

**NANYANG
TECHNOLOGICAL
UNIVERSITY**
SINGAPORE

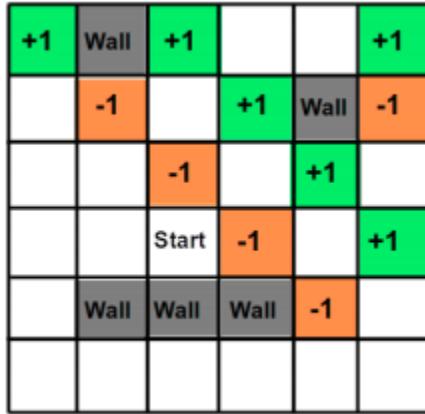
Student Details

| | |
|----------------------|----------------------|
| Name | Acharya Atul |
| Matriculation Number | U1923502C |
| Email | ATUL001@e.ntu.edu.sg |
| Course Code | CZ4046 |
| Course Name | Intelligent Agents |

1. Problem Overview

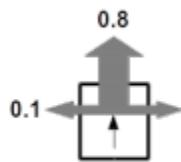
Given a maze environment, convert it to an MDP and solve using both value iteration and policy iteration.

1.1. Problem Description



Maze Environment

The maze environment is a 6x6 grid world which contains walls, rewards and penalties. All green squares have a reward of +1. All orange/red squares have a penalty of -1. All white squares have a reward of -0.04. All grey squares are walls and these are unreachable states. Hence the reward for walls need not be defined and has been set to 0 for this experiment.



Transition Model

The transition model is as follows: the intended outcome occurs with probability 0.8, and with probability 0.1 the agent moves at either right angle to the intended direction. If the move would make the agent walk into a wall, the agent stays in the same place as before. The rewards for the white squares are -0.04, for the green squares are +1, and

for the brown squares are -1. Note that there are no terminal states; the agent's state sequence is infinite.

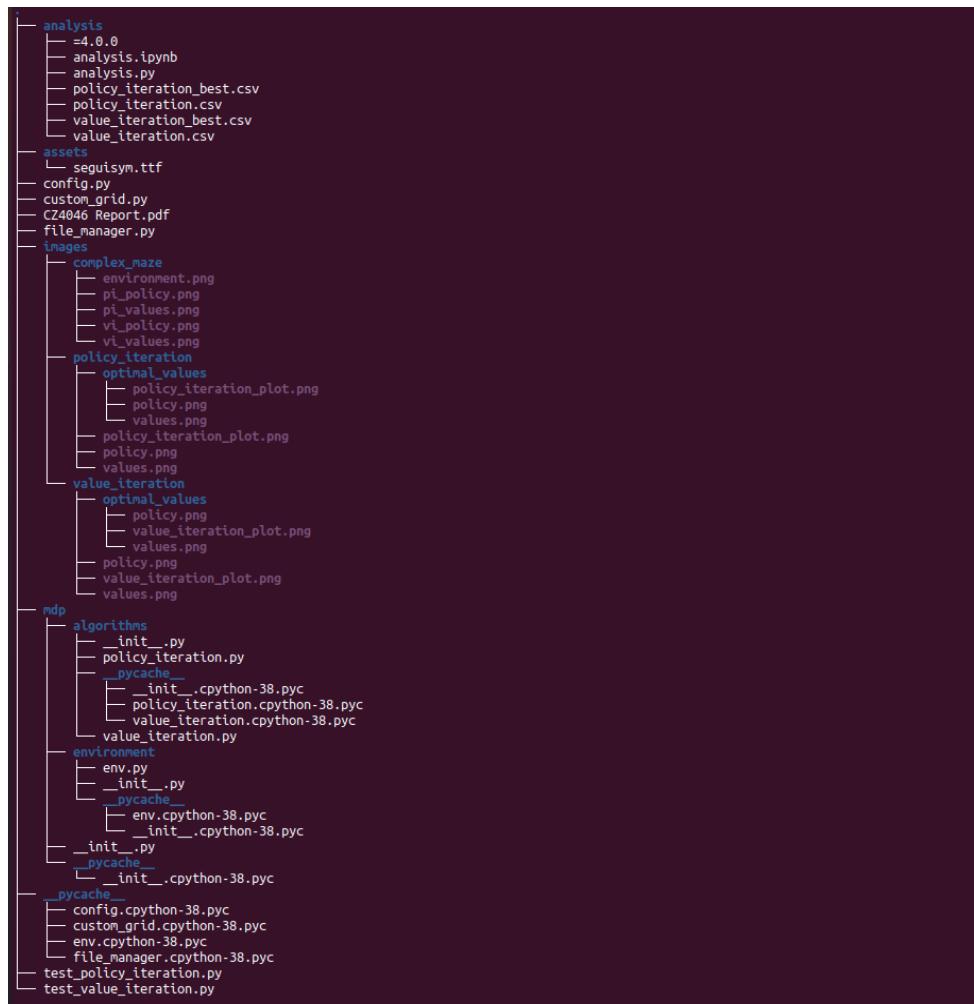
1.2. Expected Solution

From the grid world description, the agent will try to transition to state (0, 0). In this state the agent will choose the action NORTH (goes up). There are walls on all 3 possible directions that the agent can transition to and hence the agent will remain in the same state. The reward for this state is +1. Hence, the agent can accumulate infinite reward by staying in state (0, 0).

2. File Structure and Details

2.1. File Tree

The project folder contains 5 python files and 2 packages. The 5 python files contain code to test the algorithms, configure the maze environment and save data as a .csv file. The project folder consists of 2 other packages - mdp and analysis. The mdp package contains 2 packages - algorithms and environment. The algorithms package contains 2 python files which implement the value and policy iteration algorithms. The environment package contains 1 python file that contains the environment and stores it as an MDP. The analysis package contains a python notebook and 2 .csv files which renders the plot of estimated utilities over number of iterations.



Project File Tree

2.2. File Details

[env.py](#): Contains the Environment class which stores information about the MDP.

[policy_iteration.py](#): Contains the PolicyIteration class which implements the policy iteration algorithm to solve the MDP.

[value_iteration.py](#): Contains the ValueIteration class which implements the value iteration algorithm to solve the MDP.

[filemanager.py](#): Saves the utilities of every state calculated during the execution of the algorithms in each iteration as a .csv file.

[config.py](#): Contains constants for initializing the MDP.

[custom_grid.py](#): Contains variables for users to create a custom environment.

[test_value_iteration.py](#): Driver program to test the value iteration algorithm.

[test_policy_iteration.py](#): Driver program to test the policy iteration algorithm.

[analysis.ipynb](#): Plots the utilities of every state vs iterations for both algorithms.

3. Implementation of MDP

3.1. Grid World

The grid world is a 6x6 array which contains walls, rewards and penalties. It is implemented as a 2-Dimensional list as shown below.

```
grid = [
    ['G', 'W', 'G', '', '', 'G'],
    ['', 'R', '', 'G', 'W', 'R'],
    ['', '', 'R', '', 'G', ''],
    ['', '', '', 'R', '', 'G'],
    ['', 'W', 'W', 'W', 'R', ''],
    ['', '', '', '', ''],
]
```

Grid World representation in code

Cells marked as 'G' have a reward of +1. Cells marked as 'R' have a penalty of -1. Cells marked as 'W' are walls which are unreachable states. Cells that are marked with an empty string are normal cells which have a penalty of -0.04.

3.2. Actions

The possible actions that the agent can take are [UP, DOWN, LEFT RIGHT]. Each of these actions is represented by a tuple of 2 elements. The first element indicates the direction moved along the vertical direction and the second element indicates the direction moved along the horizontal direction.

The tuples of every state are given below

| Action | Tuple |
|--------|---------|
| UP | (-1, 0) |
| DOWN | (+1, 0) |
| LEFT | (0, -1) |
| RIGHT | (0, +1) |

To explain the reason for representing actions as tuples consider the Current State of the agent to be (1,1). Each state is represented in a (row, column) format.

If the agent moves UP, the next state of the agent should be (0, 1).

$$\text{Current State} + \text{UP} = (1, 1) + (-1, 0) = (0, 1)$$

If the agent moves DOWN, the next state of the agent should be (2, 1).

$$\text{Current State} + \text{DOWN} = (1, 1) + (+1, 0) = (2, 1)$$

If the agent moves LEFT, the next state of the agent should be (1, 0).

$$\text{Current State} + \text{UP} = (1, 1) + (0, -1) = (1, 0)$$

If the agent moves RIGHT, the next state of the agent should be (1, 2).

$$\text{Current State} + \text{UP} = (1, 1) + (0, +1) = (1, 2)$$

3.3. Rewards

In this problem, the reward only depends on the current state. Therefore, the rewards can be stored the same as the grid world.

Every green colored cell has a reward of +1. Every red colored cell has a reward of -1. Every white cell has a penalty of -0.04. The reward of the wall states does not matter since it is an unreachable state.

The reward is stored as a 2-Dimensional List as follows.

```
reward_map = {'G': +1, 'R': -1, 'W': 0, ' ': -0.04}  
rewards = [[reward_map[cell] for cell in row] for row in grid]
```

Rewards Represented in Code

3.4. Transition Probabilities

The transition model is as follows: the intended outcome occurs with probability 0.8, and with probability 0.1 the agent moves at either right angle to the intended direction. If the move would make the agent walk into a wall, the agent stays in the same place as before.

The transition probability is modeled as a function in the Environment Class. The function takes 2 arguments - current state and action and returns the transition probabilities.

```
def transition_model(self, state, action):
    model = {}

    possible_directions = [action, (action[1], action[0]), (-action[1], -action[0])]

    # Probability for going in the same direction is 0.8, probability for going in each perpendicular direction is 0.1
    probability = [0.8, 0.1, 0.1]

    # Loop through all possible directions the agent can go to by taking the action
    for idx, direction in enumerate(possible_directions):

        # Calculate next state
        next_state = (state[0] + direction[0], state[1] + direction[1])

        # Check if next state is within the boundaries of the grid
        if 0 <= next_state[0] < self.grid_width and 0 <= next_state[1] < self.grid_height:

            # Check if next state is a wall
            if self.grid_world[next_state[0]][next_state[1]] == 'W':
                next_state = state
            else:
                next_state = state

        # Check if next state is already in the transition model
        # If true then add the probabilities
        if next_state in model:
            model[next_state] += probability[idx]
        else:
            model[next_state] = probability[idx]

    return model
```

Transition Model Represented in Code

Explanation of possible_directions:

The intended outcome occurs with probability 0.8, and with probability 0.1 the agent moves at either right angle to the intended direction. If the move would make the agent walk into a wall, the agent stays in the same place as before.

The actions are represented by the following tuple:

| Action | Tuple |
|--------|---------|
| UP | (-1, 0) |
| DOWN | (+1, 0) |
| LEFT | (0, -1) |
| RIGHT | (0, +1) |

possible_directions = [action, (action[1], action[0]), (-action[1], -action[0])]

Evaluating this for each action gives:

| Action | Possible Directions Tuples | Possible Directions | Transition Probabilities |
|--------|-------------------------------|---------------------|-----------------------------|
| UP | (-1, 0), (0, -1), (0, 1) | UP, LEFT, RIGHT | [0.8, 0.1, 0.1] |
| DOWN | (1, 0), (0, 1), (0, -1) | DOWN, RIGHT, LEFT | [0.8, 0.1, 0.1] |
| LEFT | (0, -1), (-1, 0), (1, 0) | LEFT, UP, DOWN | [0.8, 0.1, 0.1] |
| RIGHT | (0, 1), (1, 0), (-1, 0) | RIGHT, DOWN, UP | [0.8, 0.1, 0.1] |

4. Value Iteration

The algorithm aims to calculate the utility of every state. Once the utilities converge, the algorithm computes the optimal policy by taking actions that maximize expected utility.

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.

The starting utility of every state is initialized to 0 which are then updated using the Bellman equations. The Bellman equation relates the utility of a state to the utility of neighbouring states. The Bellman Update is performed until the utility of every state converges.

The Bellman update is given by the following formula

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

4.1. Pseudo-code for Value Iteration

The Pseudo-code from “Artificial Intelligence: A Modern Approach” by S. Russell and P. Norvig. Prentice-Hall, third edition, 2010 was used to implement the algorithm.

```

function VALUE-ITERATION(mdp,  $\epsilon$ ) returns a utility function
  inputs: mdp, an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ ,
           rewards  $R(s)$ , discount  $\gamma$ 
            $\epsilon$ , the maximum error allowed in the utility of any state
  local variables:  $U$ ,  $U'$ , vectors of utilities for states in  $S$ , initially zero
                      $\delta$ , the maximum change in the utility of any state in an iteration

  repeat
     $U \leftarrow U'$ ;  $\delta \leftarrow 0$ 
    for each state  $s$  in  $S$  do
       $U'[s] \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s']$ 
      if  $|U'[s] - U[s]| > \delta$  then  $\delta \leftarrow |U'[s] - U[s]|$ 
    until  $\delta < \epsilon(1 - \gamma)/\gamma$ 
  return  $U$ 

```

Pseudo-code followed

4.2. Implementation

The Value Iteration algorithm is implemented in `value_iteration.py`. When `test_value_iteration.py` is run, the following steps occur:

1. The environment (grid world) is created
2. Constants for Value Iteration algorithm are initialized
3. The MDP is solved using Value Iteration
4. A plot of the optimal policy and final utilities is displayed
5. The data required for analysis is saved as a .csv file

A set of constants are fixed in `test_value_iteration.py`

| | |
|-------------------------|------|
| GAMMA (Discount Factor) | 0.99 |
| C | 0.1 |
| MAX_REWARD (Rmax) | 1.0 |

The discount factor is set to 0.99. The maximum reward for any set is +1. The value of C is set to 0.1.

Epsilon is calculated using these values. Epsilon is the maximum error allowed in the utility of any state.

$$\text{epsilon} = C * R_{\max} = 0.1 * 1 = 0.1$$

The threshold is calculated using the formula:

$$\text{threshold} = \text{epsilon} * (1 - \gamma) / \gamma = 0.1 * (1 - 0.99) / 0.99 = 0.00101$$

A 2-Dimensional array is used to store the utility of every state. It is initialized to 0 and every cell is updated using the Bellman update

4.2.1. Bellman Update

The solve method of the Value Iteration Class implements the value iteration algorithm. It loops through all the states and updates the expected utility of every state using the bellman update. “action_values” is a list that stores the expected utility for every possible action. The state_value is updated using the statement:

$$\text{state_value} = \text{reward} + \gamma * \max(\text{action_values})$$

Below is a code snippet that implements the Bellman update. The following Bellman update is performed for every state. The algorithm loops through all the states and the following code is a wrapper in this loop.

```

for action in actions:

    action_value = 0

    # Get the transition model
    transition_model = mdp.transition_model(cur_state, action)

    # Loop through all the next states in the transition model
    # This loop calculates expected utility for taking the action
    for next_state in transition_model:
        utility = utilities[next_state[0]][next_state[1]]
        probability = transition_model[next_state]
        expected_utility = probability * utility
        action_value += expected_utility

    action_values.append(action_value)

    # Reward for the current state
    reward = mdp.receive_reward(cur_state)

    # Bellman Update
    state_value = reward + self.gamma * max(action_values)
    new_utilities[i][j] = state_value

    # Update the value of delta
    delta = max(delta, abs(new_utilities[i][j] - utilities[i][j]))

```

Bellman Update for each state

4.2.2. Convergence Condition

The algorithm runs until the utilities of every state converges. The delta value is originally set to 0 and is updated as:

$\text{delta} = \max(\text{delta}, \text{new_utilities(state)} - \text{old_utilities(state)})$

The values converge once delta is less than the threshold (0.00101).

Below is a code snippet that shows the stopping condition of the algorithm.

```

# Check for convergence
if delta < threshold:
    break

```

Convergence Criteria

4.2.3. Calculation of Optimal Policy

Once the utilities of every state converge, the algorithm can then find the optimal policy. The optimal policy is calculated by selecting the action that maximizes the expected utility. The optimal policy is calculated in the greedify method of the Value Iteration Class. For every state the action value for every action is calculated and the action that has the action value is deemed to be the best action to take for that state.

The following code snippet calculates the optimal policy. The method loops through every state and selects the best action using the code snippet as shown below.

```
# Loop through all actions
for action in actions:
    action_value = 0
    # Get transition model
    transition_model = mdp.transition_model(cur_state, action)

    # Loop through all the next states in the transition model
    # Calculate the expected utility for taking the action
    for next_state in transition_model:
        utility = utilities[next_state[0]][next_state[1]]
        probability = transition_model[next_state]
        expected_utility = probability * utility
        action_value += expected_utility
    action_values[action] = action_value
best_action = None
best_action_value = -math.inf

# Choose best action
for action in action_values:
    if action_values[action] > best_action_value:
        best_action = action
        best_action_value = action_values[action]
policy[i][j] = best_action
```

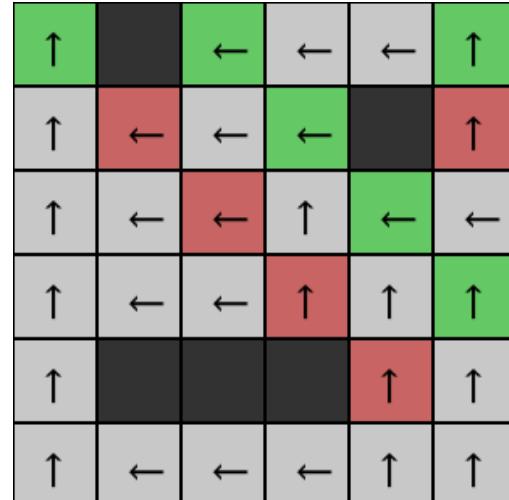
Greedification to select optimal actions

4.3. Final Utilities and Optimal Policy

Plots of final utilities and optimal policy are shown below.

| | | | | | |
|--------|--------|--------|--------|--------|--------|
| 99.901 | | 94.946 | 93.776 | 92.555 | 93.229 |
| 98.294 | 95.784 | 94.446 | 94.298 | | 90.819 |
| 96.849 | 95.487 | 93.195 | 93.077 | 93.003 | 91.696 |
| 95.455 | 94.353 | 93.133 | 91.016 | 91.715 | 91.789 |
| 94.213 | | | | 89.449 | 90.467 |
| 92.838 | 91.629 | 90.436 | 89.257 | 88.470 | 89.198 |

Plot of Final Utilities



Plot of Optimal Policy

The final utilities when $C = 0.1$ are shown. Every state is represented as (col, row).

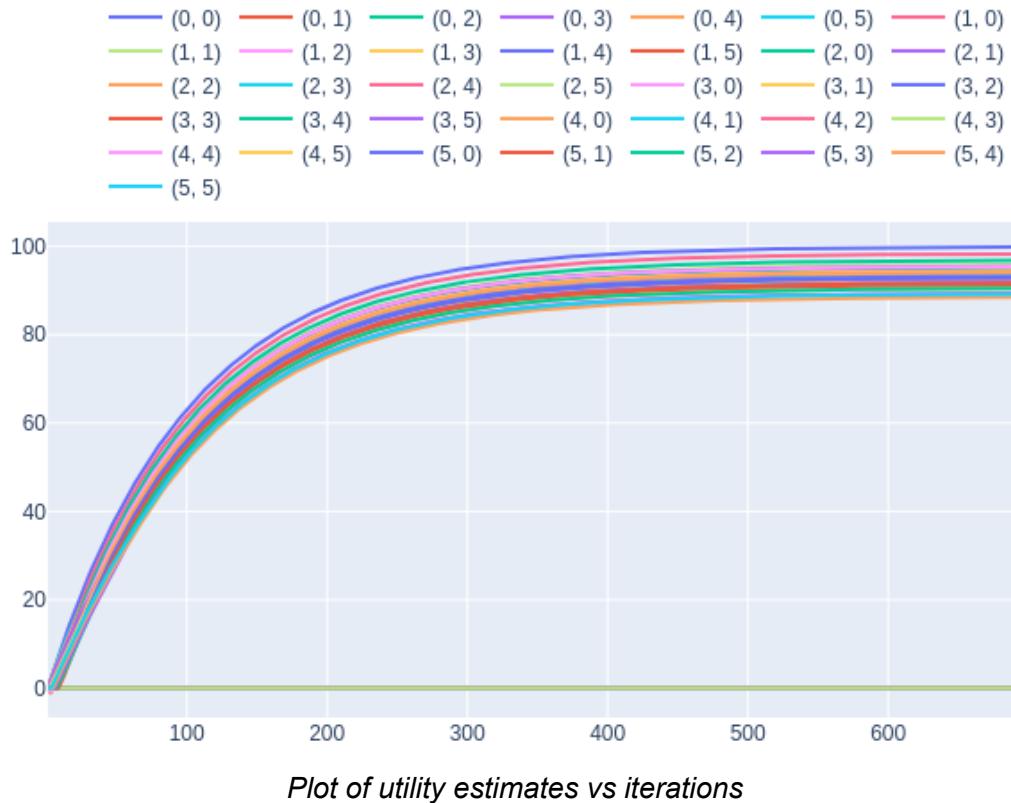
Number of iterations: 688

- (0, 0): 99.9006852204078
- (1, 0): 98.29404673107541
- (2, 0): 96.84918540226495
- (3, 0): 95.45452432190575
- (4, 0): 94.21320463236087
- (5, 0): 92.83815953740235
- (0, 1): 0.0
- (1, 1): 95.78370260537261
- (2, 1): 95.48711297196897
- (3, 1): 94.35317902261747
- (4, 1): 0.0
- (5, 1): 91.62946285015803
- (0, 2): 94.94614245455695
- (1, 2): 94.44568358931949
- (2, 2): 93.19511283494086
- (3, 2): 93.13323064237099
- (4, 2): 0.0
- (5, 2): 90.43583719392673
- (0, 3): 93.77568621731811

(1, 3): 94.29840005646918
 (2, 3): 93.07695824252093
 (3, 3): 91.01594175092443
 (4, 3): 0.0
 (5, 3): 89.25709465061603
 (0, 4): 92.55529967122197
 (1, 4): 0.0
 (2, 4): 93.00305429226383
 (3, 4): 91.71509229342458
 (4, 4): 89.44909832058741
 (5, 4): 88.46978429758704
 (0, 5): 93.22918825371615
 (1, 5): 90.8186084153956
 (2, 5): 91.69555629320722
 (3, 5): 91.7887697818793
 (4, 5): 90.46745089144353
 (5, 5): 89.19837580875523

4.4. Plot of Estimated Utilities vs Iterations

A plot of Utility estimates as a function of iterations is shown below



The above graph was plotted with $C = 0.1$. As the number of iterations increases, the utility estimates also increases. However, the utility estimates slowly converge and the algorithm takes 688 iterations to find the optimal policy.

From the above plots, it is clear that the algorithm has found the correct policy which is to reach the state $(0, 0)$. The utility of state $(0, 0)$ is also the highest.

4.5. Calculation of best value for C

In order to reduce the number of Bellman updates required to find the optimal policy, the value of C was varied until the algorithm could not find the optimal policy

Results:

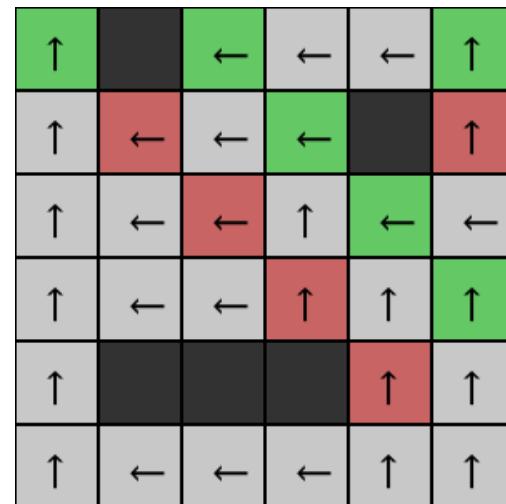
| C | Iterations | Optimal Policy |
|----------|-------------------|-----------------------|
| 0.1 | 688 | Yes |
| 1 | 459 | Yes |
| 10 | 230 | Yes |
| 25 | 138 | Yes |
| 40 | 92 | Yes |
| 45 | 80 | Yes |
| 50 | 69 | No |

Plots:

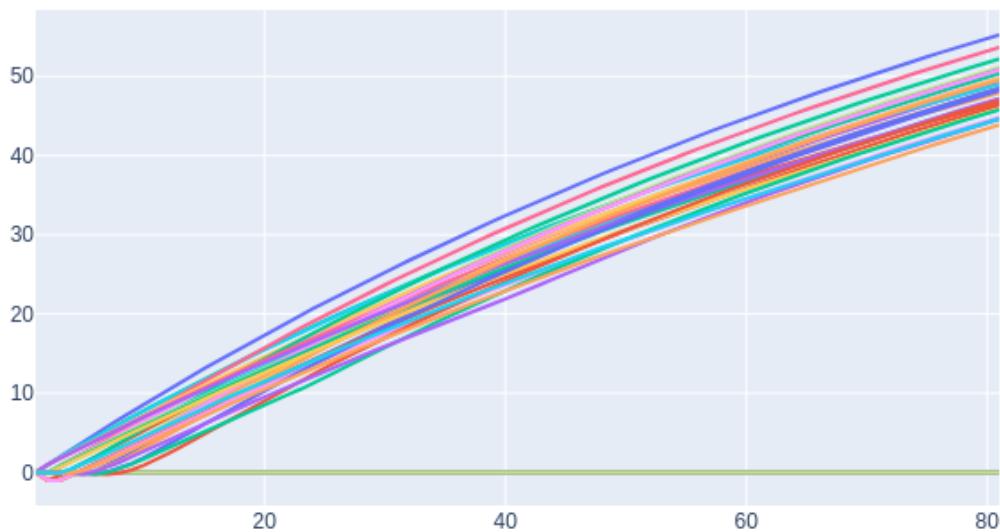
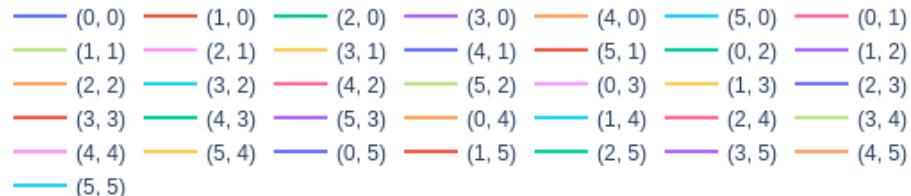
Here are the corresponding plots of optimal policy and final utilities with C = 45

| | | | | | |
|--------|--------|--------|--------|--------|--------|
| 55.248 | | 50.326 | 49.156 | 47.940 | 48.999 |
| 53.641 | 51.131 | 49.796 | 49.653 | | 46.642 |
| 52.196 | 50.834 | 48.543 | 48.432 | 48.360 | 47.111 |
| 50.802 | 49.700 | 48.480 | 46.372 | 47.079 | 47.205 |
| 49.560 | | | | 44.821 | 45.884 |
| 48.185 | 46.976 | 45.783 | 44.604 | 43.846 | 44.616 |

Plot of Final Utilities



Plot of Optimal Policy



Plot of utility estimates vs iterations

5. Policy Iteration

The algorithm aims to calculate the optimal policy. The algorithm starts with an initial policy and alternates between 2 steps: policy evaluation and policy improvement.

The policy improvement step calculates the utilities of every state using the current policy. Once the Bellman update is performed 'k' times, the algorithm moves on to the policy improvement step. The current policy is updated using the utilities calculated during the policy evaluation step by taking the action that maximizes the expected utility. These 2 steps are alternated and the algorithm stops once the current policy is stable.

The policy evaluation step uses the following Bellman Update:

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

The policy improvement step uses the following equation:

$$\pi^{i+1}(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a)U^{\pi_i}(s')$$

5.1. Pseudo-code for Policy Iteration

The Pseudo-code from "Artificial Intelligence: A Modern Approach" by S. Russell and P. Norvig. Prentice-Hall, third edition, 2010 was used to implement the algorithm.

```

function POLICY-ITERATION(mdp) returns a policy
  inputs: mdp, an MDP with states  $S$ , actions  $A(s)$ , transition model  $P(s' | s, a)$ 
  local variables:  $U$ , a vector of utilities for states in  $S$ , initially zero
     $\pi$ , a policy vector indexed by state, initially random

  repeat
     $U \leftarrow \text{POLICY-EVALUATION}(\pi, U, mdp)$ 
    unchanged?  $\leftarrow$  true
    for each state  $s$  in  $S$  do
      if  $\max_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s'] > \sum_{s'} P(s' | s, \pi[s]) U[s']$  then do
         $\pi[s] \leftarrow \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s' | s, a) U[s']$ 
        unchanged?  $\leftarrow$  false
    until unchanged?
  return  $\pi$ 

```

5.2. Implementation

The Policy Iteration algorithm is implemented in `policy_iteration.py`. When `test_policy_iteration.py` is run, the following steps occur:

1. The environment (grid world) is created
2. Constants for Policy Iteration algorithm are initialized
3. The MDP is solved using Policy Iteration
4. A plot of the optimal policy and final utilities is displayed
5. The data required for analysis is saved as a .csv file

A set of constants are fixed in `test_value_iteration.py`

| | |
|-------------------------|------|
| GAMMA (Discount Factor) | 0.99 |
| K | 100 |

The value of 'k' represents the number of iterations of Bellman updates in the policy evaluation step. A large value of k can slow down the algorithm but it has a good chance to find the optimal policy. A small value of k can speed up the algorithm but it may not find the optimal policy. In this assignment, the value of k is set to 100.

5.2.1. Initial Policy

The starting policy can be any policy. In this case, the initial policy is set to ‘UP’ for every state. The choice of starting policy affects the number of iterations. The algorithm initializes the starting policy in the *get_starting_policy* method.

```
policy = [[(-1, 0) for _ in range(mdp.grid_width)] for _ in range(mdp.grid_height)]
return policy
```

Initialization of starting policy

5.2.2. Solve method

The MDP is solved using the *solve* method of the *PolicyIteration* class. The method takes an MDP as an argument and performs a series of policy evaluation and policy improvement steps until the policy is stable. The *policy_evaluation* method returns the new utility values and this is passed as a parameter to the *policy_improvement* method. The *policy_improvement* method returns the new policy and a boolean to determine if the policy is stable. The policy is said to be stable if the new policy is the same as the old policy.

```
def solve(self, mdp):
    # Initialize the starting policy
    policy = self.get_starting_policy(mdp)

    # Initialize the starting utilities
    utilities = np.zeros((mdp.grid_height, mdp.grid_width), dtype=np.float)

    # Initialize the analysis data
    for i in range(utilities.shape[0]):
        for j in range(utilities.shape[1]):
            cur_state = (j, i)
            self.data[f'{cur_state}'] = [0]

    # Initialize the total_iterations
    total_iterations = 0

    # Loop control variable
    is_policy_stable = False

    # Loop while the policy is not stable
    while not is_policy_stable:
        # Policy Evaluation
        utilities, iterations = self.policy_evaluation(policy, utilities, mdp)
        total_iterations += iterations

        # Policy Improvement
        policy, is_policy_stable = self.policy_improvement(policy, utilities, mdp)

    # Return utilities, policy and iterations as a dictionary
    return {"utilities": utilities, "policy": policy, "iterations": total_iterations}
```

Policy Iteration algorithm

5.2.3. Policy Evaluation

The policy evaluation step is performed in the *policy_evaluation* method of the PolicyIteration Class. This method takes the current policy, current utilities and the MDP as arguments and performs updates to the utilities of every state using the current policy. The method loops through all the states and updates the utility of every state by taking the action according to its policy. The algorithm performs the update 'k' times and returns the new utilities. In this assignment the value of 'k' is set to 100.

Below is the code snippet explaining the Bellman update performed in policy evaluation.

```
# Get action from current policy
action = policy[i][j]

# Get transition model of the MDP
transition_model = mdp.transition_model(cur_state, action)

action_value = 0

# Loop through all the next states in the transition model
# This loop calculates expected utility for taking the action
for next_state in transition_model:
    utility = utilities[next_state[0]][next_state[1]]
    probability = transition_model[next_state]
    expected_value = probability * utility
    action_value += expected_value

# Reward of current state
reward = mdp.receive_reward(cur_state)

# Bellman Update
utility = reward + self.gamma * action_value
new_utilities[i][j] = utility

# Update analysis data
self.data[f'{state_format}'].append(utility)

utilities = new_utilities.copy()

return utilities, iteration
```

Policy Improvement using Bellman update

5.2.4. Policy Improvement

The policy improvement step is performed in the *policy_improvement* method of the PolicyIteration Class. This method takes the current policy, current utilities and the MDP as arguments and performs updates to the policy using the current utilities. The algorithm updates the policy by taking the greedy action (action that maximizes expected utility) and returns the new policy and a boolean depicting whether the policy

is stable. The policy is said to be stable when the old policy is the same as the new policy. Once the policy is stable, the policy iteration algorithm stops and returns the final policy and utilities.

Below is the code snippet that explains the policy evaluation algorithm. The first explains the policy update step and the second snippet explains the logic to check is the policy is stable

```
# Get all possible actions
actions = mdp.actions
action_values = {}

# Loop through all possible actions
# This loop calculates action values for all actions
for action in actions:
    action_value = 0

    # Get transition model
    transition_model = mdp.transition_model(cur_state, action)

    # Loop through all the next states in the transition model
    # This loop calculates expected utility for taking the action
    for next_state in transition_model:
        utility = utilities[next_state[0]][next_state[1]]
        probability = transition_model[next_state]
        expected_utility = probability * utility
        action_value += expected_utility

    # Set action value for the action
    action_values[action] = action_value

best_action = None
best_action_value = -math.inf

# This loop chooses the action that maximizes expected utility
for action in action_values:
    if action_values[action] > best_action_value:
        best_action = action
        best_action_value = action_values[action]
new_policy[i][j] = best_action
```

Policy Update by selecting Greedy Action

```
# Checks if the old policy is same as the new policy
# If the old policy is same as the new policy, the policy is stable
for i in range(len(policy)):
    for j in range(len(policy[i])):
        if policy[i][j] != new_policy[i][j]:
            return new_policy, False

return new_policy, True
```

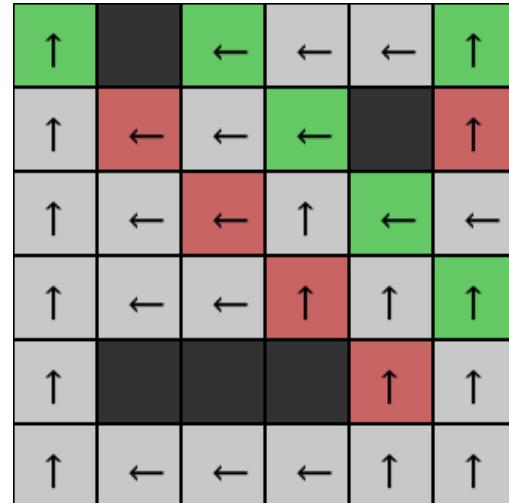
Check is Policy is Stable

5.3. Final Utilities and Optimal Policy

Plots of the final utilities and Optimal Policy are shown below.

| | | | | | |
|--------|--------|--------|--------|--------|--------|
| 99.343 | | 94.388 | 93.218 | 91.998 | 92.671 |
| 97.736 | 95.226 | 93.888 | 93.741 | | 90.261 |
| 96.291 | 94.929 | 92.637 | 92.519 | 92.445 | 91.138 |
| 94.897 | 93.795 | 92.575 | 90.458 | 91.157 | 91.231 |
| 93.655 | | | | 88.891 | 89.910 |
| 92.280 | 91.072 | 89.878 | 88.699 | 87.912 | 88.641 |

Plot of Final Utilities



Plot of Optimal Policy

The final utilities when $k = 100$ are shown. Every state is represented as (col, row).

Number of iterations: 500

(0, 0): 99.34295169575849

(1, 0): 97.73631320642609

(2, 0): 96.29145187761563

(3, 0): 94.89679079725641

(4, 0): 93.65547110771155

(5, 0): 92.28042601275303

(0, 1): 0.0

(1, 1): 95.2259690807233

(2, 1): 94.92937944731963

(3, 1): 93.79544549796816

(4, 1): 0.0

(5, 1): 91.07172932550871

(0, 2): 94.38840892990764

(1, 2): 93.88795006467016

(2, 2): 92.63737931029154

(3, 2): 92.57549711772167

(4, 2): 0.0

(5, 2): 89.87810366927741

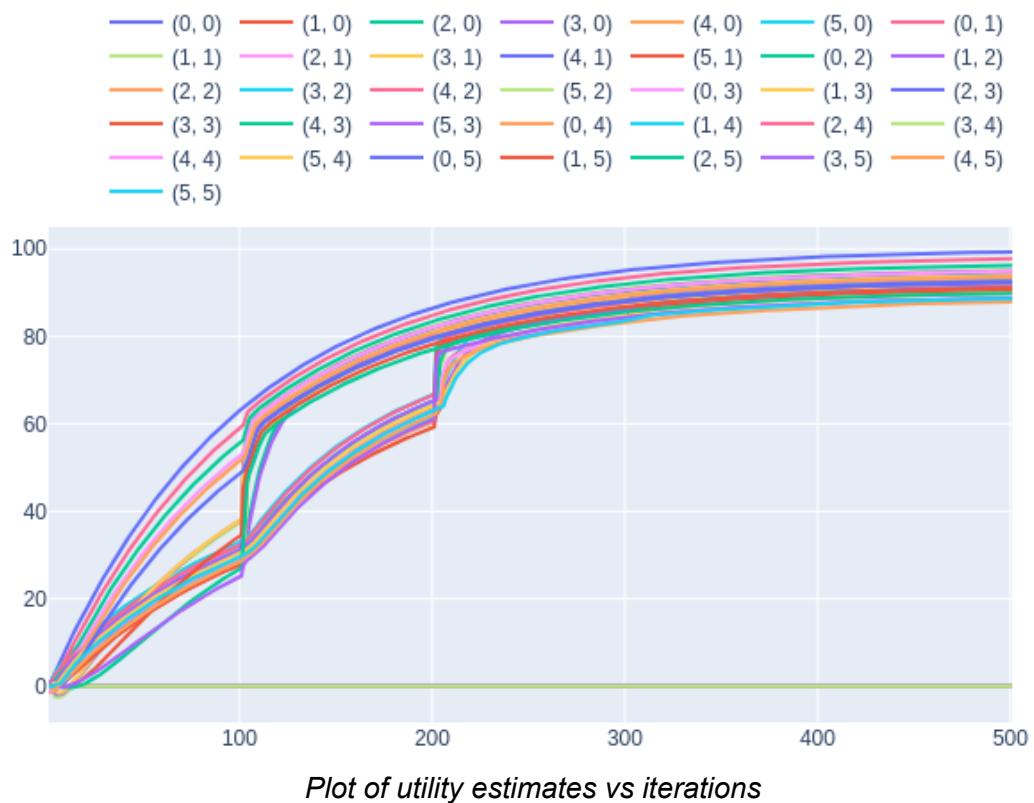
(0, 3): 93.21795269266882

(1, 3): 93.74066653181987

(2, 3): 92.5192247178716
 (3, 3): 90.45820822627492
 (4, 3): 0.0
 (5, 3): 88.69936112596669
 (0, 4): 91.99756614657267
 (1, 4): 0.0
 (2, 4): 92.44532076761429
 (3, 4): 91.15735876877382
 (4, 4): 88.89136479593525
 (5, 4): 87.91205077293387
 (0, 5): 92.67145472899253
 (1, 5): 90.2608748906614
 (2, 5): 91.13782276854705
 (3, 5): 91.23103625721895
 (4, 5): 89.90971736678283
 (5, 5): 88.64064228409407

5.4. Plot of Estimated Utilities vs Iterations

A plot of Utility estimates as a function of iterations is shown below



From the above plots, it is clear that the algorithm found the correct policy which is to reach the state (0, 0) and receive infinite reward. The plot of utility estimates vs iterations shows that the policy changes after every 100 iterations and there is a sharp point at positions where iterations = {100, 200}. After iteration 300, the graph suggests that the policy doesn't change much and it is slowly becoming stable. This could suggest that the smoothness of the curve after every policy improvement step could suggest that the policy is becoming more stable.

5.5. Calculation of best value for K

In order to reduce the number of Bellman updates required to find the optimal policy, the value of K was varied until the algorithm could not find the optimal policy

Results:

| K | Iterations | Optimal Policy |
|-----|------------|----------------|
| 100 | 500 | Yes |
| 75 | 375 | Yes |
| 50 | 250 | Yes |
| 10 | 60 | Yes |
| 5 | 25 | No |

Note

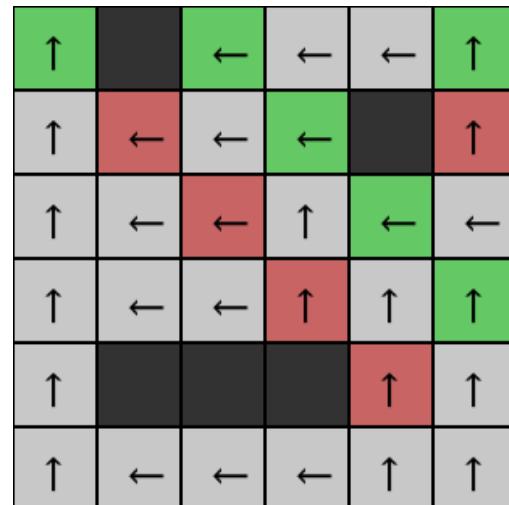
The above results may differ and will depend on the initial policy chosen. When the starting policy was to take the action 'DOWN' for every state, the results differed and the lowest value of 'K' was 13. The number of Bellman updates also increased and this suggests that the choice of starting policy affects the algorithm.

Plots:

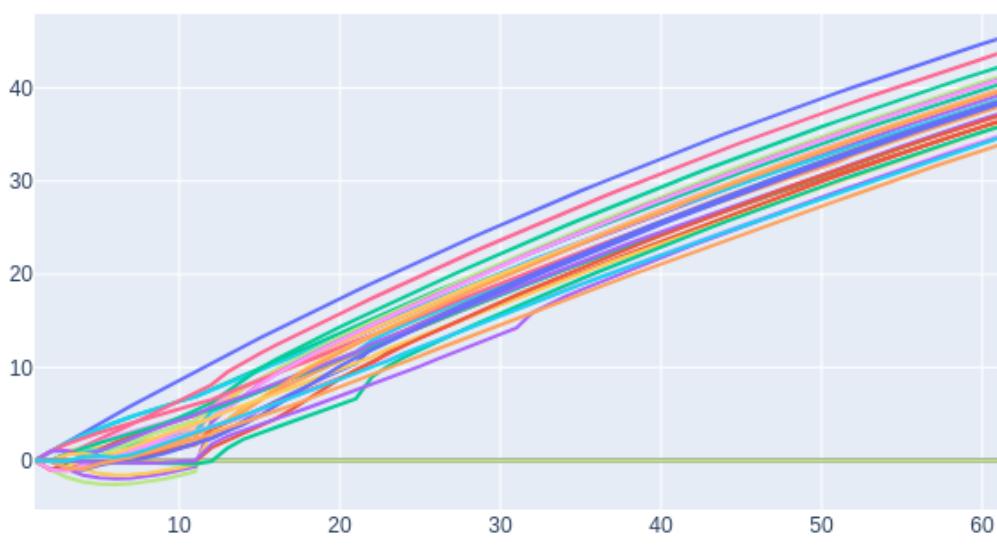
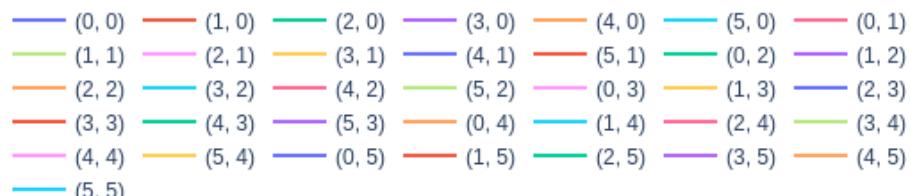
Here are the corresponding plots of optimal policy and final utilities with k = 10

| | | | | | |
|--------|--------|--------|--------|--------|--------|
| 45.284 | | 40.358 | 39.188 | 37.971 | 38.749 |
| 43.678 | 41.167 | 39.832 | 39.689 | | 36.352 |
| 42.233 | 40.871 | 38.579 | 38.467 | 38.394 | 37.104 |
| 40.838 | 39.737 | 38.517 | 36.406 | 37.108 | 37.197 |
| 39.597 | | | | 34.844 | 35.875 |
| 38.222 | 37.013 | 35.819 | 34.641 | 33.864 | 34.603 |

Plot of Final Utilities



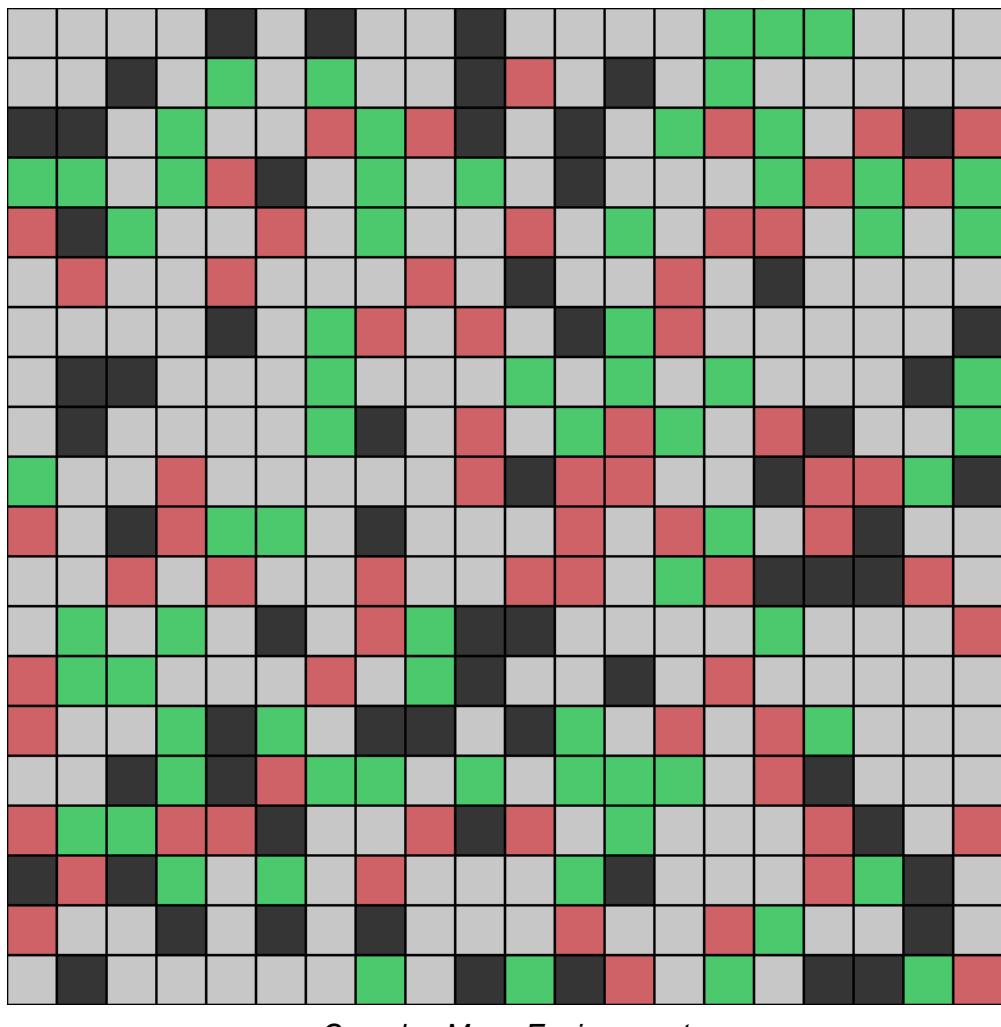
Plot of Optimal Policy



Plot of utility estimates vs iterations

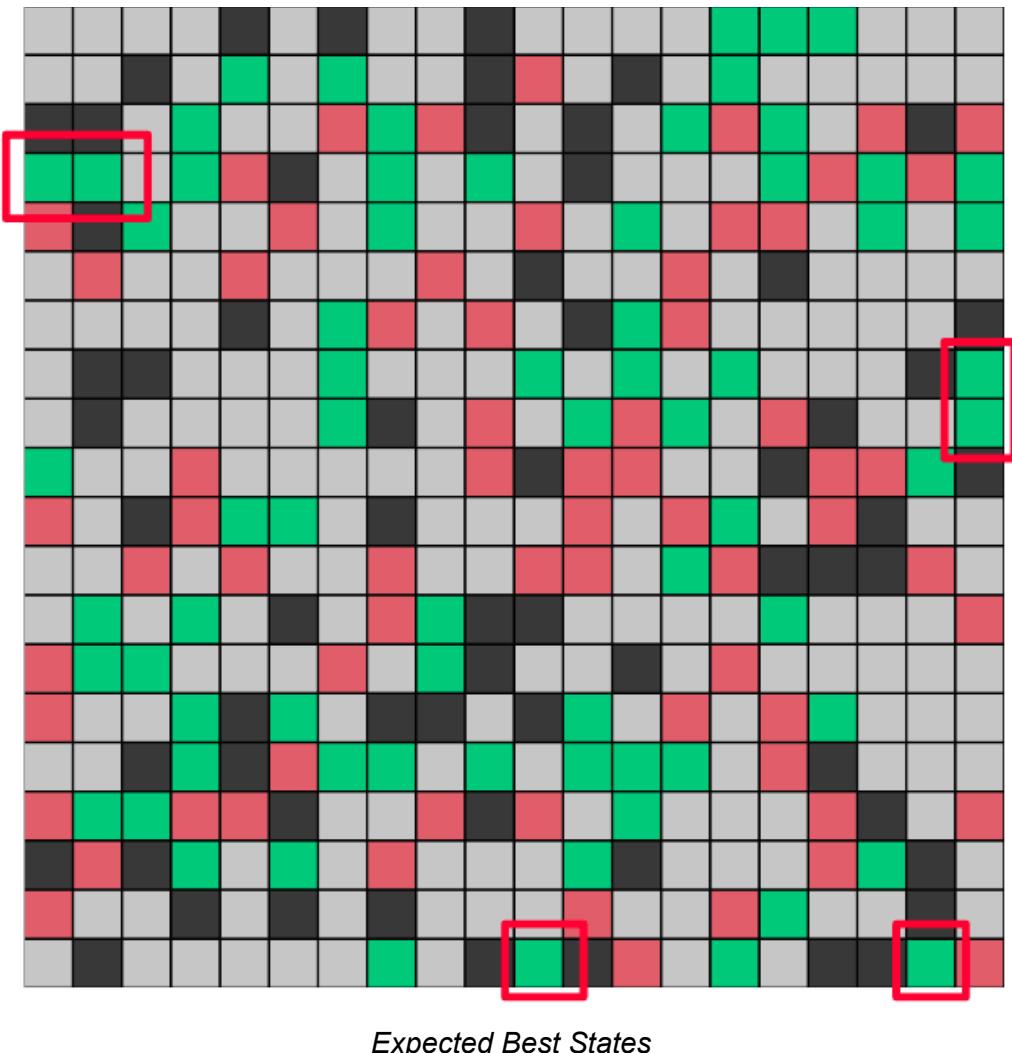
6. Complex Environment

A more complex maze environment was designed to test the value iteration and policy iteration algorithm algorithms. The complex environment can be found in *custom_grid.py*. The size of the grid is 20x20 and the transition probabilities and rewards are the same. Increasing the dimensions of the grid results in more states and this implies longer calculation time and more memory required for the 2 algorithms. The maze environment designed is shown below.



6.1. Observations

The maze has states where the agent can receive infinite reward. Some of the states are marked below



In the states highlighted above, the agent can either take an action that results in the agent staying in the same state and receiving rewards over and over again or it can take actions that toggle between such states.

6.2. Value Iteration Final Value

The final utilities calculated by the value iteration algorithm are shown below.

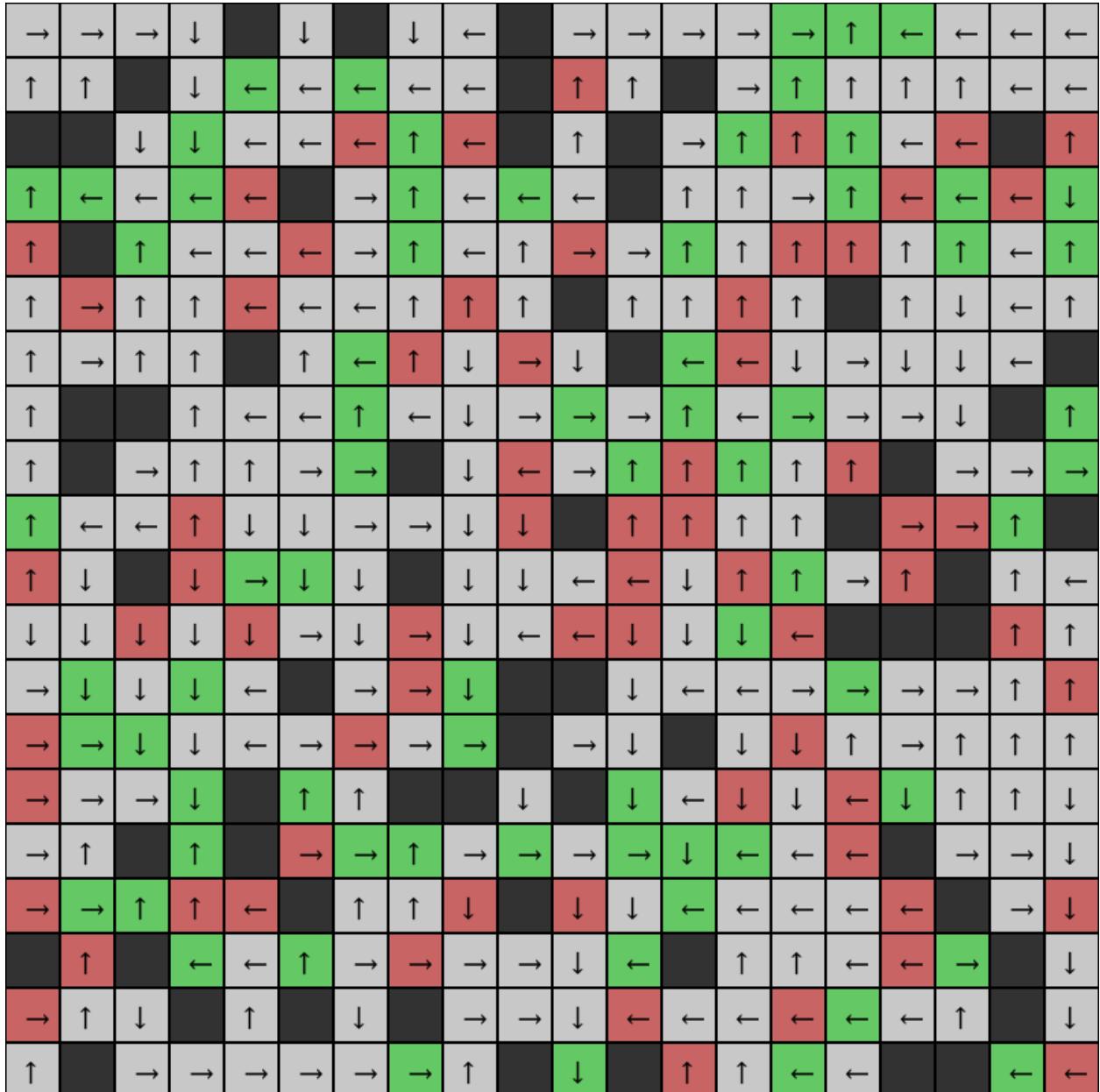
| | | | | | | | | | | | | | | | | | | | |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 91.381 | 92.719 | 94.108 | 95.348 | | 94.137 | | 92.389 | 91.306 | | 92.797 | 94.288 | 95.716 | 96.977 | 98.281 | 98.292 | 98.140 | 96.713 | 95.335 | 93.970 |
| 90.321 | 91.381 | | 96.759 | 96.784 | 95.377 | 95.178 | 93.743 | 92.280 | | 90.656 | 92.797 | | 96.754 | 98.028 | 97.089 | 96.750 | 95.475 | 94.354 | 92.955 |
| | | 97.257 | 98.029 | 96.670 | 95.409 | 93.107 | 93.498 | 91.287 | | 89.474 | | 95.138 | 96.529 | 95.780 | 96.834 | 95.537 | 93.421 | | 90.548 |
| 99.901 | 99.901 | 98.440 | 98.272 | 95.858 | | 92.133 | 93.272 | 91.936 | 91.892 | 90.526 | | 94.035 | 95.122 | 94.977 | 96.402 | 94.012 | 93.991 | 91.660 | 92.724 |
| 97.407 | | 98.299 | 97.005 | 95.525 | 92.984 | 91.769 | 93.024 | 91.592 | 90.736 | 90.205 | 92.567 | 93.923 | 93.780 | 92.785 | 93.740 | 92.982 | 93.791 | 92.384 | 92.766 |
| 95.918 | 94.114 | 96.593 | 95.598 | 93.418 | 92.157 | 91.048 | 91.504 | 89.464 | 89.543 | | 91.483 | 92.461 | 91.492 | 91.568 | | 91.860 | 92.577 | 91.599 | 91.563 |
| 94.586 | 93.896 | 95.106 | 94.437 | | 90.969 | 91.066 | 89.246 | 88.490 | 87.925 | 89.867 | | 92.940 | 90.715 | 91.168 | 91.811 | 92.929 | 94.010 | 92.656 | |
| 93.356 | | | 93.042 | 91.688 | 90.474 | 90.962 | 89.718 | 89.635 | 89.756 | 91.297 | 91.492 | 92.715 | 91.417 | 92.347 | 92.557 | 94.159 | 95.553 | | 99.901 |
| 92.141 | | 90.278 | 91.548 | 90.532 | 89.825 | 90.957 | | 90.792 | 88.746 | 90.295 | 91.361 | 90.521 | 91.380 | 91.080 | 90.257 | | 96.986 | 98.570 | 99.901 |
| 92.074 | 90.751 | 89.647 | 89.228 | 89.370 | 89.604 | 90.049 | 91.045 | 92.246 | 90.440 | | 88.911 | 88.402 | 89.984 | 89.903 | | 93.191 | 95.951 | 98.297 | |
| 89.686 | 89.752 | | 88.630 | 90.539 | 90.761 | 90.683 | | 93.839 | 92.619 | 91.414 | 89.124 | 88.683 | 87.932 | 89.763 | 89.454 | 90.635 | | 96.846 | 95.437 |
| 89.942 | 90.945 | 89.593 | 90.774 | 88.987 | 90.683 | 91.870 | 92.976 | 95.228 | 93.839 | 91.420 | 90.129 | 89.893 | 89.875 | 87.874 | | | | 94.371 | 94.215 |
| 91.004 | 92.439 | 91.671 | 92.343 | 90.940 | | 93.092 | 94.460 | 96.937 | | | 92.399 | 91.053 | 89.996 | 88.969 | 90.193 | 90.206 | 91.505 | 92.859 | 91.910 |
| 90.188 | 92.620 | 92.701 | 92.507 | 91.263 | 91.979 | 93.181 | 95.783 | 97.209 | | 92.566 | 93.786 | | 91.163 | 89.667 | 89.113 | 89.337 | 90.327 | 91.434 | 90.792 |
| 89.139 | 91.421 | 92.644 | 93.858 | | 92.067 | 91.979 | | | 92.748 | | 95.175 | 94.078 | 92.553 | 91.945 | 89.636 | 89.982 | 89.326 | 90.097 | 89.687 |
| 88.947 | 90.089 | | 93.934 | | 90.715 | 92.956 | 93.139 | 92.640 | 93.971 | 94.049 | 95.252 | 95.218 | 94.871 | 93.370 | 90.858 | | 88.705 | 89.799 | 90.820 |
| 89.013 | 91.409 | 91.766 | 91.300 | 89.094 | | 91.762 | 91.894 | 91.760 | | 94.340 | 95.136 | 95.198 | 93.910 | 92.631 | 91.266 | 88.808 | | 90.820 | 92.146 |
| | 89.022 | | 92.089 | 90.563 | 91.111 | 90.771 | 91.760 | 94.166 | 95.580 | 96.696 | 96.480 | | 92.550 | 91.434 | 90.386 | 88.222 | 89.478 | | 94.739 |
| 84.980 | 87.456 | 86.653 | | 89.383 | | 91.451 | | 95.185 | 96.720 | 98.135 | 95.753 | 94.207 | 92.803 | 90.619 | 90.670 | 89.349 | 88.425 | | 95.987 |
| 83.870 | | 87.699 | 88.988 | 90.164 | 91.451 | 92.658 | 94.030 | 93.956 | | 99.901 | | 91.765 | 91.611 | 91.594 | 90.431 | | 99.901 | 97.251 | |

Plot of Final Utilities

The values of the states highlighted earlier are the highest. All of these states have a utility of 99.901 and this could be an indication that the agent found the right policy.

6.3. Value Iteration Policy

The final policy calculated by the value iteration algorithm is shown below.



Plot of Optimal Policy

The algorithm finds the optimal policy which is to reach the previously mentioned states. However, the optimal policy can vary because for state (19, 7) ((col, row) format) all four actions will result in the same expected utility. Therefore, the optimal policy can vary.

6.4. Policy Iteration Final Value

The final utilities calculated by the policy iteration algorithm are shown below.

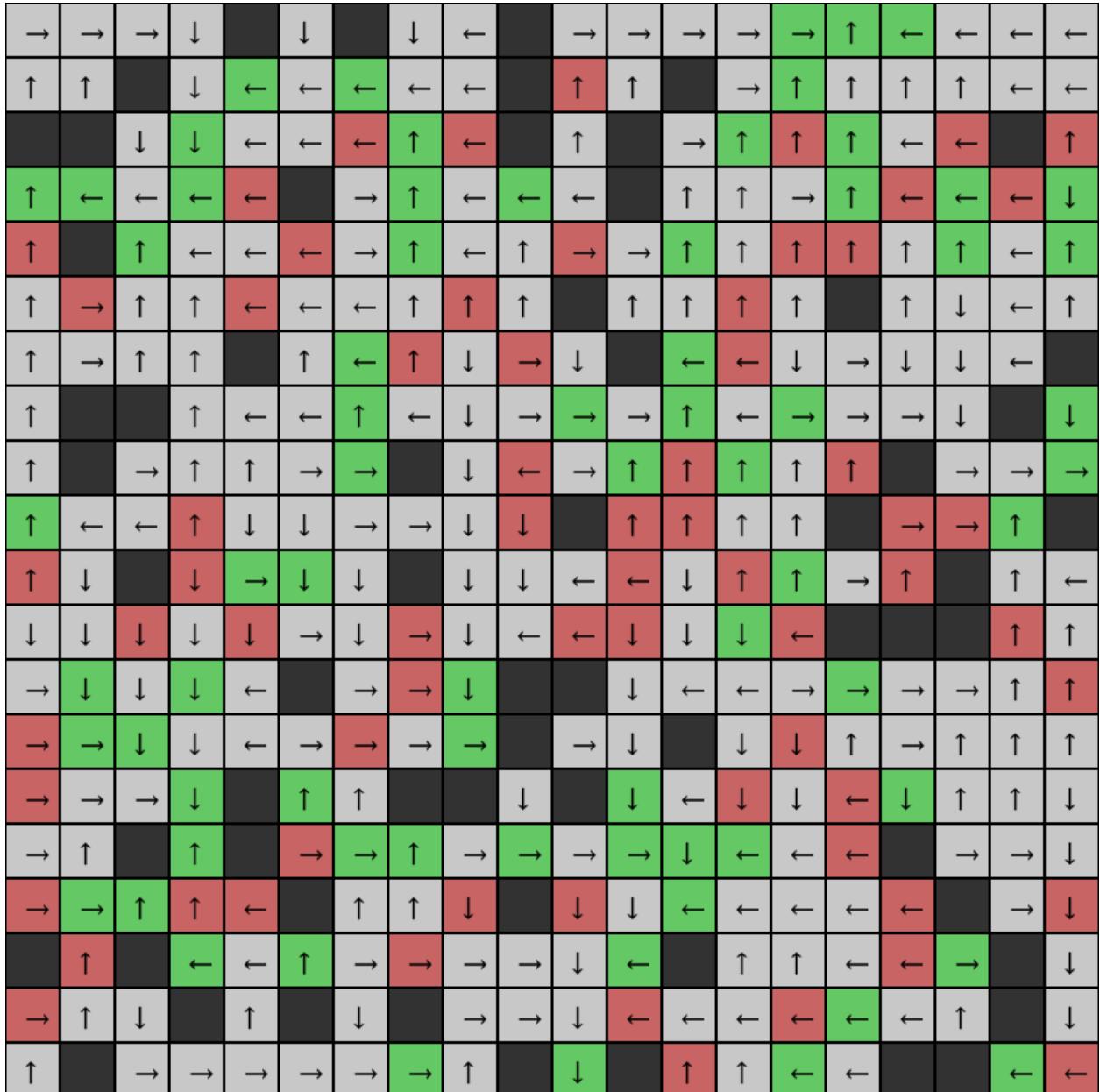
| | | | | | | | | | | | | | | | | | | | |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 91.474 | 92.812 | 94.202 | 95.442 | | 94.230 | | 92.482 | 91.399 | | 92.891 | 94.382 | 95.811 | 97.071 | 98.375 | 98.386 | 98.235 | 96.807 | 95.429 | 94.064 |
| 90.414 | 91.474 | | 96.852 | 96.877 | 95.470 | 95.271 | 93.836 | 92.373 | | 90.750 | 92.891 | | 96.849 | 98.122 | 97.184 | 96.845 | 95.569 | 94.449 | 93.049 |
| | | 97.350 | 98.123 | 96.763 | 95.503 | 93.200 | 93.591 | 91.380 | | 89.569 | | 95.233 | 96.624 | 95.875 | 96.929 | 95.632 | 93.516 | | 90.642 |
| 99.994 | 99.994 | 98.534 | 98.365 | 95.951 | | 92.227 | 93.366 | 92.029 | 91.986 | 90.619 | | 94.130 | 95.217 | 95.072 | 96.496 | 94.107 | 94.086 | 91.755 | 92.819 |
| 97.500 | | 98.392 | 97.098 | 95.618 | 93.077 | 91.862 | 93.117 | 91.685 | 90.829 | 90.299 | 92.662 | 94.018 | 93.875 | 92.880 | 93.835 | 93.077 | 93.886 | 92.479 | 92.861 |
| 96.012 | 94.207 | 96.686 | 95.691 | 93.511 | 92.251 | 91.141 | 91.597 | 89.557 | 89.637 | | 91.577 | 92.555 | 91.586 | 91.662 | | 91.955 | 92.675 | 91.696 | 91.658 |
| 94.679 | 93.989 | 95.199 | 94.531 | | 91.062 | 91.159 | 89.339 | 88.582 | 88.019 | 89.962 | | 93.034 | 90.810 | 91.265 | 91.909 | 93.027 | 94.108 | 92.754 | |
| 93.449 | | | 93.135 | 91.781 | 90.567 | 91.055 | 89.812 | 89.727 | 89.851 | 91.391 | 91.586 | 92.810 | 91.511 | 92.445 | 92.655 | 94.256 | 95.651 | | 99.998 |
| 92.234 | | 90.371 | 91.641 | 90.625 | 89.917 | 91.050 | | 90.884 | 88.838 | 90.389 | 91.456 | 90.615 | 91.474 | 91.177 | 90.354 | | 97.083 | 98.668 | 99.998 |
| 92.167 | 90.843 | 89.739 | 89.321 | 89.461 | 89.696 | 90.141 | 91.137 | 92.338 | 90.531 | | 89.006 | 88.497 | 90.079 | 90.000 | | 93.289 | 96.049 | 98.395 | |
| 89.778 | 89.834 | | 88.712 | 90.630 | 90.853 | 90.775 | | 93.931 | 92.711 | 91.506 | 89.217 | 88.778 | 88.028 | 89.860 | 89.552 | 90.733 | | 96.944 | 95.535 |
| 90.023 | 91.026 | 89.674 | 90.855 | 89.070 | 90.775 | 91.962 | 93.068 | 95.320 | 93.931 | 91.512 | 90.224 | 89.989 | 89.970 | 87.970 | | | 94.469 | 94.313 | |
| 91.084 | 92.520 | 91.752 | 92.424 | 91.021 | | 93.184 | 94.552 | 97.029 | | | 92.495 | 91.148 | 90.091 | 89.066 | 90.291 | 90.303 | 91.603 | 92.957 | 92.007 |
| 90.269 | 92.701 | 92.781 | 92.588 | 91.344 | 92.071 | 93.273 | 95.875 | 97.301 | | 92.661 | 93.881 | | 91.259 | 89.762 | 89.211 | 89.434 | 90.425 | 91.532 | 90.889 |
| 89.220 | 91.502 | 92.725 | 93.939 | | 92.159 | 92.071 | | 92.844 | | 95.270 | 94.173 | 92.648 | 92.040 | 89.732 | 90.079 | 89.424 | 90.195 | 89.783 | |
| 89.028 | 90.170 | | 94.014 | | 90.810 | 93.050 | 93.234 | 92.735 | 94.066 | 94.144 | 95.347 | 95.313 | 94.966 | 93.465 | 90.953 | | 88.802 | 89.896 | 90.916 |
| 89.094 | 91.490 | 91.847 | 91.380 | 89.175 | | 91.856 | 91.989 | 91.855 | | 94.436 | 95.231 | 95.293 | 94.005 | 92.727 | 91.361 | 88.903 | | 90.916 | 92.242 |
| | 89.103 | | 92.170 | 90.644 | 91.199 | 90.866 | 91.855 | 94.261 | 95.675 | 96.791 | 96.575 | | 92.645 | 91.529 | 90.481 | 88.318 | 89.560 | | 94.835 |
| 85.062 | 87.538 | 86.745 | | 89.464 | | 91.547 | | 95.280 | 96.815 | 98.230 | 95.848 | 94.302 | 92.898 | 90.714 | 90.765 | 89.444 | 88.508 | | 96.083 |
| 83.952 | | 87.792 | 89.082 | 90.257 | 91.547 | 92.753 | 94.126 | 94.052 | | 99.996 | | 91.860 | 91.707 | 91.689 | 90.526 | | 99.997 | 97.347 | |

Plot of Final Utilities

The values of the states highlighted earlier are the highest. All of these states have a utility of 99.998 and this could be an indication that the agent found the right policy.

6.5. Policy Iteration Policy

The final policy calculated by the policy iteration algorithm is shown below.



Plot of Optimal Policy

The algorithm finds the optimal policy which is to reach the previously mentioned states.

However, the optimal policy can vary because for state (19, 7) ((col, row) format) all four actions will result in the same expected utility. Therefore, the optimal policy can vary.

6.6. Impact of number of states affecting convergence

As the number of states increases, the time and space required for the 2 algorithms also increases. However, the discounted reward equation always converges as long as the discount factor (GAMMA) is less than 1. If gamma was equal to 1, then the equation would not converge since the limit of gamma tending to infinity would not be 0.

The algorithms will converge but it may take a longer time. Note that the choice of programming language also affects the time taken to run the algorithm.

$$\begin{aligned} U([s_0, s_1, s_2, \dots]) &= R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \\ &= \sum_{t=0}^{\infty} \gamma^t R(s_t) \leq \frac{R_{\max}}{1-\gamma} \quad (0 < \gamma < 1) \end{aligned}$$

Discounted Reward

6.7. Impact of maze complexity affecting policy

Both policy and value iteration algorithms are guaranteed to converge to an optimal policy for discounted finite MDPs. Since increasing the number of states does not change the fact that the MDP is finite, the algorithms will still converge for this question. If the maze was to be made in such a way that it had an infinite state space and infinite action space, then the MDP would not guarantee to converge. Since the utilities calculated by these 2 algorithms will converge, the policy learnt by the agent will also be optimal (greedy action).