

---

## Homework 2

Please write the following on your homework:

- Name
- Collaborators (write none if no collaborators)
- Source, if you obtained the solution through research, e.g. through the web.

While you may collaborate, you **must write up the solution yourself**. While it is okay for the solution ideas to come from discussion, it is considered as plagiarism if the solution write-up is highly similar to your collaborator's write-up or to other sources. For the programming assignment, we will be running an **automated plagiarism detector**. Highly similar programs will be investigated.

Your solution should be submitted by **Sunday 11 October 11.59pm**. Scanned handwritten solutions are acceptable but must be legible.

*Late Policy:* A late penalty of 20% per day will be imposed for the written assignment (no submission accepted after 5 late days) unless prior permission is obtained. Late submissions will not be accepted for the programming assignment.

---

### 1. Markov Decision Process

- Assume that you are given a directed graph with  $n$  vertices and positive weights and would like to compute the shortest path from every vertex to a goal vertex  $g$ . Describe how to model the problem as a MDP so that it can be solved with an appropriately initialized value iteration algorithm, i.e. describe the state space, action space, transition function, and reward function.
- Assume that we now have  $M$  agents on the same graph and the reward is the sum of the rewards of the agents. Can value iterations be used to solve the problem efficiently (as a function of  $M$ ) in the following two cases? Why?
  - There is no restriction on the number of agents that can be at each vertex.
  - Only one agent can be at each vertex at any time, i.e. having more than one agent at a vertex incurs a very large penalty. Note that every agent must move to another vertex at every time instance.
- For the problem with  $M$  agents, the number of actions grows exponentially with  $M$ . Suggest one way to allow each trial of Monte Carlo Tree Search with UCT to run efficiently.

- Search Tree** Consider a 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 possibly at some 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 of 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).

### 3. Programming Assignment

In the first homework, we learnt to solve deterministic planning problems with pddl solvers. In this task, we will learn to solve planning problems in non-deterministic environments with two separate algorithms, Value Iteration and Monte Carlo Tree Search.

We will be using the same `gym_grid_environment` (<https://github.com/cs4246/gym-grid-driving>) to simulate the solutions. All dependencies have been installed in the docker image “cs4246/base” which you have already downloaded for Homework 1.

By now, you should be familiar with the `gym_grid_environment`. Go through the IPython Notebook file on Google Colab [here](#), if you have already not done so in the previous programming assignment. You are now ready to begin solving the two tasks which are listed below.

- (a) **Value Iteration:** In the crossing the road task in homework 1, we assumed all the cars move with a constant speed of -1. Most of the times, this assumption does not hold. There is inherent noise in the car speeds due to reasons like bumpy roads, uneven driving etc. This non-deterministic behaviour is not captured in the PDDL format.

For problems that are not very large, planning algorithms like Value Iteration can be used to handle non-deterministic behaviours. In this task, we model the problem as a Markov Decision Process and use Value Iteration algorithm to find the optimal policy.

We have provided the script which models the problem as an MDP. Every state in this MDP is represented as a tuple of features

$\{\text{Agent}_x, \text{Agent}_y, \text{Car}_x^1, \text{Car}_y^1, \dots, \text{Car}_x^N, \text{Car}_y^N, \text{isTerminal}\}$ ,  
where

- $\text{Agent}_x, \text{Agent}_y$  denote the  $x$  and  $y$  coordinates agent.
- $\text{Car}_x^N, \text{Car}_y^N$  denote the  $x$  and  $y$  coordinates of  $\text{Car}^N$ .
- `isTerminal` is a boolean variable denoting whether the state is terminal.

Your task is to fill up the function that implements the Value Iteration algorithm which is marked with “FILL ME”.

You can test your code with the test configurations given in the script by running `python __init__.py` (present in the .zip file) on the docker (using the `docker run ...` command).

HINT : You can use Matrix multiplication to optimise the code, but refrain from using the function `np.matmul()/np.multiply()/np.dot()` as these functions interfere with the evaluation script. Instead you can use function `matmul()` defined in the same script.

#### (b) Monte Carlo Tree Search :

Unfortunately, we moved to a bigger parking lot and Value Iteration is not feasible as it does not scale well with large state spaces. However, we can use Monte Carlo Tree Search (MCTS) to handle such cases.

We have provided a separate script which solves the problem with MCTS, Your task is to complete the script by filling up 2 code snippets inside the functions marked with “FILL ME”.

The functions required to be filled up are,

- `backpropagate()` : This function should implement the backpropagation step of MCTS.
- `chooseBestActionNode()` : Populate the list `bestNodes` with all children having maximum value. The value of the  $i^{th}$  child node can be calculated with the formula  $\frac{v_i}{n_i} + c\sqrt{\frac{\ln N}{n_i}}$ . where,
  - $v_i$  : sum of returns from the  $i^{th}$  child node
  - $n_i$  : Number of visits of the  $i^{th}$  child node
  - $N$  : Number of visits of the current node
  - $c$  : The exploration constant. We set it to be 1.

For more details on MCTS, it might be helpful to look [here](#). You are also encouraged to look through the rest of the code to see how the entire MCTS algorithm is implemented.

You can test your code with the test configurations given in the script by running `python __init__.py` (present in the .zip file) on the docker (using the `docker run ...` command).

#### Submission Instructions

Please follow the instructions listed [here](#). You have to make 2 separate submissions as mentioned below :

- i. For Task (a) : You have been provided with a `vi_agent.zip` file in the correct submission format. Complete the functions marked with “FILL ME” in `__init__.py`. Your submission zip file should have the exact same structure as the zip file you received.
- ii. For Task (b) : Modify the `__init__.py` in the `mcts_agent.zip` by completing the functions with “FILL ME” and make a separate submission.

Note :

- i. Please remember to set the variable `SUBMISSION` in `__init__.py` to `False` when testing locally, and to `True` before making a submission.
- ii. Please do not print anything to the console, as it might interfere with the grading script.