# Tutorial Week 7: MDP

**Guidelines**

- You can discuss the content of the questions with your classmates.
- However, everyone should work on and be ready to present ALL the solutions.
- Your attendance is marked in the tutorial and participation noted to award class participation marks.

## Problem 1: Online Search for Markov Decision Process

Consider an MDP where the state is described using $M$ variables where each variable can take $n$ values. The MDP has 2 actions and at each state each action can only lead to 2 possible next states.

a) What is the size of the state space of this MDP? Can this MDP be efficiently solvable with value iteration as $M$ grows?

b) A search tree of depth $D$ (number of actions from the root to any leaf is $D$) is constructed from an initial state $s$. What is the size of the search tree (the number of nodes and edges) as a function of $M$ and $D$, in $O$-notation? Can online search be done efficiently as $M$ grows if $D$ is a fixed small constant?

c) MCTS is used for solving this MDP. What is the size of the search tree if $T$ trials of MTCS is performed up to a search depth of $D$, as a function of $M$, $D$ and $T$ in $O$-notation?

d) Consider a search tree where the reward is zero everywhere except at the leaves. When a MCTS trial goes through a node, we say that an action at the node wins if the trial ends in a leaf with reward $1$. Consider an MCTS simulation where a node has been visited 16 times and has two actions, A and B. Action A has a won 2 out 4 times whereas action B has won 8 out of 12 times. Which action will the MCTS algorithm chose given the exploration parameter $c$ is set to 1? Give the values of $\pi_{UCT}$ for the node (consider log base 2 in UCT bound).

## Problem 2: Value Iteration

Consider the following 2 state, 2 action MDP with discount factor 0.9.

| $P(s_1|s_1, a_1)$ | $P(s_2|s_1, a_1)$ | $P(s_1|s_2, a_1)$ | $P(s_2|s_2, a_1)$ |
|---|---|---|---|
| 0.9 | 0.1 | 0 | 1 |

| $P(s_1|s_1, a_2)$ | $P(s_2|s_1, a_2)$ | $P(s_1|s_2, a_2)$ | $P(s_2|s_2, a_2)$ |
|---|---|---|---|
| 0.1 | 0.9 | 0 | 1 |

| $R(s_1, a_1)$ | $R(s_1, a_2)$ | $R(s_2, a_1)$ | $R(s_2, a_2)$ |
|---|---|---|---|
| 1 | 0 | 3 | 3 |

1. Assume a finite horizon problem with horizon 1 (only 1 action is to be taken). What is the utility or value function and the optimal action in each state?

2. Assume a finite horizon problem with horizon 2 (2 actions to be taken). What is the utility or value function and the optimal action in each state?

3. What is the optimal infinite horizon policy?

## Problem 3: Bellman operator

**[RN 17.6]** Suppose that we view the Bellman update

$$U_{t+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_t(s')$$

as an operator $B$ that is applied simultaneously to update the utility of every state, that is,

$$U_{t+1} \leftarrow B U_t .$$

We claim that the Bellman operator $B$ is a contraction.

1. Show that, for any function $f$ and $g$,

$$|\max_a f(a) - \max_a g(a)| \le \max_a |f(a) - g(a)| .$$

2. Write out an expression for $|(BU_t - BU'_t)(s)|$ and then apply the result from part 1 to complete the proof that the Bellman operator $B$ is a contraction.