

# Assignment II - Planning Decision Process

## Introduction

This assignment involves finding an optimal policy for a given implementation will focus on iterative techniques such as value iteration and their respective variants. Additionally, the analysis will extend to various components defining an MDP on the optimal policy.

## Domain Description

Consider an autonomous car in the grid world environment. The car is allowed to run on petrol or electricity. It is capable of moving up, down, right, and left. Due to construction work, holes have been dug in some cells. Therefore, the autonomous car needs to find an optimal policy (a sequence of actions for every grid cell) to avoid falling into the holes and drive safely to the goal.

## Environment

The environment is represented as a 2D grid world with discrete states. Cells designated as 'H,' indicating the presence of a hole. The garage where the vehicle is initially parked, and the destination is marked as 'G'. All other cells are labelled as 'F,' indicating free cells.

## Transition Model

- The car can move in four directions: top ('t'), down ('d'), left ('l'), and right ('r'). The state change for each action is:-

$$right(x, y) = (x + 1, y)$$

$$left(x, y) = (x - 1, y)$$

$$top(x, y) = (x, y + 1)$$

$$down(x, y) = (x, y - 1)$$

- Because of rain, the roads in the grid world become slippery. The effects of actions are stochastic. If the car moves to an intended cell with a probability  $p$  and randomly in one of the other cells with probability  $(1 - p)$ . Transitions outside the grid boundary result in a position outside the grid boundary.
- It is important to note that moving into a hole or reaching the goal is a terminal *state* in this context.

## Reward Model

- The car gets a reward after taking an action and reaching a new state.
- For each action that does not lead to a hole or the goal, the car receives a `living_reward`. In terminal states, where the agent encounters a hole or the goal, it receives the respective rewards, namely `hole_reward` and `goal_reward`. The reward model  $R(s, a, s')$  is defined as

$$R(s, a, s') = \begin{cases} \text{living\_reward}, & \text{if } s' \text{ is not a terminal state} \\ \text{hole\_reward}, & \text{if } s' \text{ is a hole} \\ \text{goal\_reward}, & \text{if } s' \text{ is the goal} \end{cases}$$

- The discount factor,  $\gamma$ , represents the autonomous car's preference for immediate rewards over accumulating delayed rewards.

## Your Implementation

### Part A: Solving for optimal policy

Your task is to use iterative methods to solve the given MDP. Your implementation should be general enough to be able to work with two maps: `small_map` and `large_map` to test your implementation. The default values for the parameters of the transition and reward model are:

- The transition probability  $p$  is 0.8, meaning that 80% of the time the car moves to the intended cell, while there's a 20% chance of reaching a random cell.
- The `living_reward` and `hole_reward` are 0. The `goal_reward` is 10.

- The discount factor,  $\gamma$ , is 0.9.

## Value Iteration

Implement a vanilla value iteration. As discussed in the lecture, vary the distance between consequent value updates to determine the effect of different epsilon values. Experiment with varying epsilon values until the final value estimate yields a sensible optimal policy.

## Asynchronous Updates

Implement value iteration with asynchronous updates, which are performed in-place within each iteration. Note that the order in which states are updated influences the convergence rate in the case of asynchronous updates following sweeping strategies.

- **Row-major sweep:** Cells are updated row by row, from top-left to the next.
- **Prioritized sweep:** Use the magnitude of Bellman residual as a priority. That is, the states whose successors change the most are updated first, defined as

$$\left| \max_a \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t] - v(S_t) \right|$$

Compare and report the convergence speed, measured in terms of the number of updates, for synchronous and asynchronous versions of value iteration. The difference becomes notable for larger grids.

## Policy Iteration

Implement policy iteration, assuming a random initial policy. Use the value iteration algorithm to determine the convergence of the policy iteration.

For all the algorithms, report the following information for

1. Report the wall time taken for convergence of both the iterations (count of policy iteration cycles for policy iteration, value iteration) till convergence in the default settings. Are they scalable for larger environments?
2. Plot the estimated value of the starting state for every map, one for each map. Each plot should contain four plots - vanilla value iteration, row-major sweep, prioritized sweep, and the evolution of the value of the starting state.
3. Generate a heat map of the values obtained after the iteration for each policy as arrows (denoting the direction of movement), and report if the results in a better optimal policy?

## Part B: Analysis

In the subsequent parts, you will be tasked to change the value iteration implementation and then observe, analyze and report the results and value estimates. Unless otherwise specified, assume that you use the vanilla value iteration implementation to derive the optimal policy. Also, report the results for the

`small_map.`

### B1. Living reward analysis

For this part, assume that the car uses your vanilla value iteration implementation to derive the optimal policy. Also, assume that the transition function is deterministic.

- The petrol prices are fluctuating due to demand in the current market (the expense will increase if it uses petrol to power its engine, which is the living reward). Indicate how your policy will change as the living reward fluctuates between  $-0.9$  and  $+0.9$ . In both cases observe and explain the differences between the results and that of the default settings.
- To decrease the inflation in petrol prices, the grid world game is modified to use electric vehicles. So, using the electric engine will provide a reward of  $+0.001$ . Assume a discount factor ( $\gamma$ ) of  $0.999$ . Run the value iteration on the new setting and report your observations along with the results.

- Plot the heat maps of the values after convergence and the policies obtained under the above settings. Analyse the policies obtained and report the results.

## B2. Changing transition probabilities

- The rain, coupled with snowfall, makes the road exceptionally slippery. Consequently, every action has an equal probability of moving in the intended direction or in a direction perpendicular to the intended direction. If the car is moving right, it may equally likely end up moving right, up, or down, with a probability of 1/3.
- Plot the heat map and the policy and highlight the differences from that of default settings.

# Implementation Guidelines

## Maps and Visualizer

- Please download the package from Moodle. The package contains the following files:

```
A2 |— maps # The layouts of different environments
    |   |— large_map.csv |— plot.py
```

- **Maps:** There are two maps. `small_map` is a grid of size 4 × 4 and `large_map` is a grid of size 50 × 50. Each line in the CSV file corresponds to the layout of the map. The entries are separated by commas, and each entry denotes the type of cell: 'S' (start), 'G' (goal), 'H' (hole), or 'F' (free).
- **Visualizer:** We have provided you with a `plot` function (located in `plot.py`) to visualize the simulation of your policy. However, this is only for visualization purposes and is optional. You are also free to make your own visualization pipeline.

```
def plot(frames, frame_delay = 500, cell_size = 50):
    """
    Visualize the simulation of a policy.
    frames: a list of numpy arrays where each array is a state of the grid world.
    frame_delay: delay between consecutive frame in milliseconds while rendering.
    cell_size: size of each cell in grid world.
    """
    # color coding:
    # cyan: start, green: destination/goal, yellow: hole
```

## Evaluation

You will be evaluated based on the correctness of the implementation, the understanding of the concepts and the validity of your analysis/findings. A viva will

## Submission Guidelines

- This assignment is to be done **individually**.
- Please submit it on Moodle as a single zip file named **<A2>**. The zip file should contain a single python file - **A2.py** - containing your code and a **report.pdf** - that details your algorithm, results and the implementation. Upload plots to the drive (if required), and include the drive link at the end of the report.

|— A2.py |— report.pdf

- **The submission date is 6:00 pm on Wednesday, February 14, 2024.** Late submissions will incur a 10% penalty up to a single day.
- The assignment must be done only from your own original work. Do not reuse previous implementations done by others. Do not violate the university policy on honour code violations discussed in class).