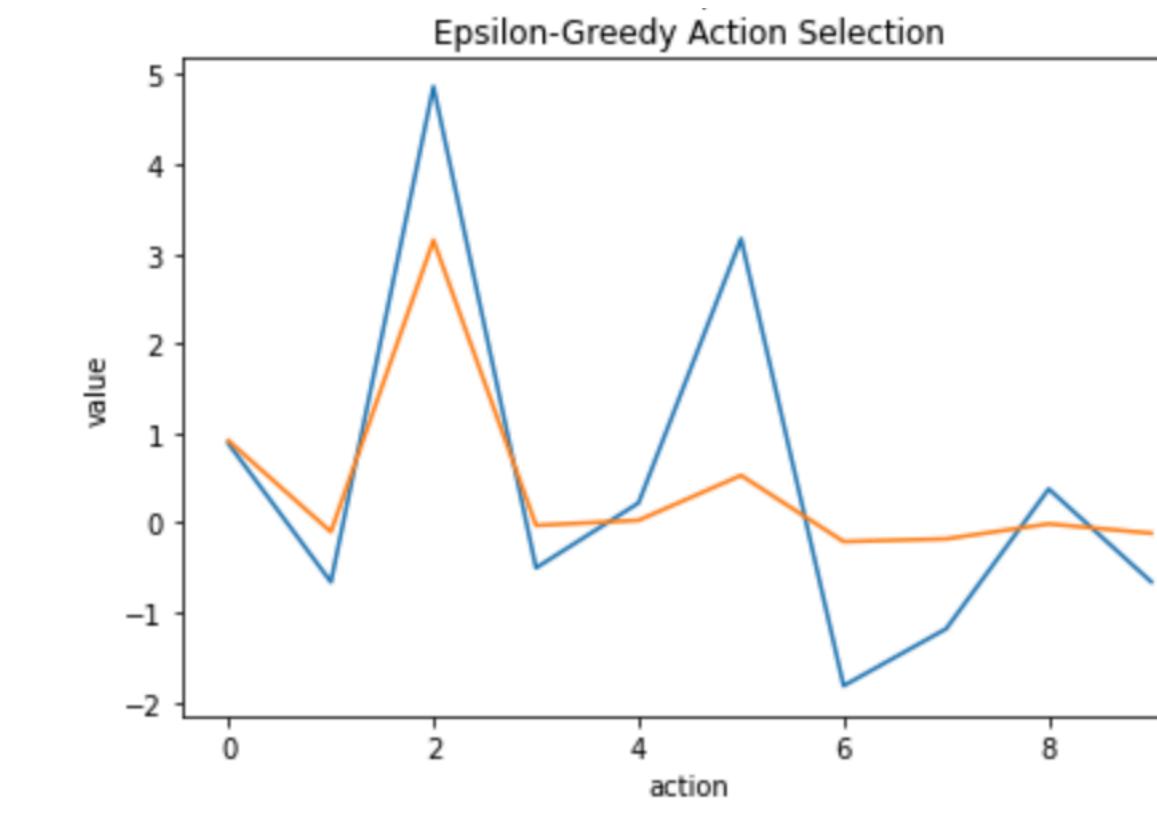
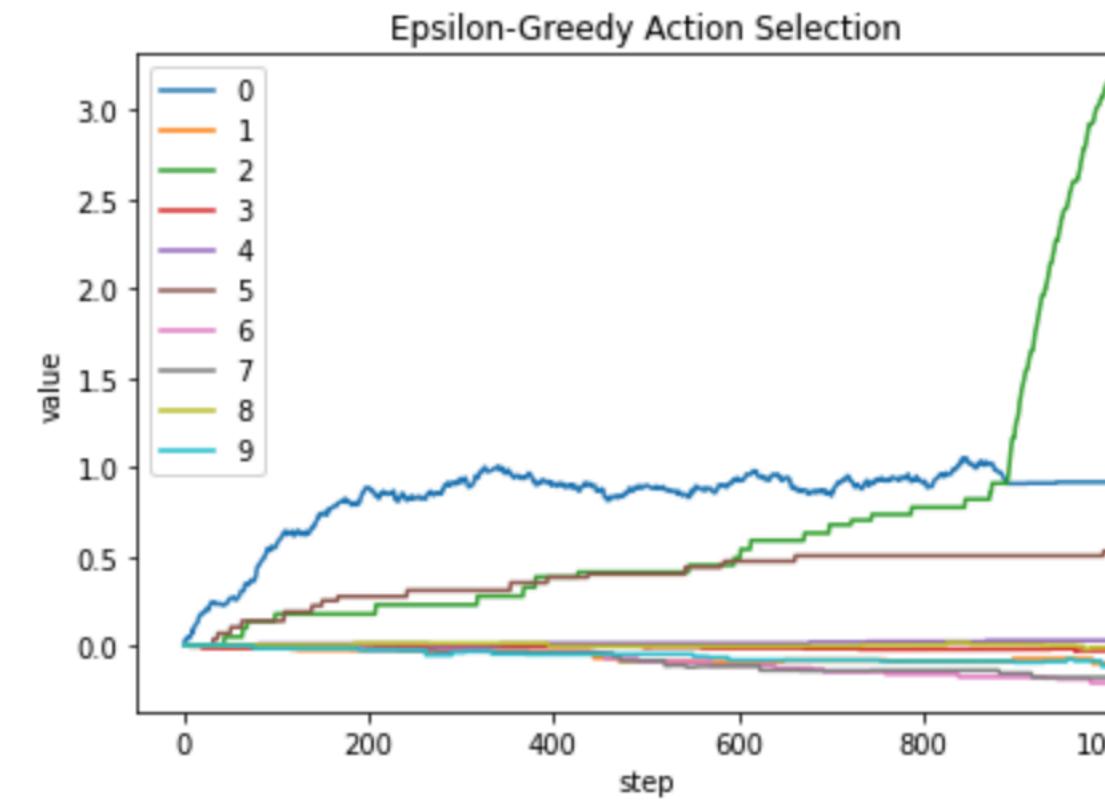
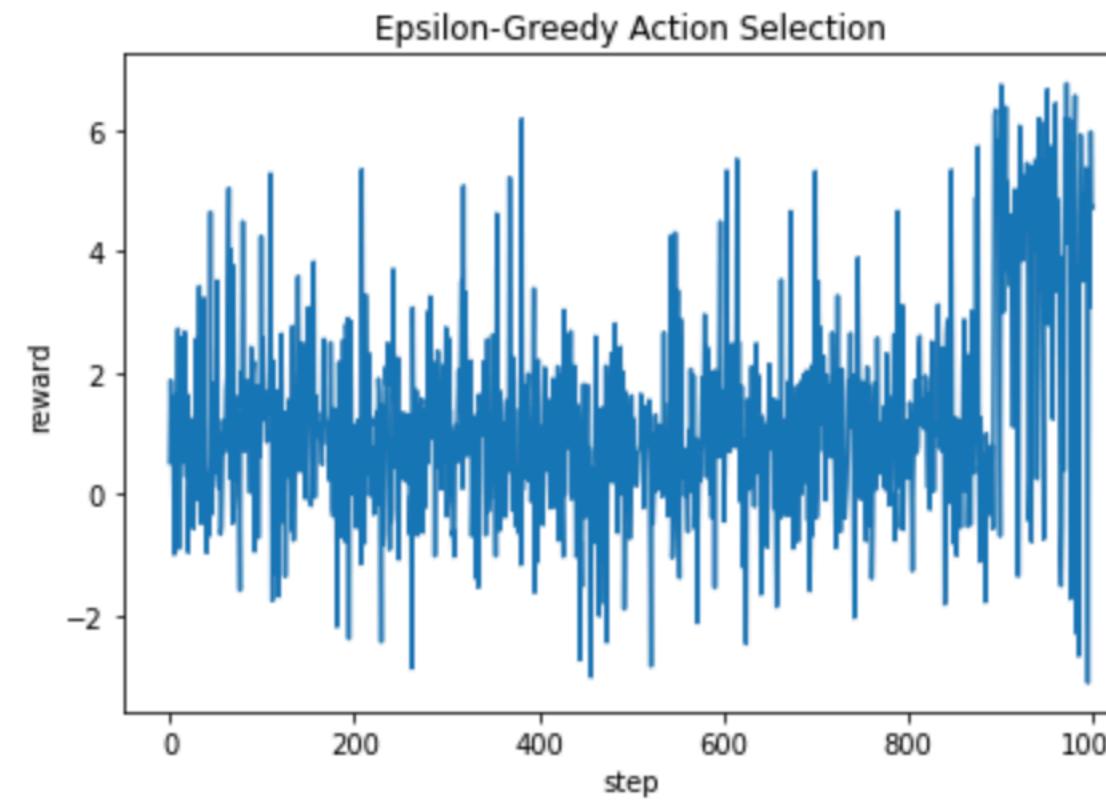
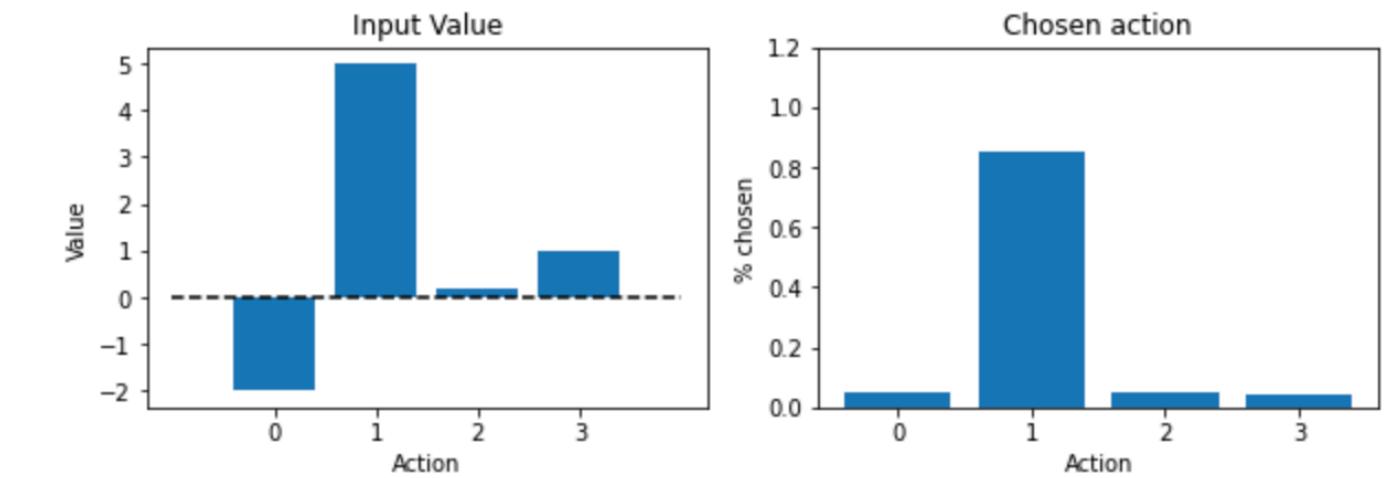
Performance comparison



N arms = 10
N simulations = 1000

Epsilon-greedy action selection:

$$P(a_t = a) = \begin{cases} 1 - \epsilon & \text{if } a_t = \operatorname{argmax}_a V_t(a) \\ \epsilon/N & \text{otherwise} \end{cases}$$

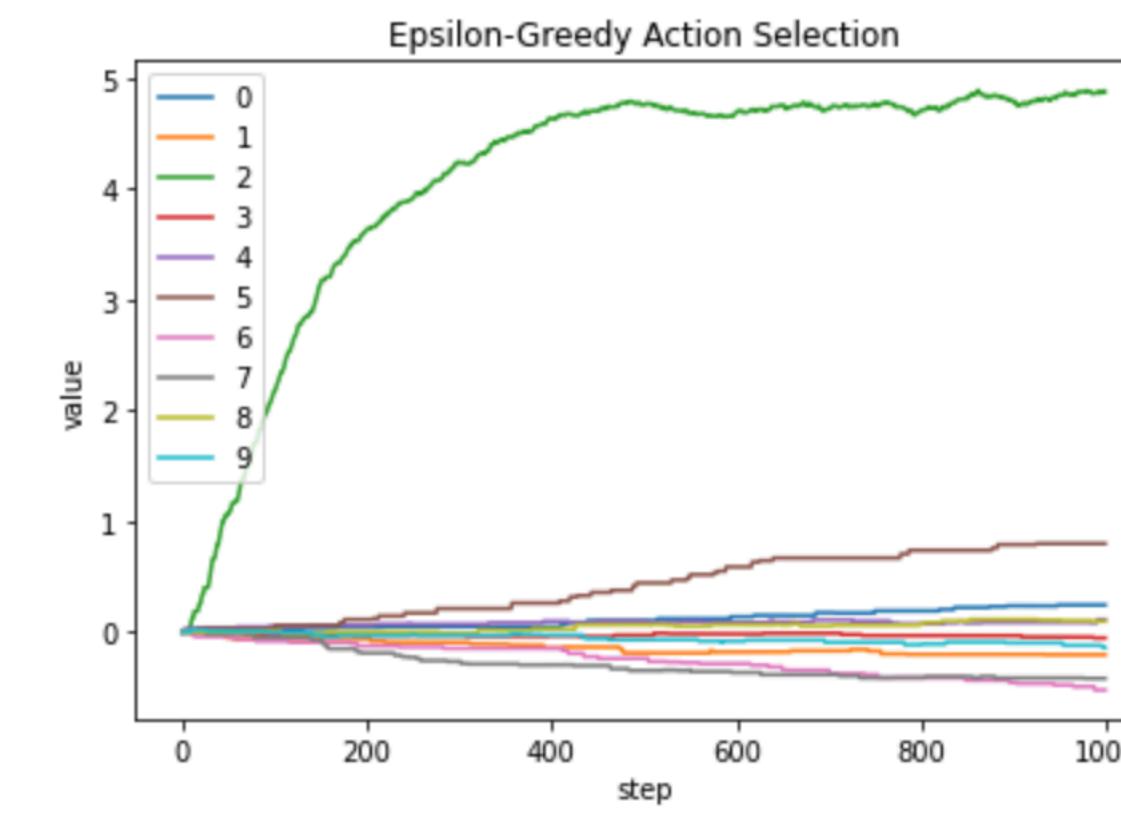
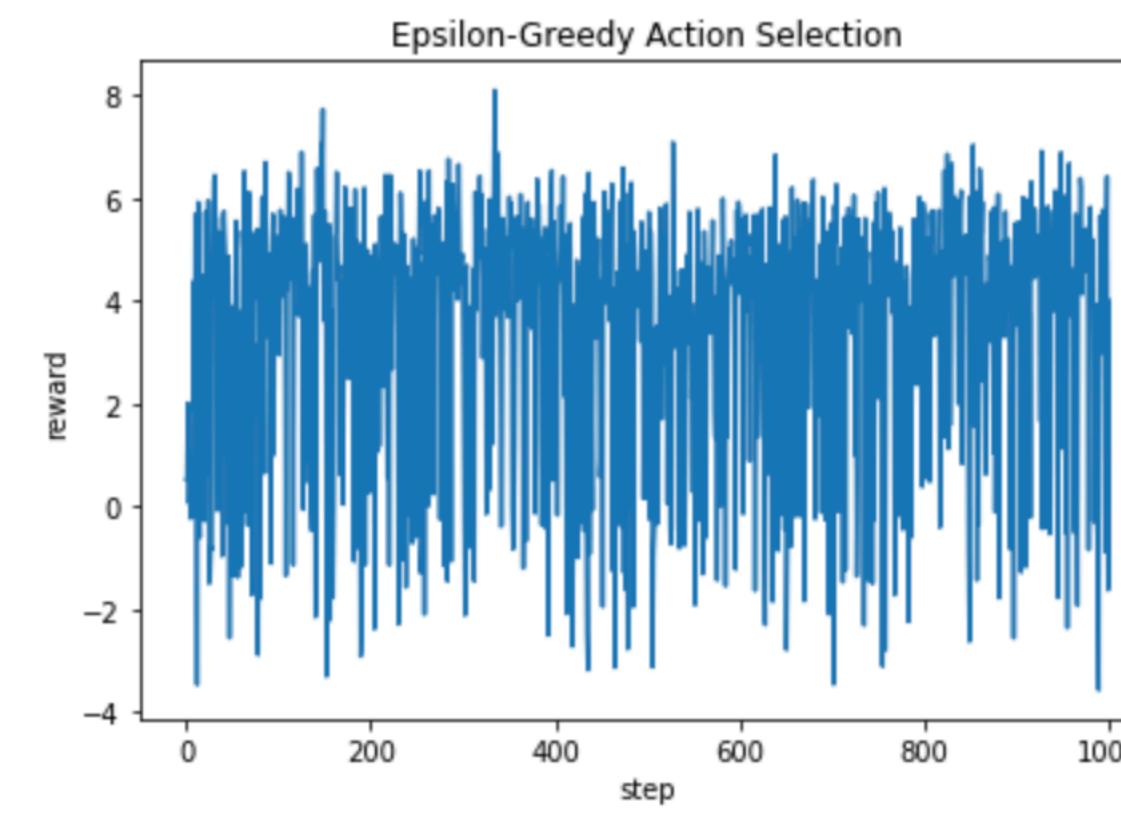


How will your choice of ϵ affect these plots? (Here: $\epsilon = 0.15$)

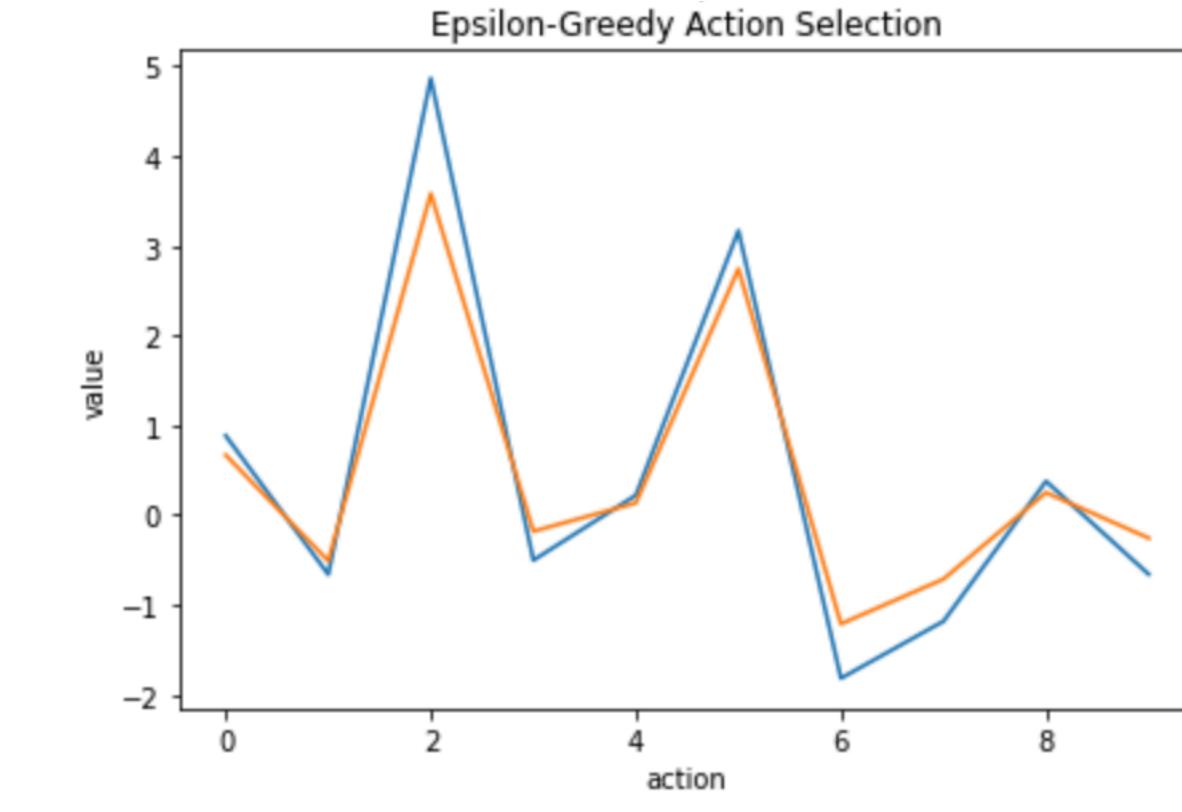
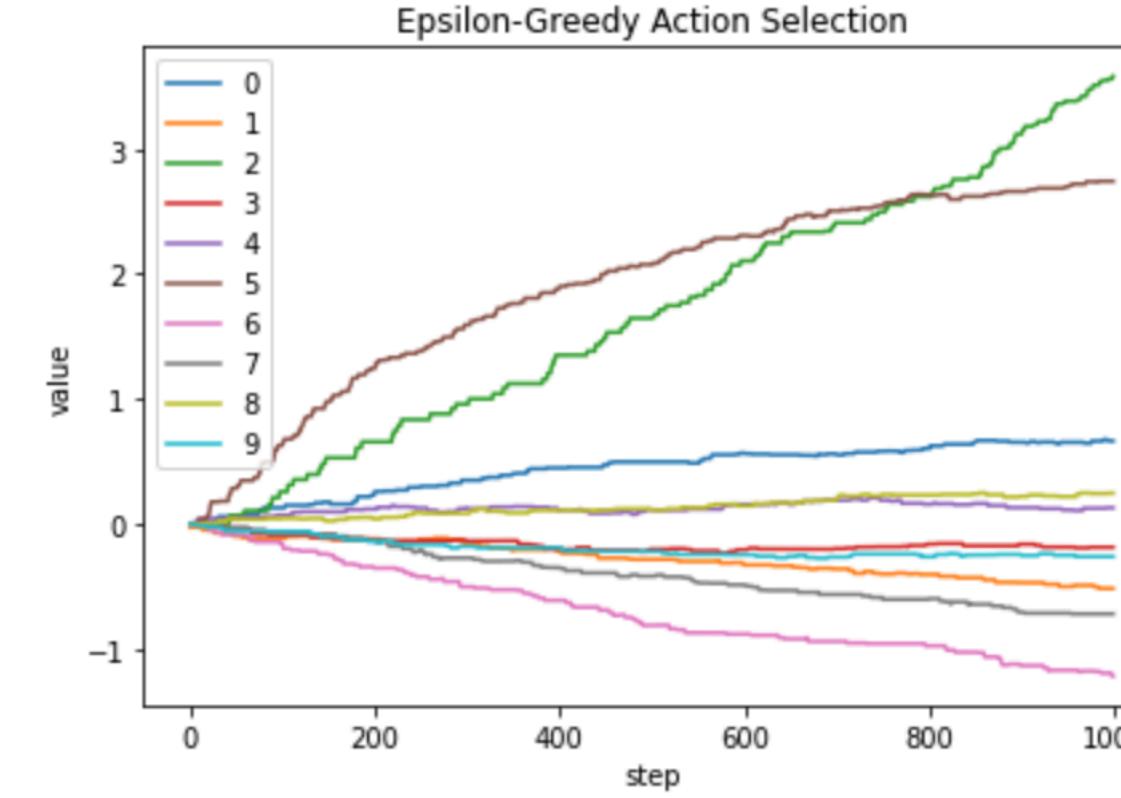
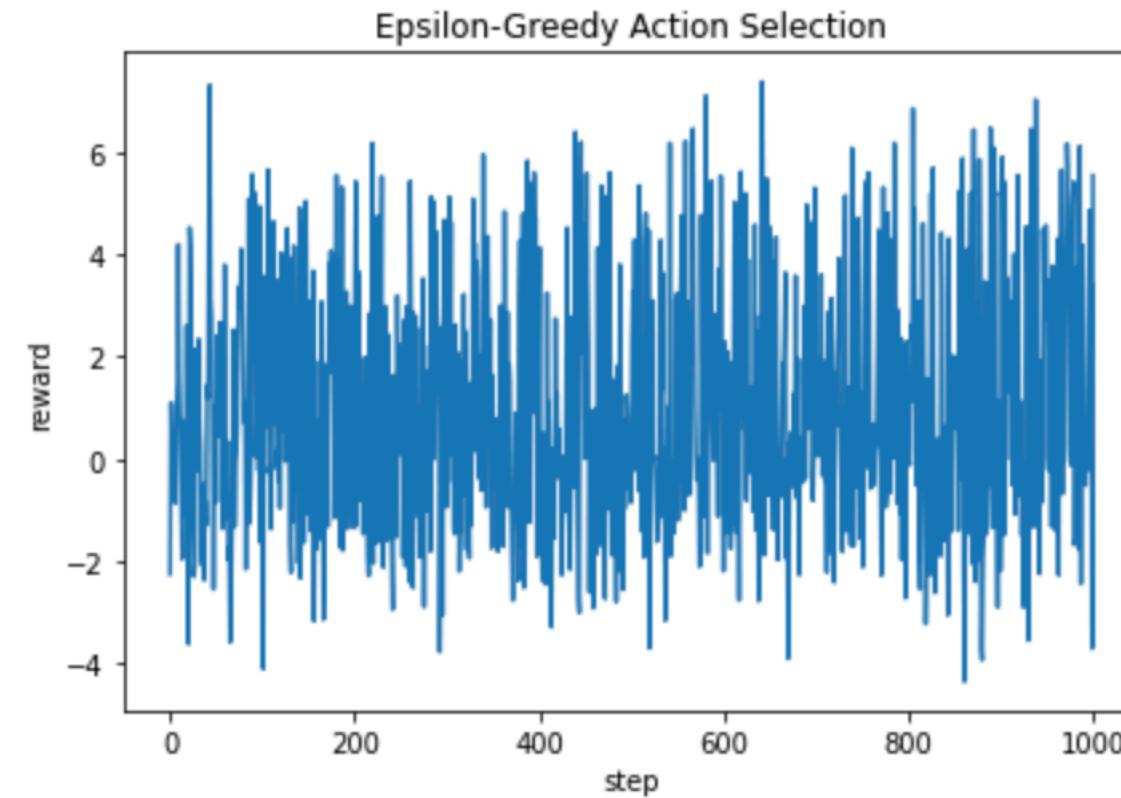
Perfomance comparison

$$P(a_t = a) = \begin{cases} 1 - \epsilon & \text{if } a_t = \operatorname{argmax}_a V_t(a) \\ \epsilon/N & \text{otherwise} \end{cases}$$

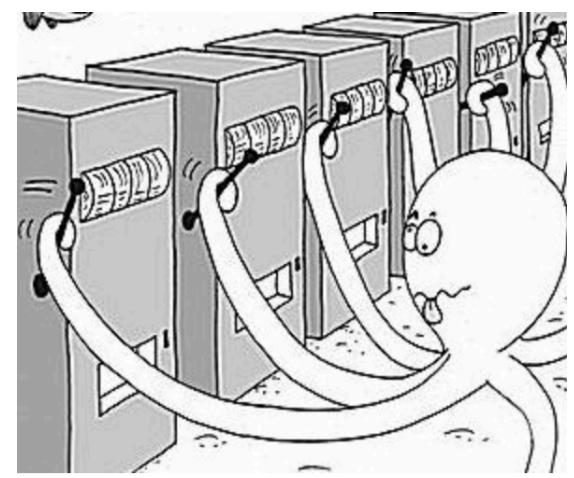
$\epsilon = 0.35$



$\epsilon = 0.85$



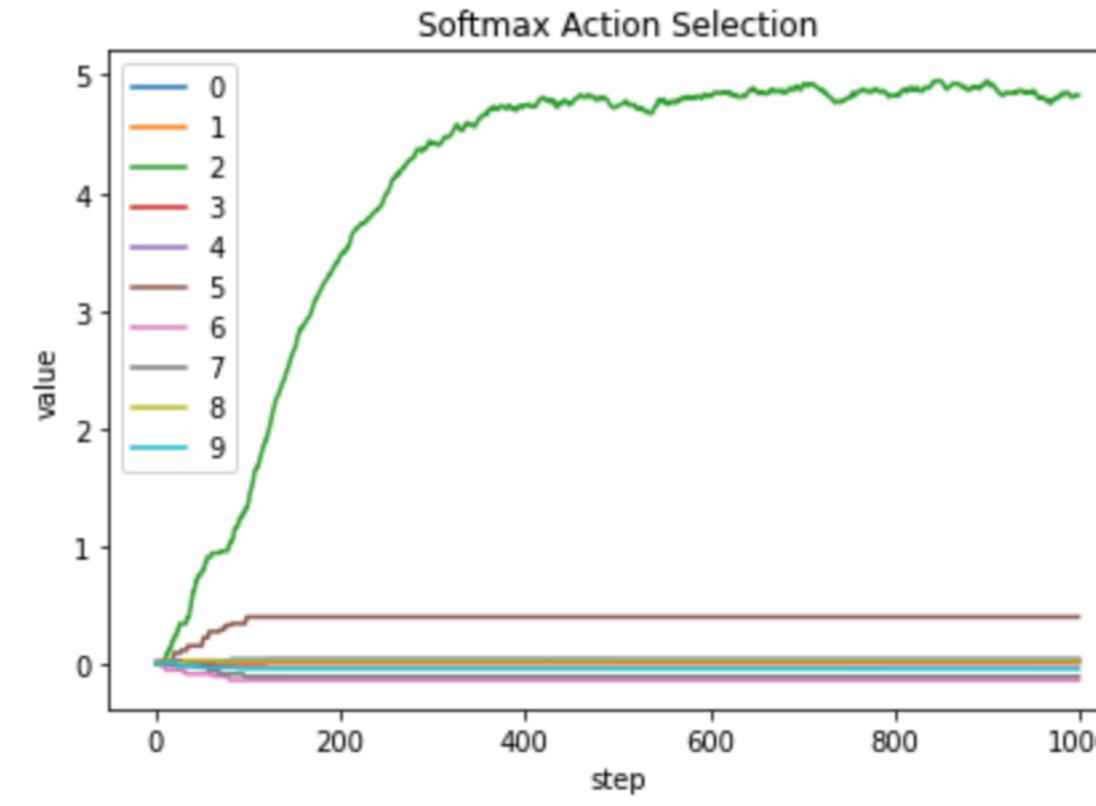
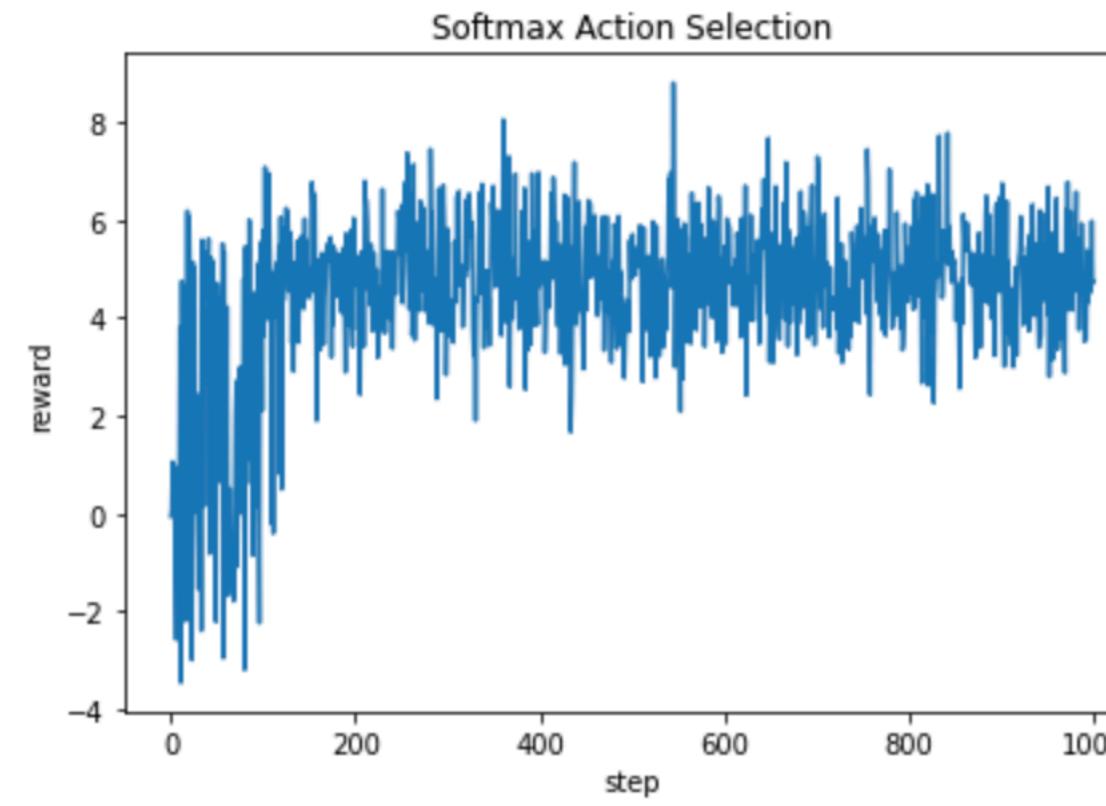
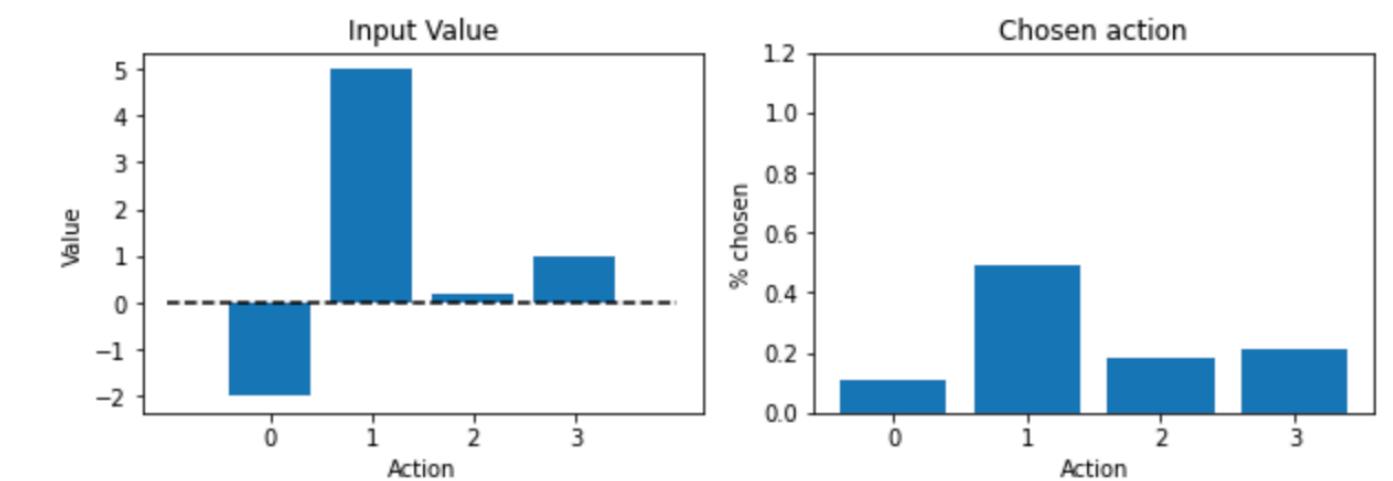
Performance comparison



N arms = 10
N simulations = 1000

Softmax action selection:

$$P(a_t = a) = \frac{e^{V_t(a) \cdot \beta}}{\sum_{i=1}^N e^{V_t(a_i) \cdot \beta}}$$

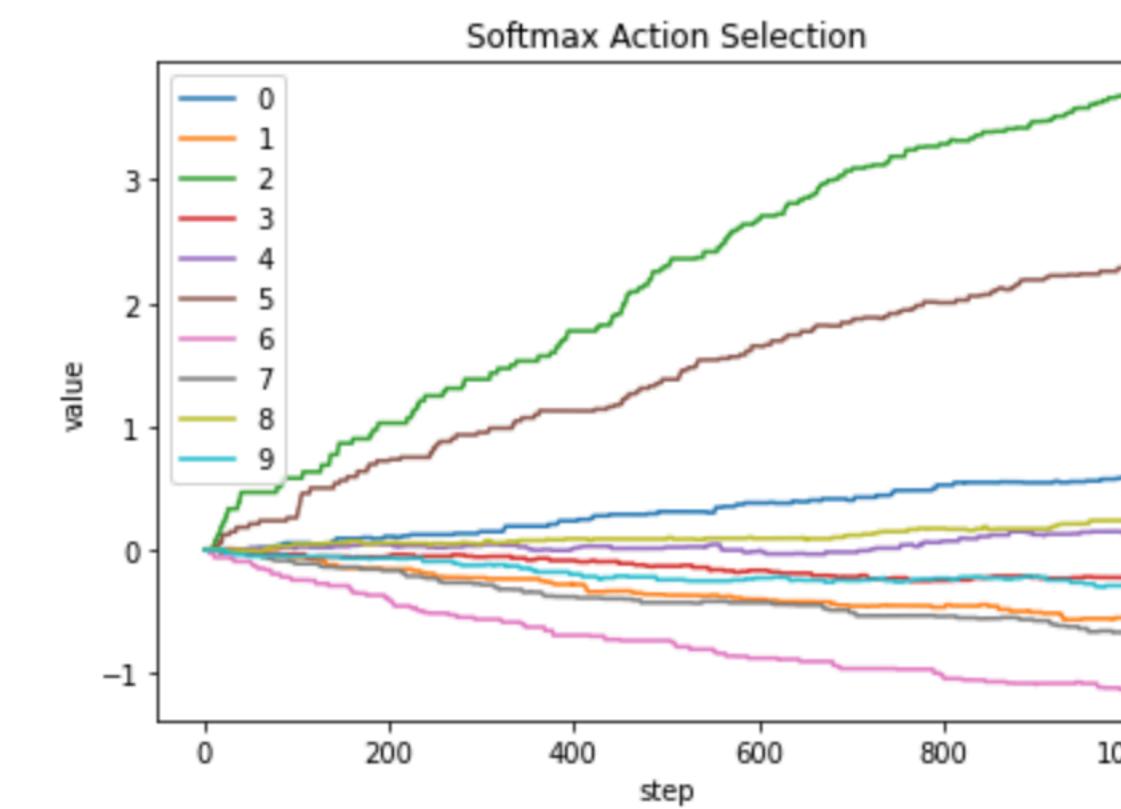
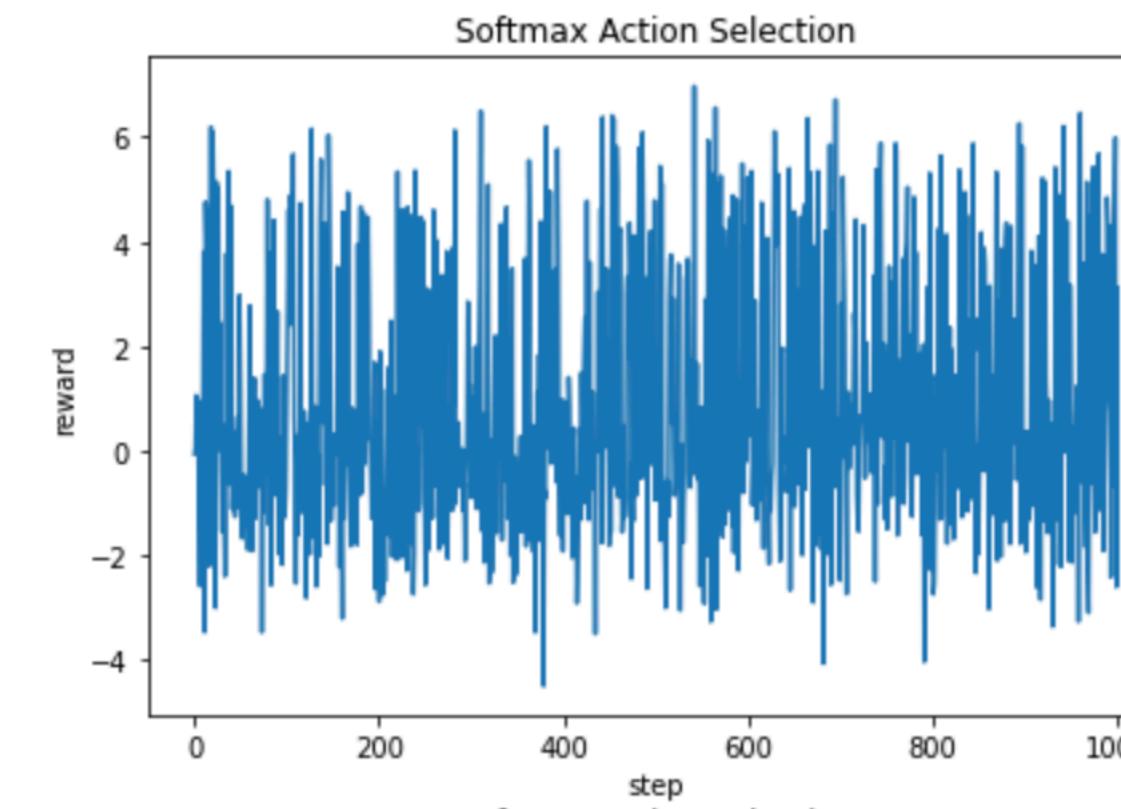


How will your choice of β affect these plots? (Here: $\beta = 2$)

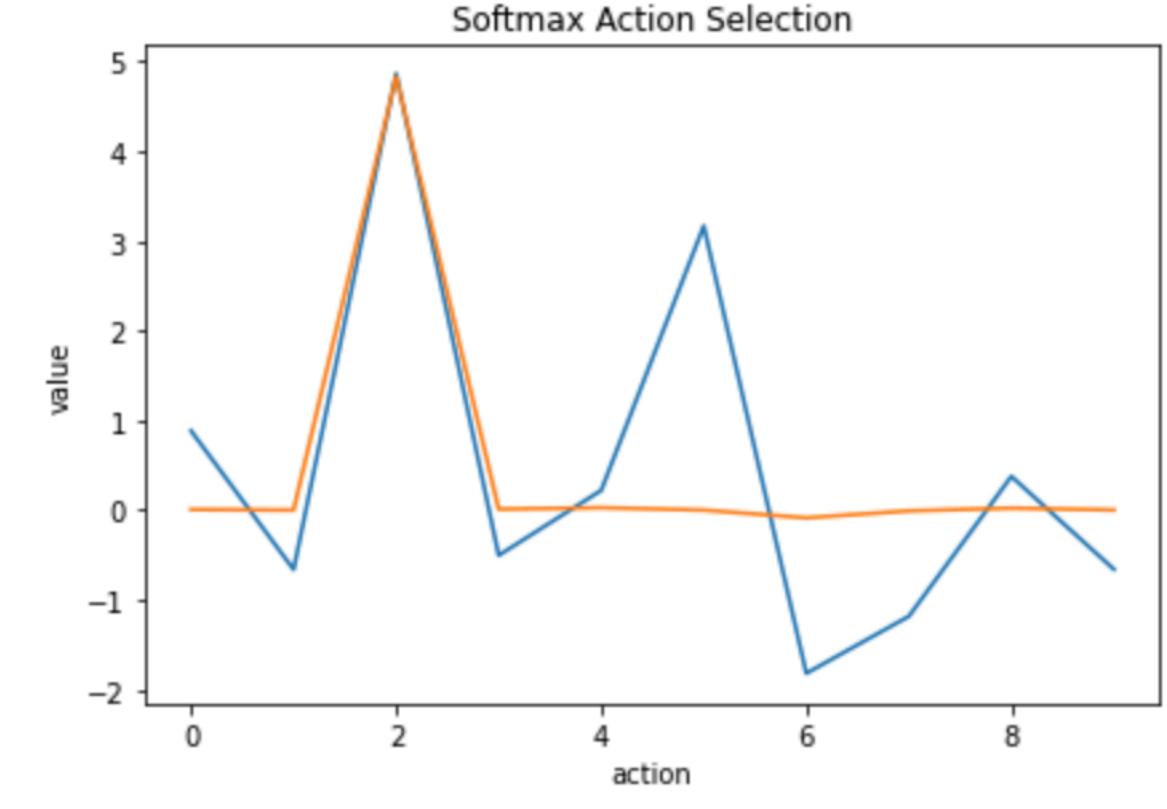
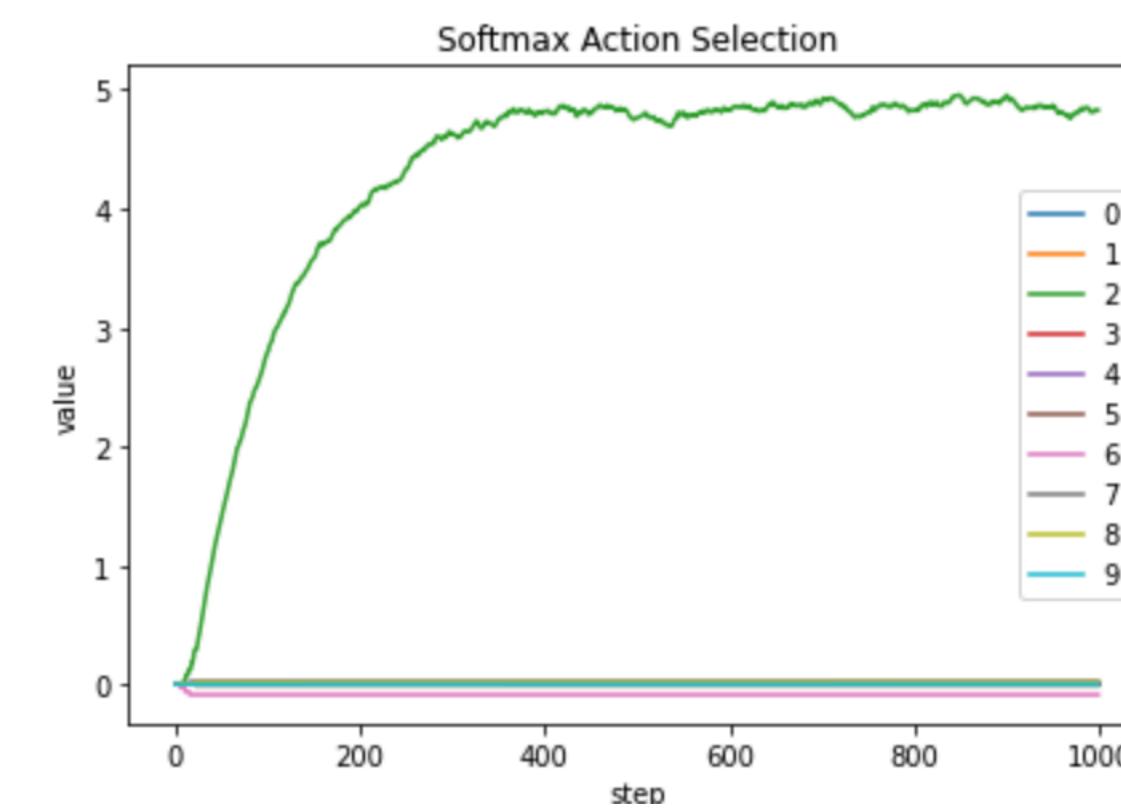
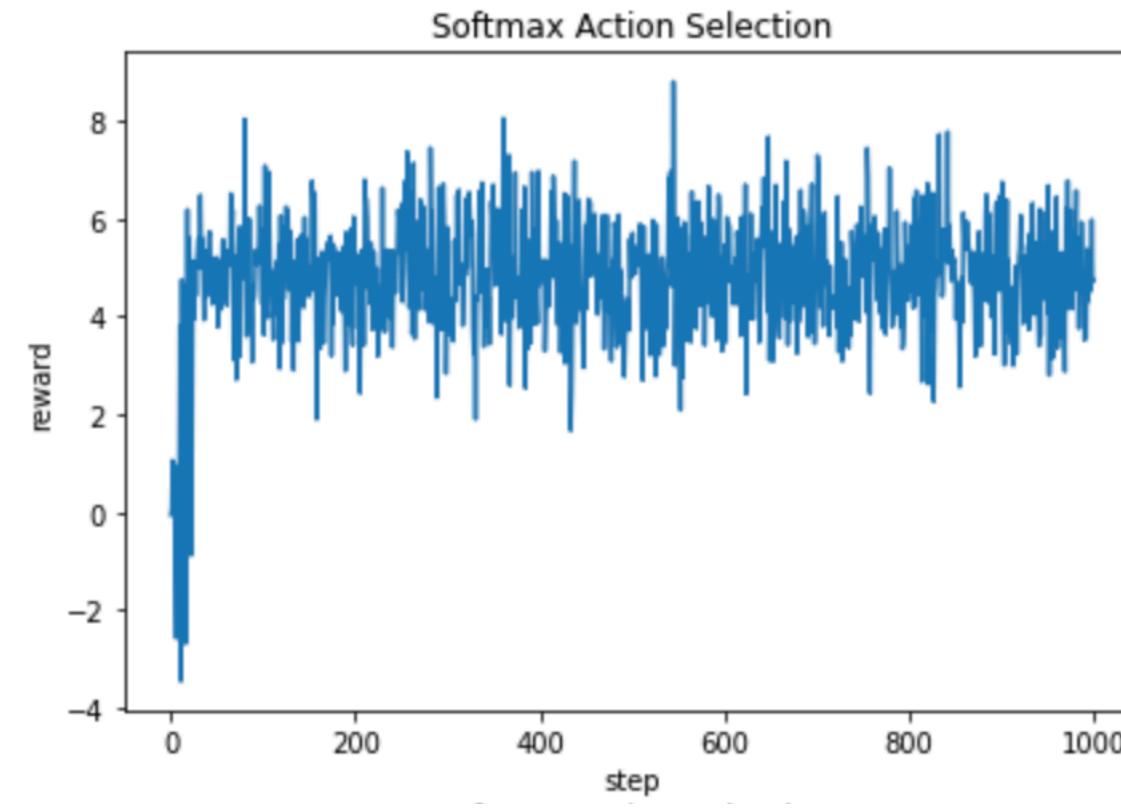
Perfomance comparison

$$P(a_t = a) = \frac{e^{V_t(a) \cdot \beta}}{\sum_{i=1}^N e^{V_t(a_i) \cdot \beta}}$$

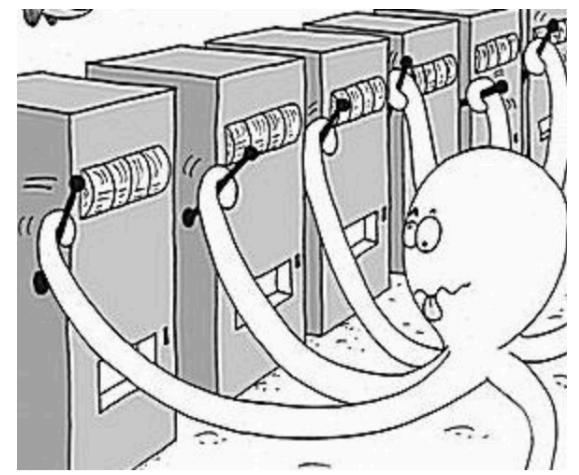
$\beta = 0.2$



$\beta = 10$



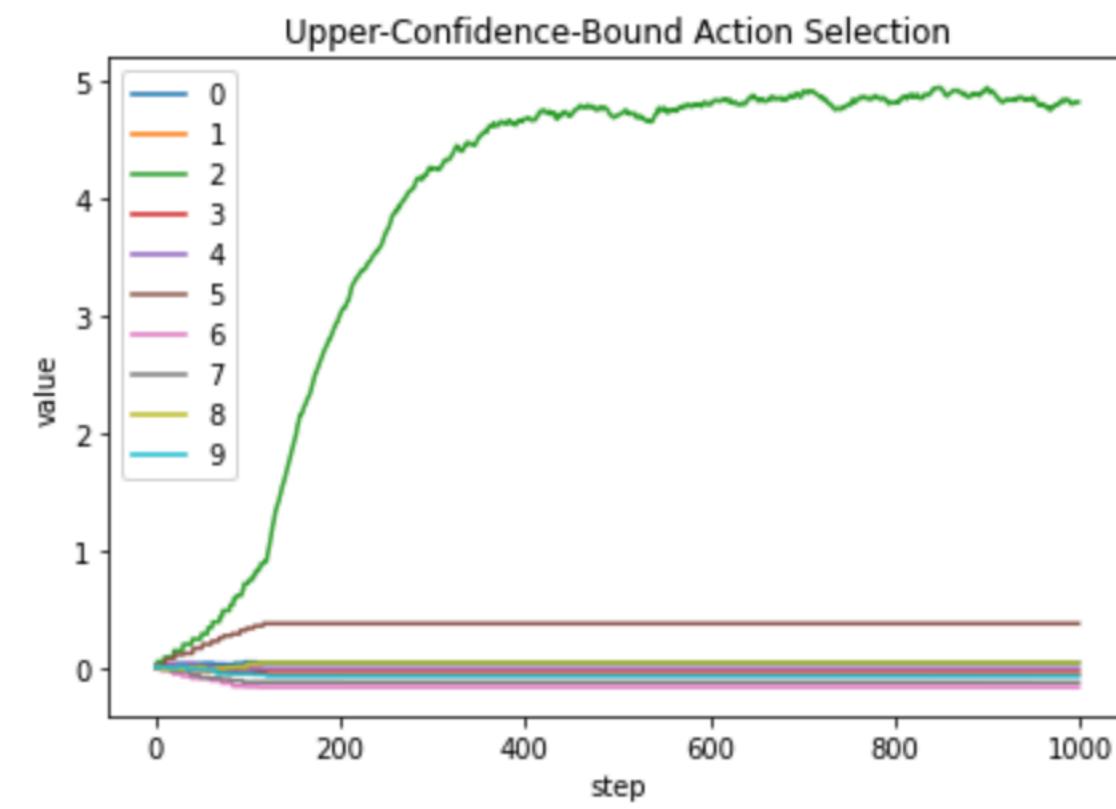
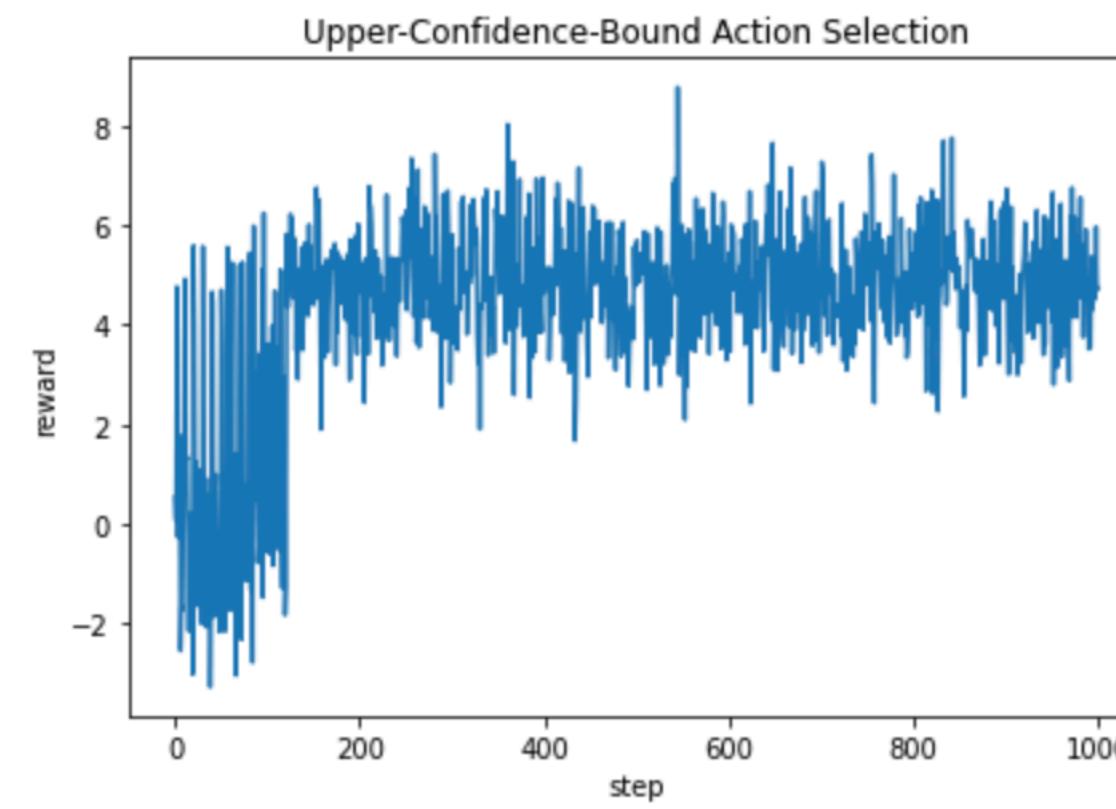
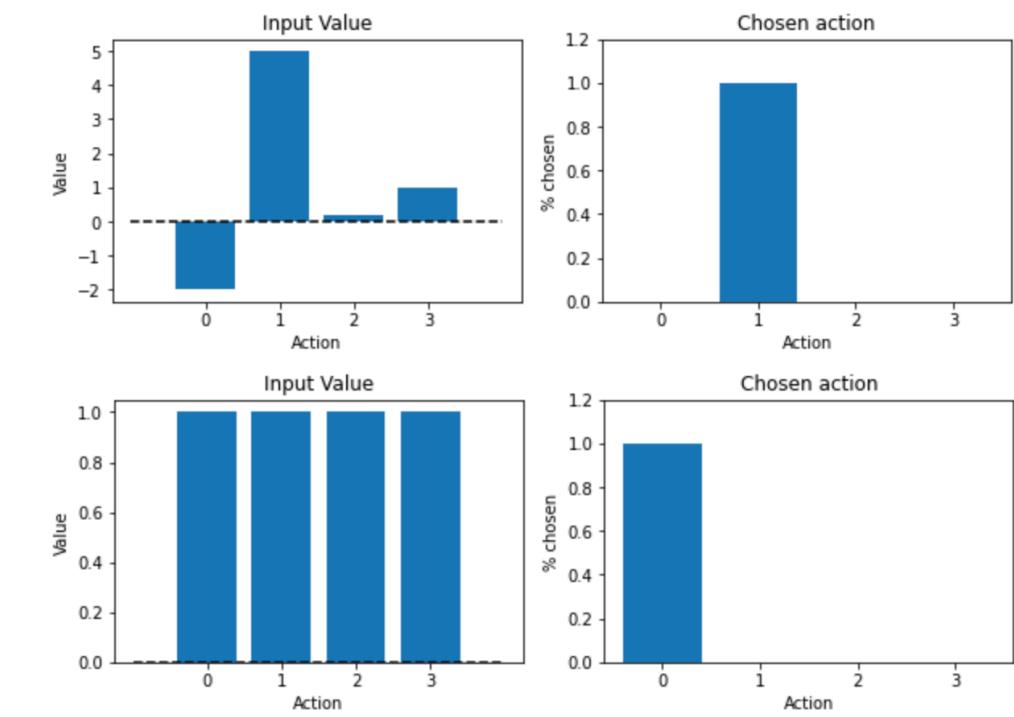
Performance comparison



N arms = 10
N simulations = 1000

**Upper-confidence-bound
(UCB) action selection:**

$$P(a_t = a) = \operatorname{argmax}_a [V_t(a) + c \cdot \sqrt{\frac{\ln t}{N_t(a)}}]$$

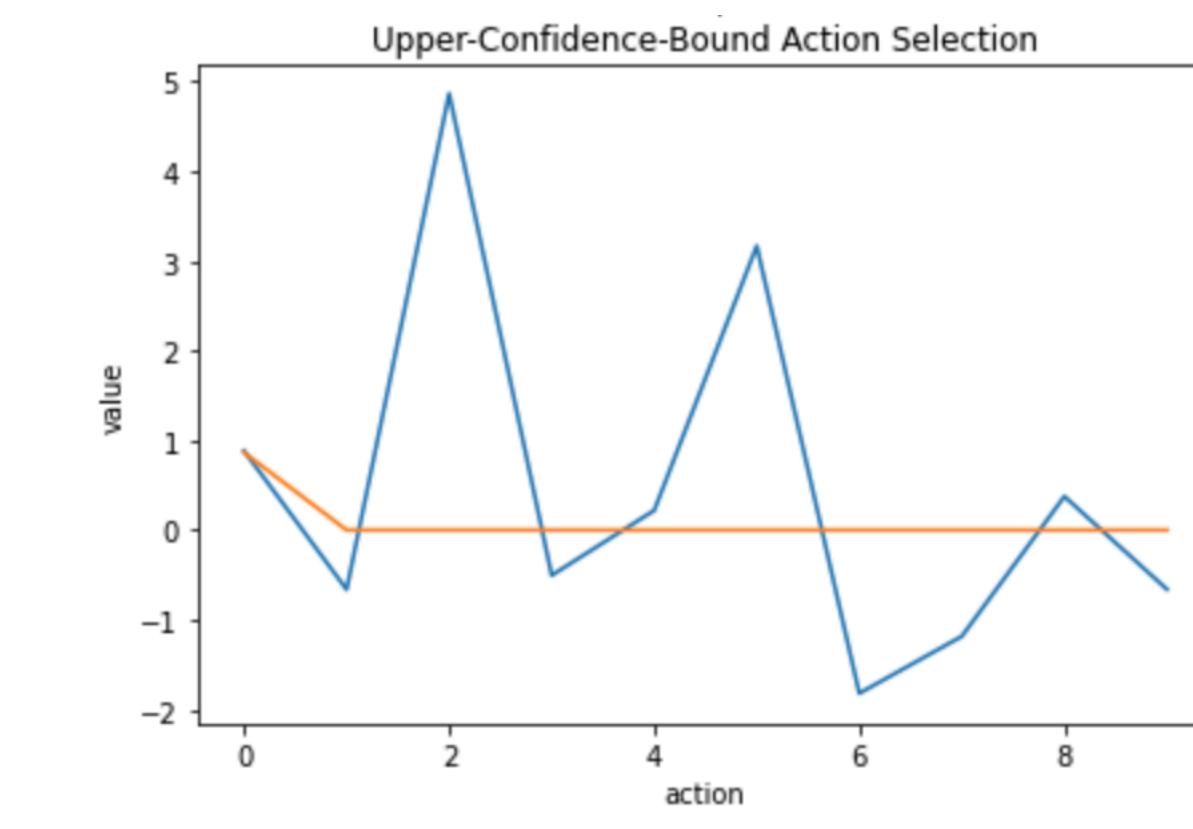
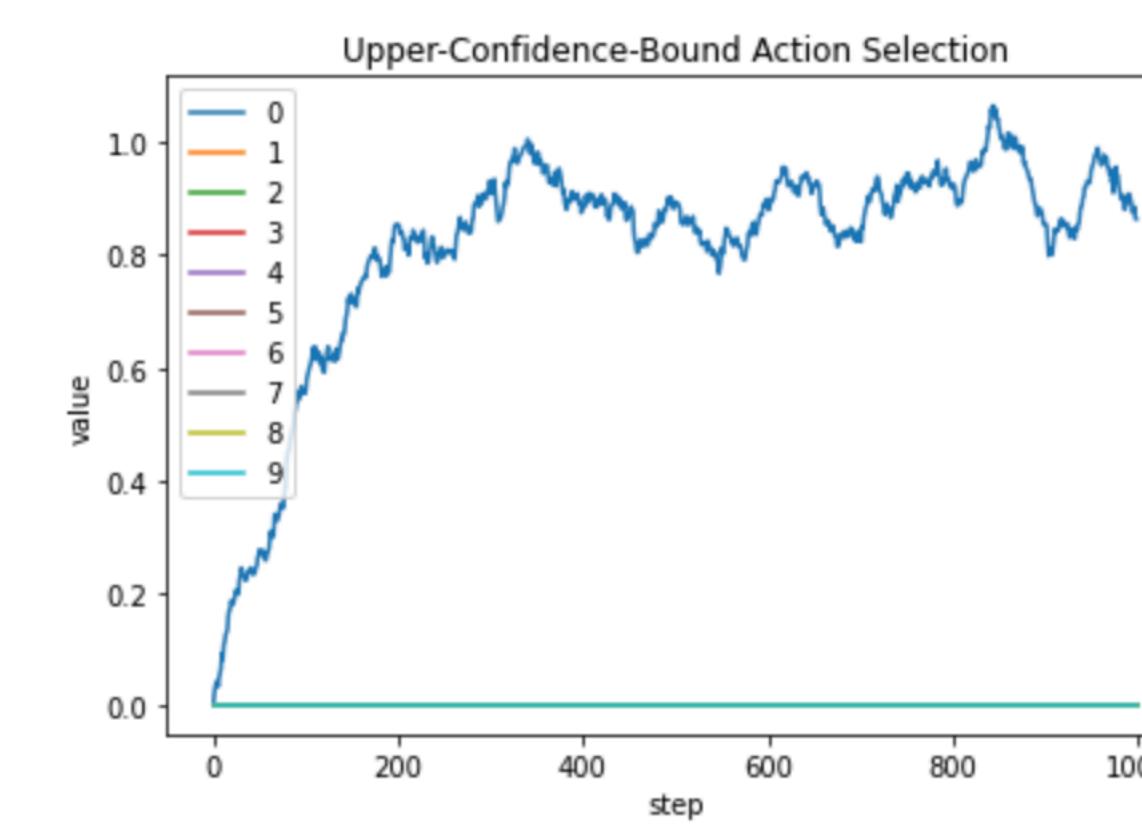
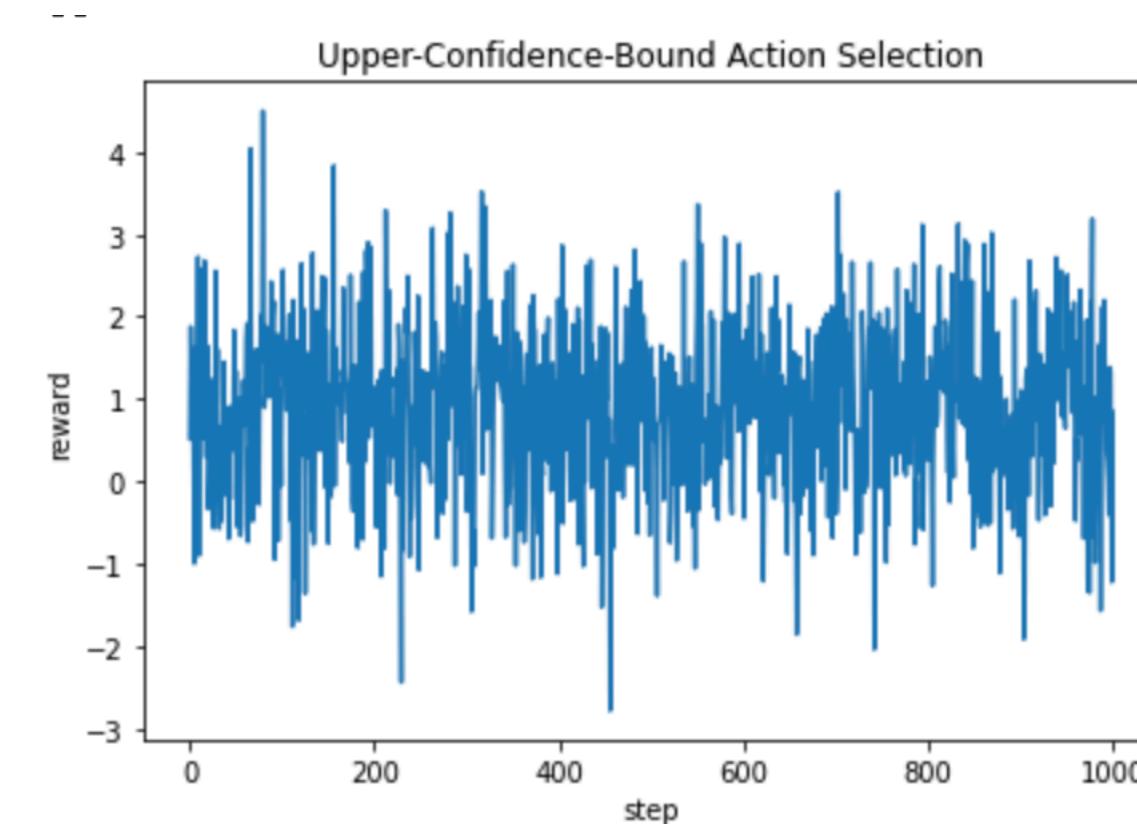


How will your choice of c affect these plots? (Here: $c = 5$)

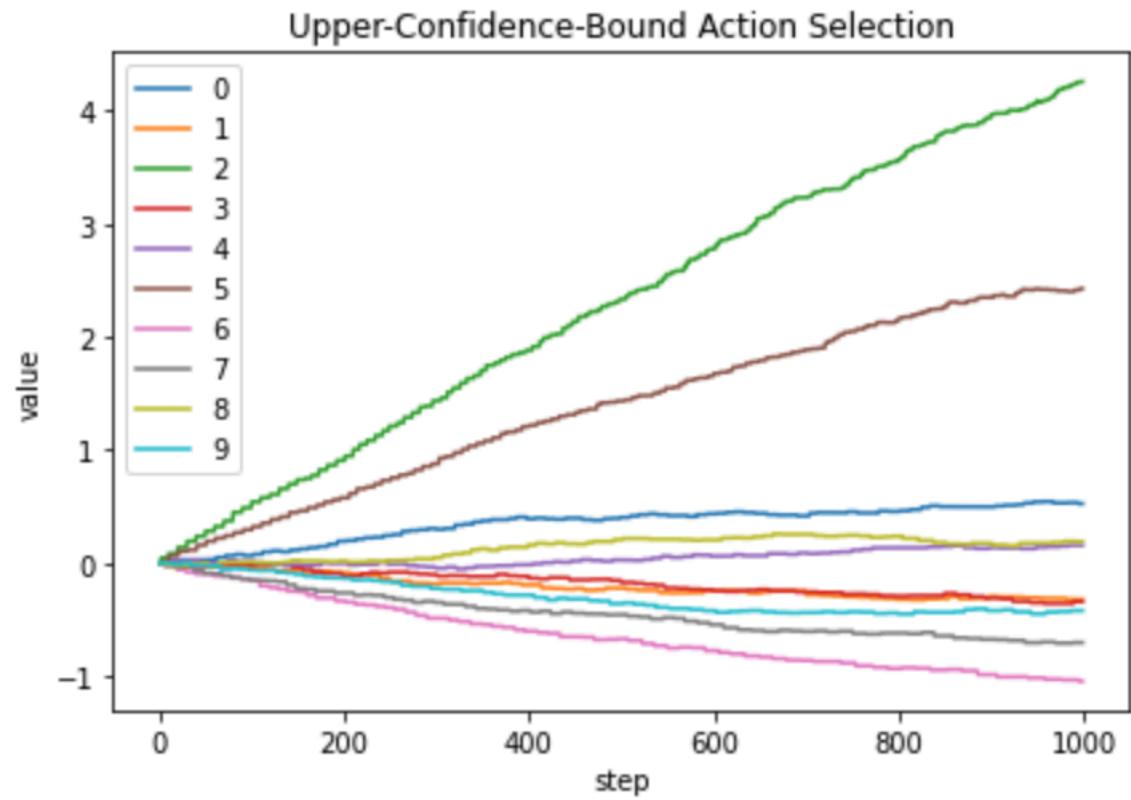
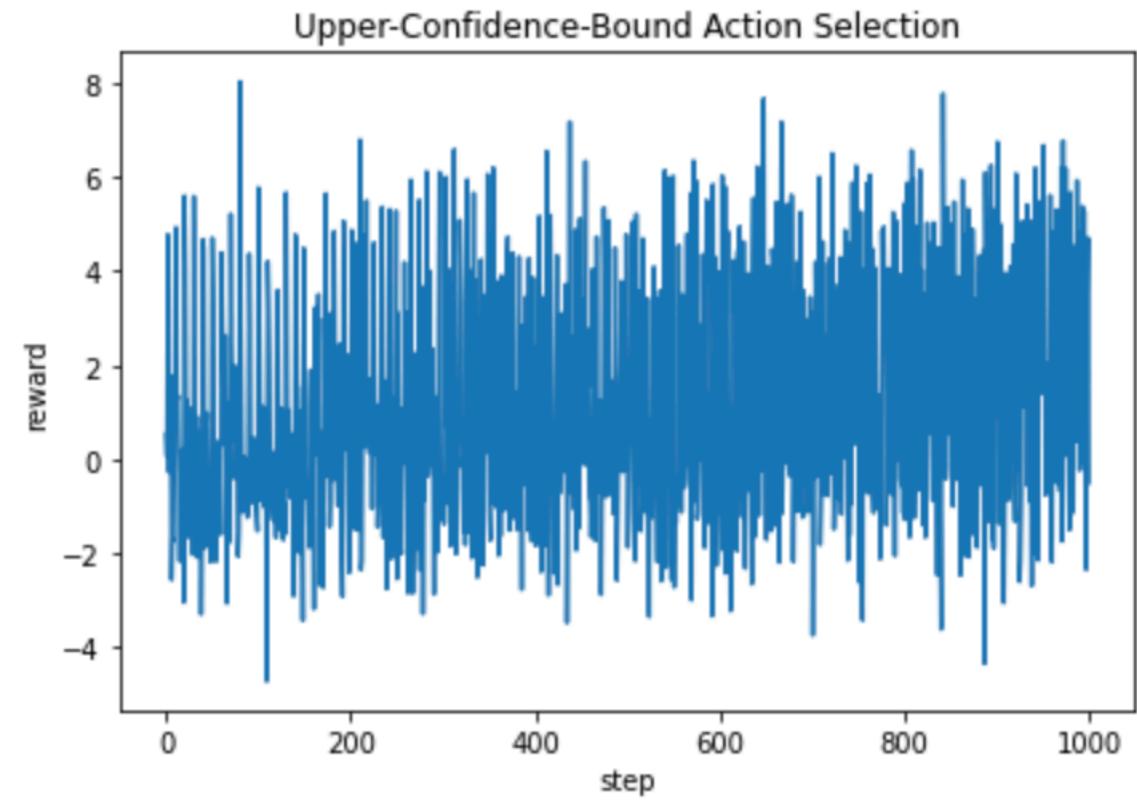
Performance comparison

$$P(a_t = a) = \operatorname{argmax}_a [V_t(a) + c \cdot \sqrt{\frac{\ln t}{N_t(a)}}]$$

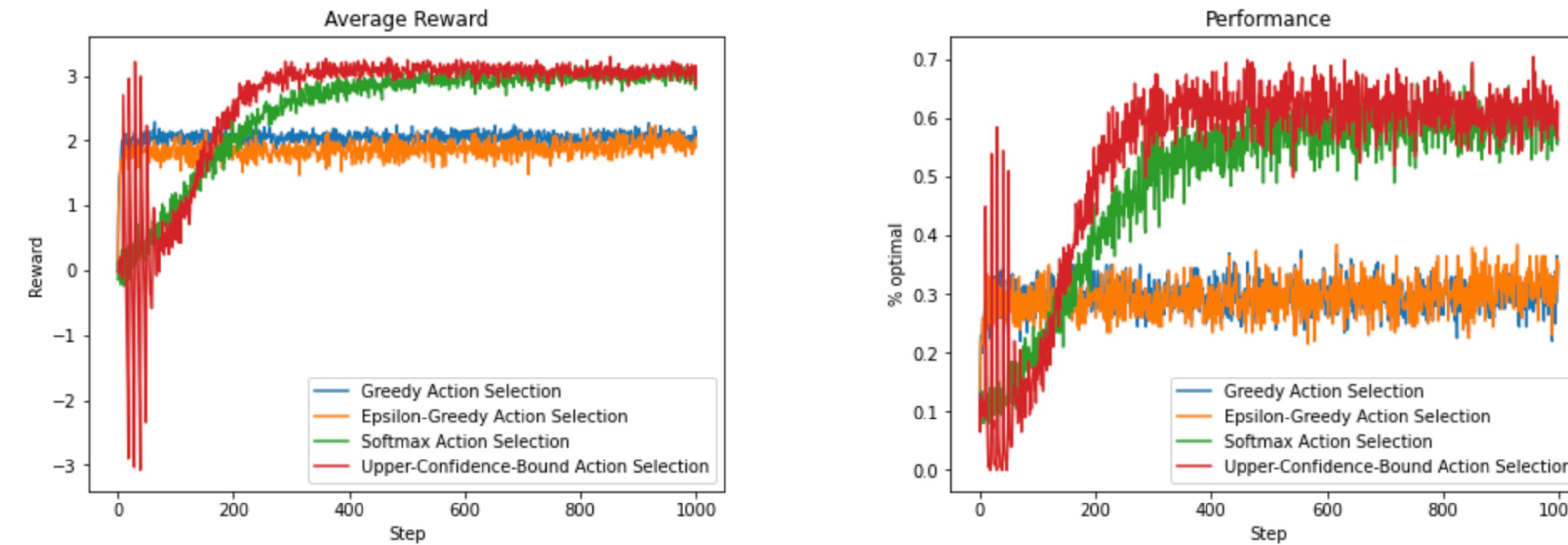
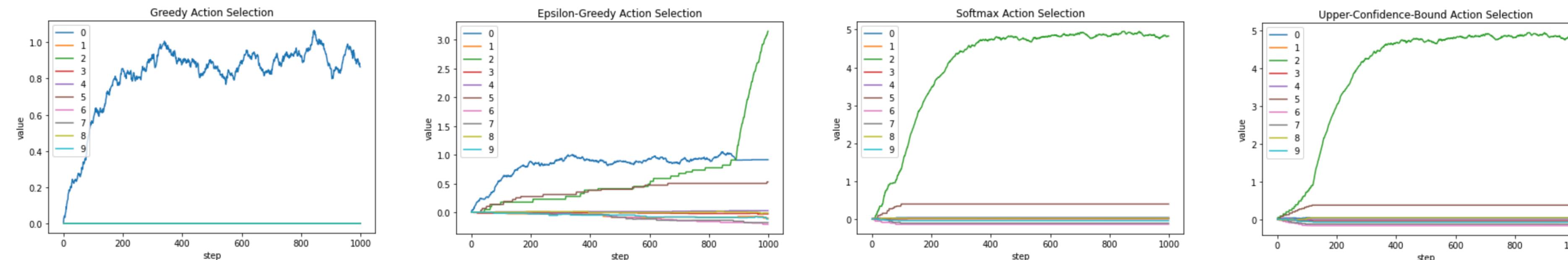
$c = 0$



$c = 50$



Performance comparison



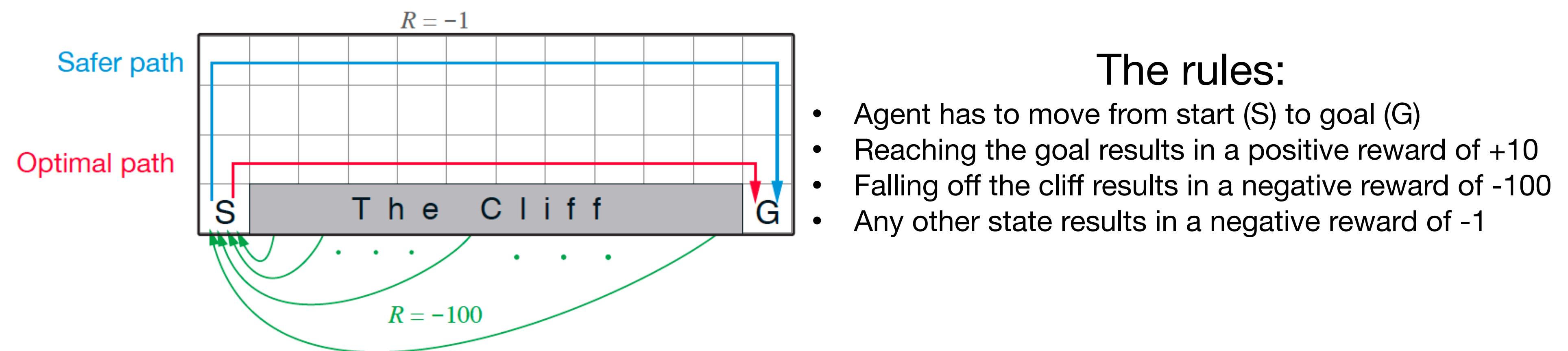
Q-Learning

Limitation of multi-armed bandit problems

Your current action does not influence what happens next!!

How can we solve sequential problems?

The textbook problem:
'Cliff-World'



What's the problem the agent has to solve here??

Note the subtle introduction of the concept of '**transition probabilities**' here
- implicit, later: explicit

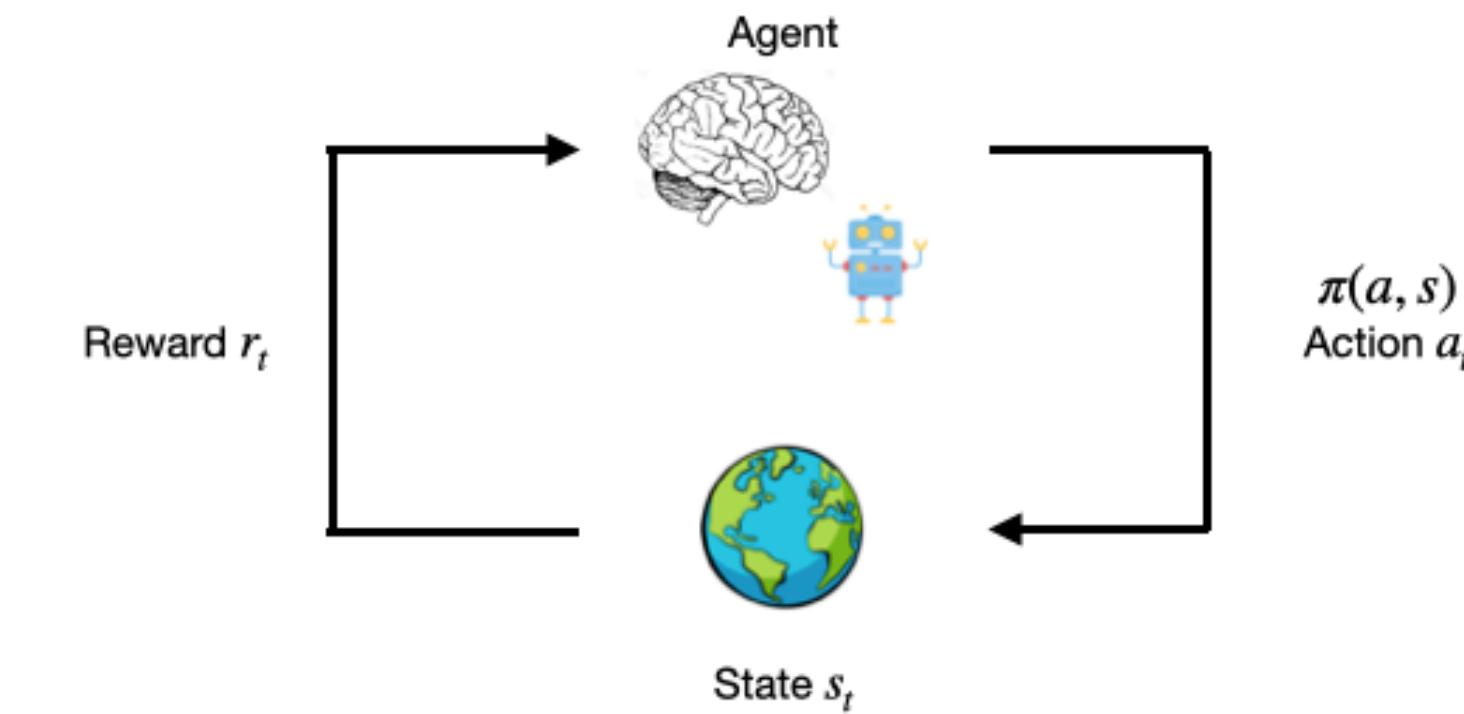
From classical to instrumental learning

TD Learning:

$$V(s_t) \leftarrow V(s_t) + \alpha \cdot (r + \gamma \cdot V(s_{t+1}) - V(s_t))$$

↑ ↑
Learning rate Discount rate

Prediction error



Q-Learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \cdot (r + \gamma \cdot \max_a Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$

↑ ↑
Learning rate Discount rate

Prediction error

What's the difference between $V(s_t)$ and $Q(s_t, a_t)$?

What's $\max_a Q(s_t, a_t)$ doing?

Note that this is just an update rule - doesn't tell us how to select an action!

Coding: Q-Learning

https://github.com/schwartenbeckph/RL-Course/tree/main/2022_06_28