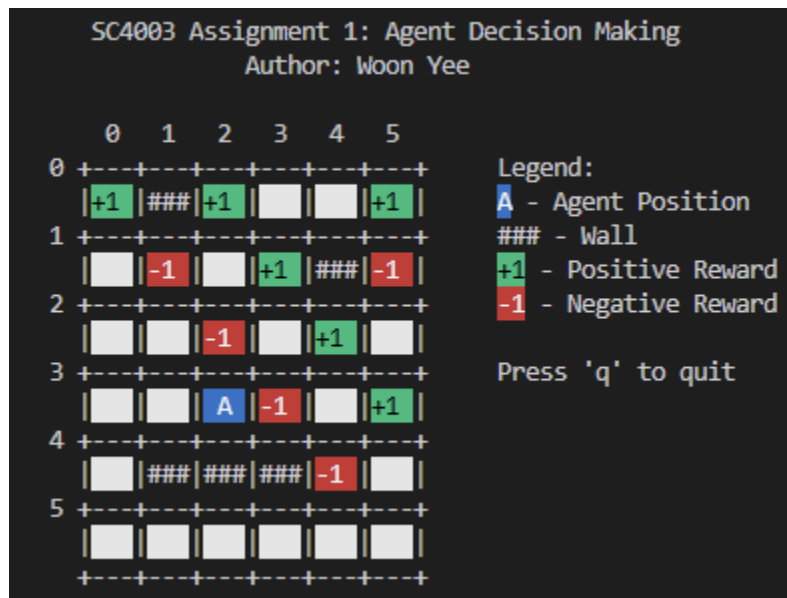Semester 2 AY 24/25

**SC4003 Intelligent Agent Assignment 1**

Completed by:
Ng Woon Yee (U2120622C)

# Initial Maze Environment:



# Part 1:

## Value iteration:

- Descriptions of implemented solutions (5 marks)

The utility of a state is given by

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a)U(s')$$

Which means the utility of a state is the immediate reward for that state plus the expected discounted utility of the next state, assuming that the agent chooses the optimal action.
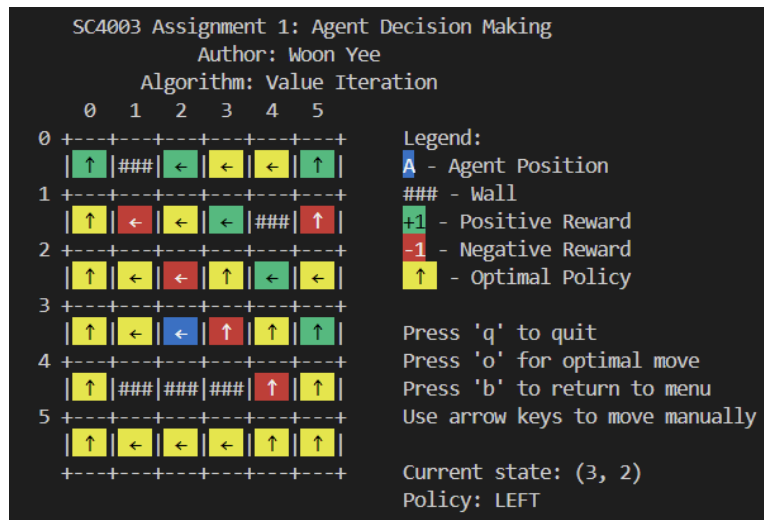
This is called the Bellman equation, and it is the basis of value iteration algorithm for solving Markov Decision Processes. If there are n possible states, that means there are n Bellman equations for the iteration. However, the problem is the Bellman equations are non-linear due to the max operator, and hence these systems cannot be solved by pure linear algebra techniques. And thus, iterative approaches. Initially, we start the utilities with 0. We let $U_i(s)$ be the utility value for state s at the i-th iteration. The iteration step is as follows:

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a)\, U_i(s')$$
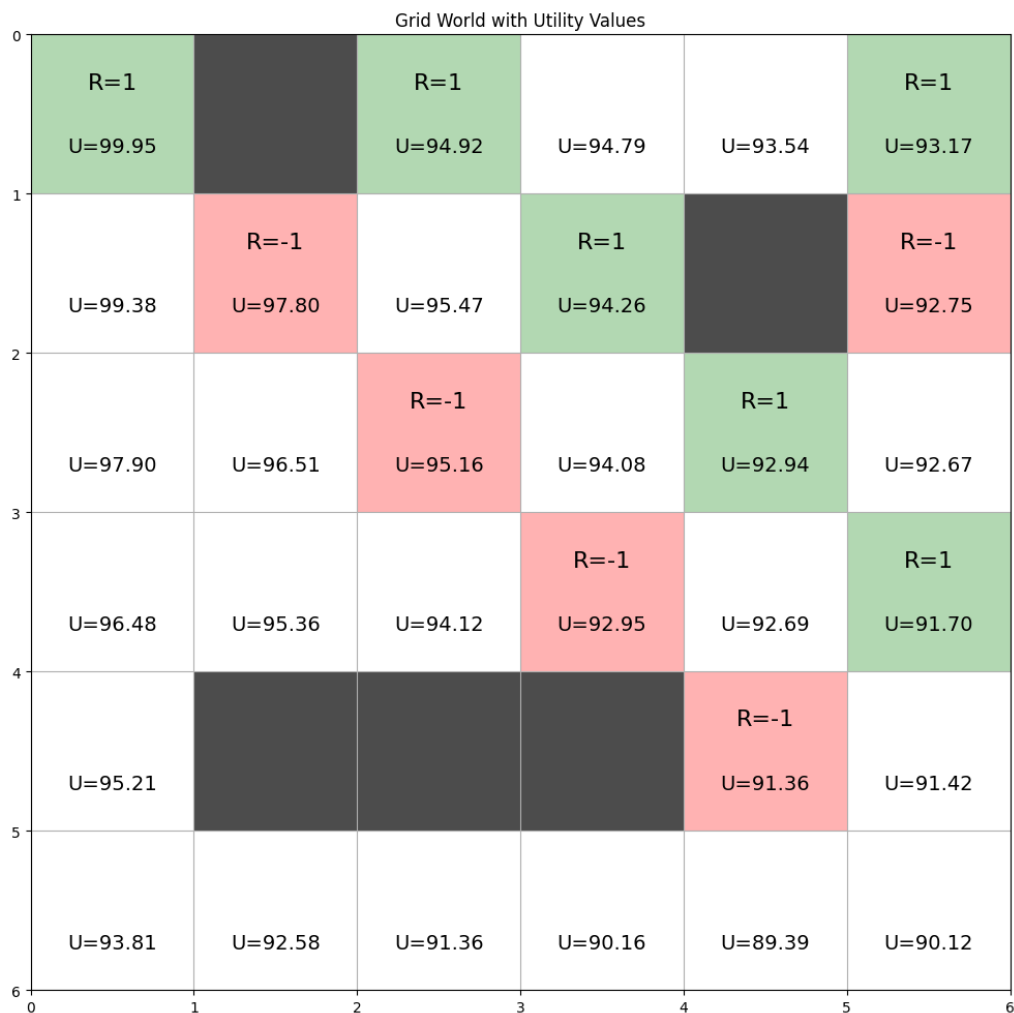
The convergence criterion is based on Bellman error bound:
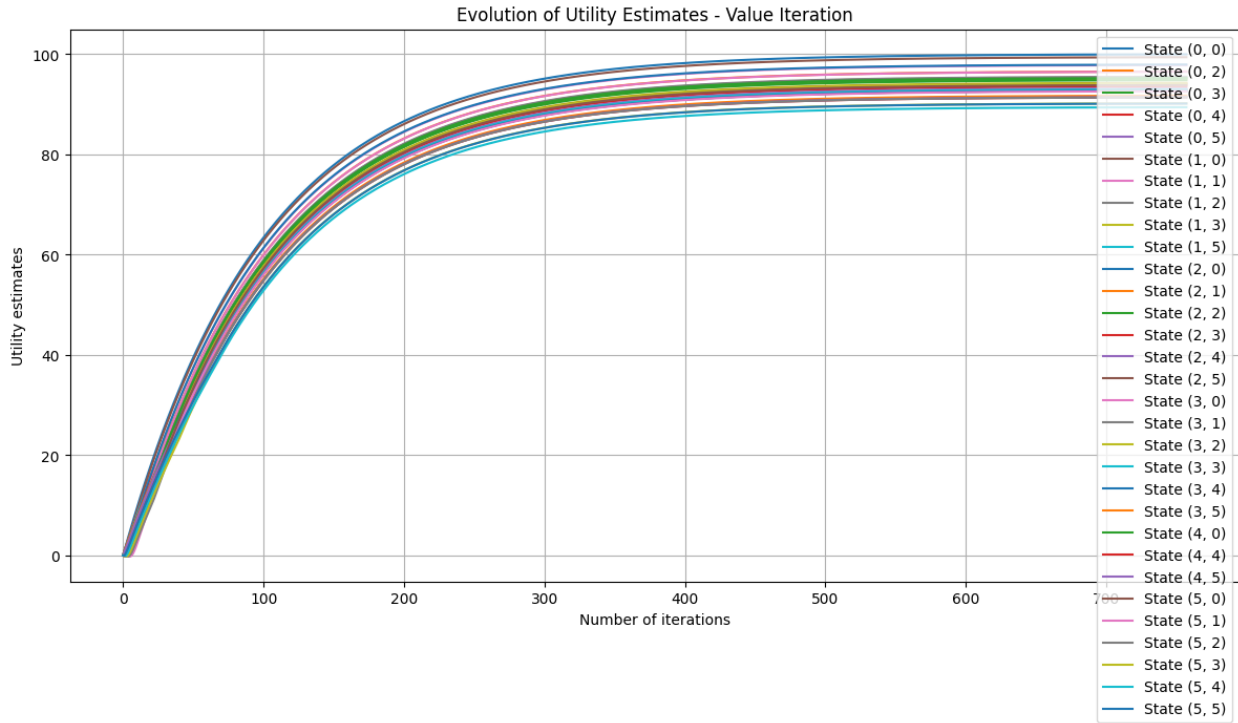
$$\delta < \epsilon(1-\gamma)/\gamma$$

- Plot of optimal policy (10 marks)



- Utilities of all states (10 marks)

- Plot of utility estimates as a function of the number of iterations (10 marks)



Evolution of Utility Estimates - Value Iteration

## Policy iteration:

- Descriptions of implemented solutions (5 marks)

In value iteration, we observed that it is possible to get an optimal policy even when the utility function estimate is accurate. If one action is clearly better than all other choices, then exact magnitude of the utility function need not to be precise. Hence, policy iteration algorithm is proposed, altering between two steps:

- Policy Evaluation: Given a policy $\pi_i$, calculate the utility of each state if $\pi_i$ were to be executed. This is simpler than the standard Bellman Equation as it removes the max operator and restore the linearity of the equations.
- Policy Improvement: Calculate a new policy $\pi_{i+1}$, using one step look ahead based on $U_i$

The utility function in policy iteration is given by

$$U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s))\, U_i(s')$$

And the policy improvement step is as follows:

$$\pi^{i+1}(s) = \underset{a \in A(s)}{\mathrm{argmax}} \sum_{s'} P(s'|s, a) U\ (s')$$

- Plot of optimal policy (10 marks)

SC4003 Assignment 1: Agent Decision Making
Author: Woon Yee
Algorithm: Policy Iteration

```
     0   1   2   3   4   5
0 +---+---+---+---+---+---+
  | ↑ |###| ← | ← | ← | ↑ |
1 +---+---+---+---+---+---+
  | ↑ | ← | ← | ← |###| ↑ |
2 +---+---+---+---+---+---+
  | ↑ | ← | ← | ↑ | ← | ← |
3 +---+---+---+---+---+---+
  | ↑ | ← | ← | ↑ | ↑ | ↑ |
4 +---+---+---+---+---+---+
  | ↑ |###|###|###| ↑ | ↑ |
5 +---+---+---+---+---+---+
  | ↑ | ← | ← | ← | ↑ | ↑ |
  +---+---+---+---+---+---+
```
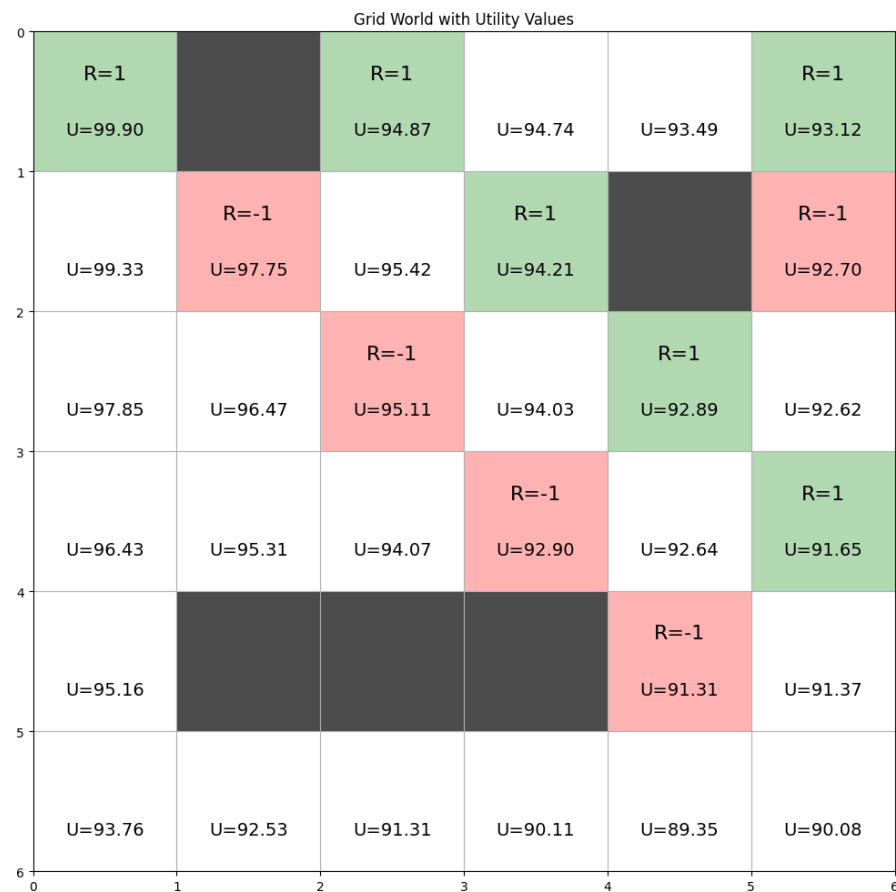
Legend:
A - Agent Position
### - Wall
+1 - Positive Reward
-1 - Negative Reward
↑ - Optimal Policy

Press 'q' to quit
Press 'o' for optimal move
Press 'b' to return to menu
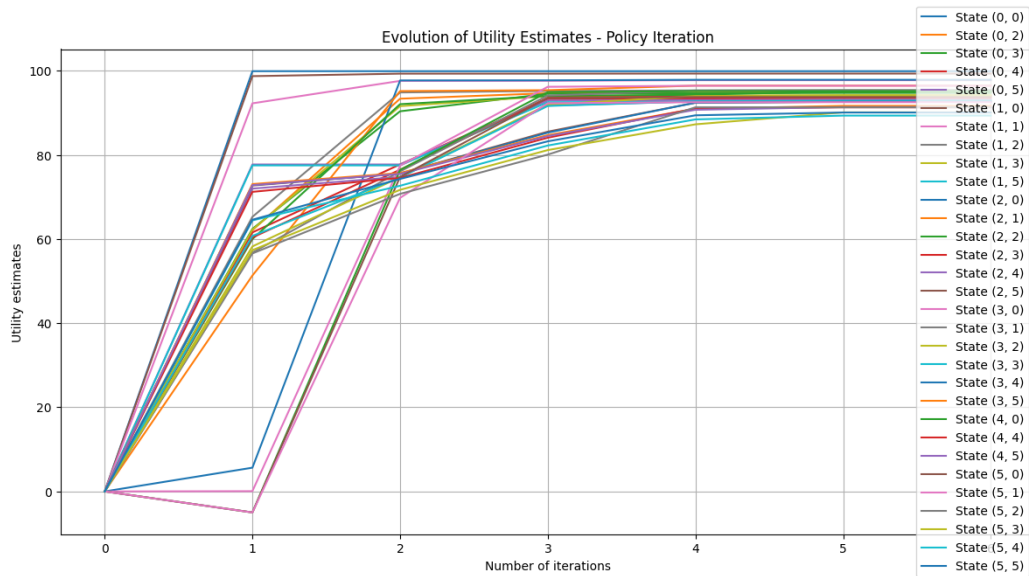Use arrow keys to move manually

Current state: (3, 2)
Policy: LEFT

- Utilities of all states (10 marks)

Grid World with Utility Values

| | | | | | |
|---|---|---|---|---|---|
| R=1 U=99.90 | (wall) | R=1 U=94.87 | U=94.74 | U=93.49 | R=1 U=93.12 |
| U=99.33 | R=-1 U=97.75 | U=95.42 | R=1 U=94.21 | (wall) | R=-1 U=92.70 |
| U=97.85 | U=96.47 | R=-1 U=95.11 | U=94.03 | R=1 U=92.89 | U=92.62 |
| U=96.43 | U=95.31 | U=94.07 | R=-1 U=92.90 | U=92.64 | R=1 U=91.65 |
| U=95.16 | (wall) | (wall) | (wall) | R=-1 U=91.31 | U=91.37 |
| U=93.76 | U=92.53 | U=91.31 | U=90.11 | U=89.35 | U=90.08 |

- Plot of utility estimates as a function of the number of iterations (10 marks)



## Source Code (Include Part 1 and 2) – 15 Marks:

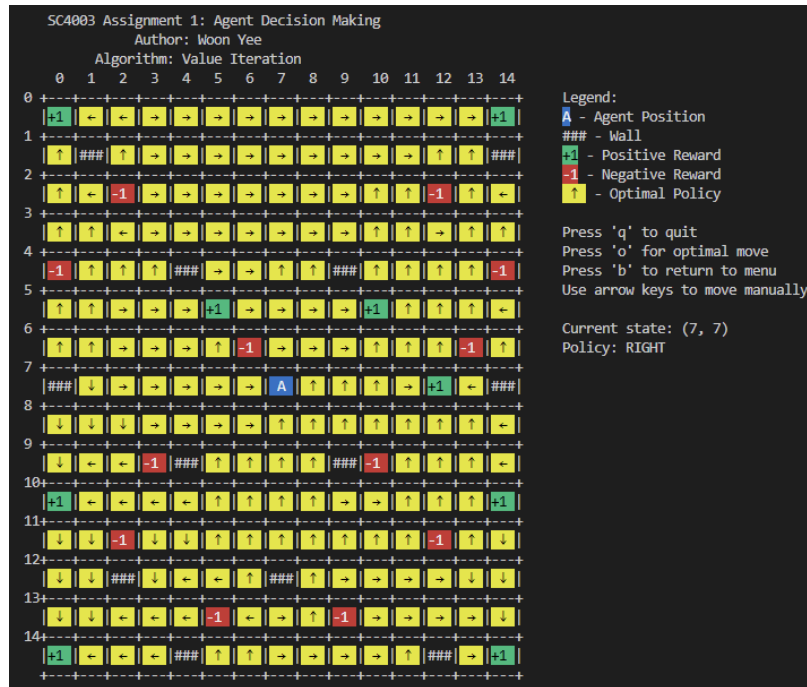https://github.com/woonyee28/SC4003-Intelligent-Agent-Assignment-1

# Part 2:

1. Design a more complicated maze environment of your own and re-run the algorithms designed for Part 1 on it.
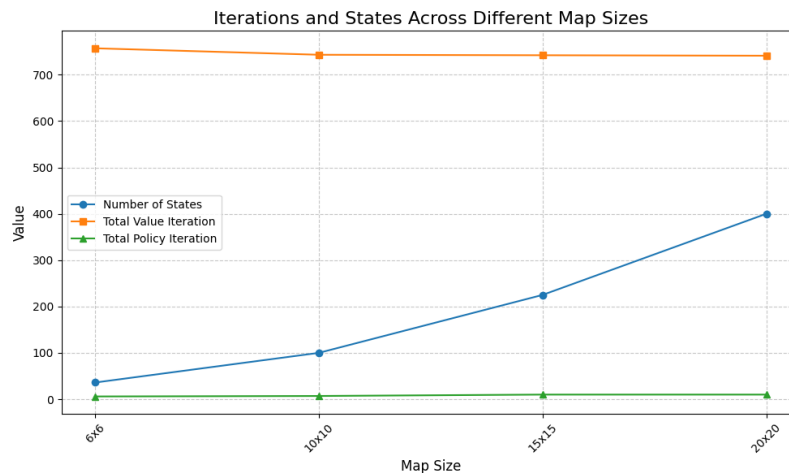
10x10:

**15x15:**



**20x20:**

2. How does the number of states and the complexity of the environment affect convergence?

| Maps | Number of States | Total Value Iteration before convergence | Total Policy Iteration before convergence |
|---|---|---|---|
| 6x6 | 36 | 757 | 6 |
| 10x10 | 100 | 743 | 7 |
| 15x15 | 225 | 742 | 10 |
| 20x20 | 400 | 741 | 10 |



Observation:

- For both Value Iteration and Policy Iteration, the amount of iteration required before convergence remains quite stable even as the map size increases. Policy Iteration constantly outperform Value Iteration in terms of computationally efficiency. I also observed that value iteration reaches the optimal policy way before the convergence, and the rest of the computations are wasteful, and that's why Policy Iteration is preferable in this context.

3. How complex can you make the environment and still be able to learn the right policy?

Observation:

- I have managed to output the map up to the size 20x20. Beyond that the algorithm takes very long to run, especially for Value Iteration. This is due to their respective complexities. Value iteration runs in O(S^2 * A) whereas policy iteration runs in O(S^2) while computing the values. S is the number of states and A is the number of possible actions.