

COMP3702 Artificial Intelligence  
Semester 2, 2023  
Tutorial 8 - Sample Solutions

**Exercise 8.1**

a) In one instance of the MAB, the actions taken and rewards received for the first six trials are given in the table below:

Trial	Action	Reward
1	$A_1$	2.66
2	$A_2$	1.25
3	$A_1$	3.21
4	$A_2$	2.34
5	$A_1$	1.87
6	$A_1$	1.69

Using  $\epsilon$ -greedy, which action is most likely to be chosen next?

- Assuming that the agent wishes to maximise its reward, under the  $\epsilon$ -greedy strategy, it will choose the action with the highest mean reward with probability  $1 - \epsilon$ , and a random action with probability  $\epsilon$ .
- Further, assume  $\epsilon < 0.5$ .
- Estimated rewards, by calculating the mean observed rewards:
  - $\hat{v}_1 = (2.66 + 3.21 + 1.87 + 1.69)/4 = 2.3575$
  - $\hat{v}_2 = (1.25 + 2.34)/2 = 1.795$
- So  $A_1$  is most likely to be chosen next, with probability of at least  $1 - \epsilon$ .

b) Given the same sample information, now consider UCB1 with upper bounds given by:

$$UCB1_a = \hat{v}_a + \sqrt{\frac{C \ln(N)}{n_a}}$$

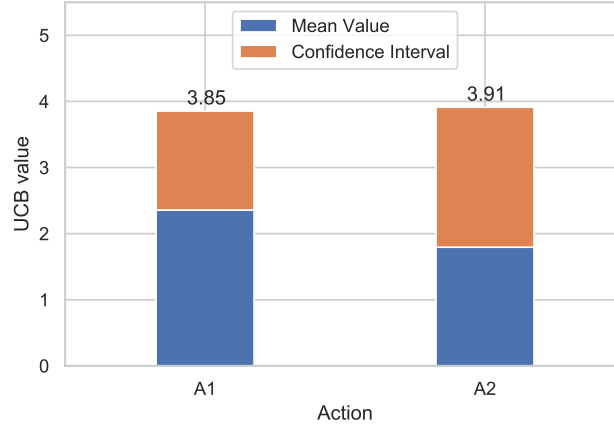
Set the tunable parameter  $C$  to 5. Using this UCB algorithm, which action is chosen next?

From the table, we have  $N = 6$  and  $n_{a_1} = 4$ ,  $n_{a_2} = 2$ .

Now, using the formula above,

- the interval for  $A_1$  is 1.4965 and the upper confidence bound for  $A_1$  is 3.8541
- the interval for  $A_2$  is 2.116 and the upper confidence bound for  $A_2$  is 3.9115
- So  $A_2$  is chosen next, as it has the greater upper confidence bound.

$A_2$  has been sampled fewer times than  $A_1$ . This implies that there is typically greater uncertainty in its mean estimate, which is reflected in the larger interval size; and this pushes its upper confidence bound slightly above that of  $A_1$ .



To see the link to UCT and MCTS, consider this restatement of the UCB1 upper bound expression:

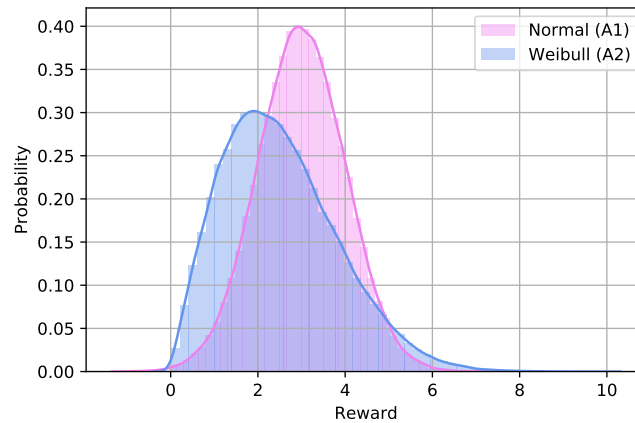
$$UCB1_a = \hat{Q}(a) + \sqrt{\frac{C \ln(N)}{n_a}}$$

Because a MAB has only one state, we can drop the state argument. However, if we explicitly include it, the connection to MCTS becomes apparent:

$$\begin{aligned} UCT_{a,s} &= \hat{Q}(a, s) + \sqrt{\frac{C \ln(N_s)}{n_{a,s}}} \\ &= \frac{R(a, s)}{n_{a,s}} + \sqrt{\frac{C \ln(N_s)}{n_{a,s}}} \end{aligned}$$

where now  $\hat{Q}$  is a table of  $(a, s)$  values indexed by the state node,  $s$ , from which the sample and roll-out simulation is taken, and the action chosen,  $a$ .

c) Plot the distributions of rewards from each arm. If the agent wishes to maximise its cumulative reward over time and knew these distributions, which would be the optimal arm to pull?



Explanation:

- The mean of the normal distribution is given as  $\mu = 3$ .
- The mean of the Weibull distribution is a complicated function of  $a$  and  $b$ , but in this case it is equal to  $2\sqrt{\pi}/2 \approx 2.506$ .

d) Set up a MAB instance with two arms described above, and consider the  $\epsilon$ -greedy exploration strategy with random sampling parameter set to  $\epsilon = 0.1$ , and the UCB bound as described in b) above. For each strategy, plot their cumulative rewards over 1000 arm trials in an MAB instance. **Questions:** Which performs better initially? Which performs better in the long run?

It is difficult to say which does better initially, as performance depends on the random realisation of rewards.

However, in the long run, we expect the UCB1 algorithm to outperform  $\epsilon$ -greedy. This is because  $\epsilon$ -greedy continues to sample random arms with probability  $\epsilon$  even long after the mean reward estimates for all have converged to very close to their true values. In contrast, UCB1 reduces the chance of exploring away from the highest-value arm by adjusting the confidence interval based on the number of samples taken.

## Exercise 8.2

a)  $V^*(s) = -|s - 4|$

b) From lecture notes on Q-Learning, we know that given a tuple  $(s, a, r, s')$  we use the following update equation:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left( r + \gamma \max_{a' \in \{-1, 1\}} Q(s', a') - Q(s, a) \right)$$

Using this equation with  $\alpha = 0.5$ ,  $\gamma = 1$ , we have:

$$Q(3, -1) \leftarrow 0 + 0.5 \left( -1 + \max_{a' \in \{-1, 1\}} Q(2, a') \right) = 0.5(-1 + 0) = -0.5$$

$$Q(2, 1) \leftarrow 0 + 0.5 \left( -1 + \max_{a' \in \{-1, 1\}} Q(3, a') \right) = 0.5(-1 + 0) = -0.5$$

$$Q(3, 1) \leftarrow 0 + 0.5 \left( -1 + \max_{a' \in \{-1, 1\}} Q(4, a') \right) = 0.5(-1 + 0) = -0.5$$

## Exercise 8.3

a)

For the Gridworld environment, we have the following dimensions of complexity:

Dimension	Values
Modularity	<b>flat</b> , modular, hierarchical
Planning horizon	non-planning, finite stage, <b>indefinite stage</b> , <b>infinite stage</b>
Representation	<b>states</b> , <b>features</b> , relations
Computational limits	<b>perfect rationality</b> , bounded rationality
Learning	knowledge is given, <b>knowledge is learned</b>
Sensing uncertainty	<b>fully observable</b> , partially observable
Effect uncertainty	deterministic, <b>stochastic</b>
Preference	goals, <b>complex preferences</b>
Number of agents	<b>single agent</b> , multiple agents
Interaction	offline, <b>online</b>

In RL in general, we may have other variants, e.g. partially observable, multi-agent RL, offline, etc.

b) In RL, what is the connection between the environment's state transition function and the exploration policy used by the agent?

In reinforcement learning problems, when the transition function is unknown/uncertain, we should have a higher exploration term to learn information about the environment. As it becomes more certain, we should lower the exploration term (similar to UCB, and “Annealing epsilon-greedy” in the supporting code, where  $\epsilon$  is reduced with time).