

National University of Singapore
School of Computing

Semester 1, AY2021-22

CS4246/CS5446

AI Planning and Decision Making

Issued: 10 Sep, 2021

Tutorial Week 6: Markov Decision Process

Guidelines

You may discuss the content of the questions with your classmates. But everyone should work on and be ready to present ALL the solutions.

Problem 1: Formulating Markov Decision Processes

Specify the following problem as a Markov decision process, *i.e.* specify the state space, the actions, the transition functions, and the reward function. What is the (approximate) size of the state space and the action space?

- Atari games. Atari games have 128 bytes of RAM, 18 actions, and 33,728 screen pixels taking values from 0-127.

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: 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 MCTS 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 won 2 out of 4 times whereas action B has won 8 out of 12 times. Which action will the MCTS algorithm choose given the exploration parameter c is set to 1? Give the values of π_{UCT} for the node (consider log base 2 in UCT bound).
-