

EE183DA  
Team Buffalo  
Lab 3: Markov Decision Processes (MDPs)



THE UNIVERSITY OF CALIFORNIA - LOS ANGELES

Iou-Sheng (Danny) Chang · UID: 804-743-003  
William Argus · UID: 004-610-455  
XianXing (Gray) Jiang · UID: 604-958-018  
Ho (Bobby) Dong · UID: 604-954-176

February 19th 2019

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction and Lab Overview</b>              | <b>3</b>  |
| <b>2</b> | <b>MDP system</b>                                 | <b>5</b>  |
| 2.1      | State space . . . . .                             | 5         |
| 2.2      | Action space . . . . .                            | 5         |
| 2.3      | Probability function . . . . .                    | 5         |
| 2.4      | Next state function . . . . .                     | 5         |
| 2.5      | MATLAB Implementation . . . . .                   | 6         |
| <b>3</b> | <b>Planning Problem</b>                           | <b>10</b> |
| 3.1      | Reward function . . . . .                         | 10        |
| 3.2      | MATLAB Implementation . . . . .                   | 10        |
| <b>4</b> | <b>Policy Iteration</b>                           | <b>11</b> |
| 4.1      | Overview . . . . .                                | 11        |
| 4.2      | Initial policy matrix . . . . .                   | 11        |
| 4.3      | Robot trajectory plotting function . . . . .      | 12        |
| 4.4      | Trajectory plot using policy $\pi_0$ . . . . .    | 12        |
| 4.5      | Policy evaluation computation . . . . .           | 12        |
| 4.6      | Value of the trajectory in section 4.3 . . . . .  | 12        |
| 4.7      | One-step lookahead optimal policy $\pi$ . . . . . | 13        |
| 4.8      | Policy Iteration . . . . .                        | 13        |
| 4.9      | Optimal policy trajectory plot . . . . .          | 13        |
| 4.10     | Computation time . . . . .                        | 13        |
| 4.11     | MATLAB Implementation . . . . .                   | 14        |
| <b>5</b> | <b>Value Iteration</b>                            | <b>27</b> |
| 5.1      | Overview . . . . .                                | 27        |
| 5.2      | Value Iteration . . . . .                         | 27        |
| 5.3      | Optimal policy trajectory plot . . . . .          | 28        |
| 5.4      | Computation time . . . . .                        | 28        |
| 5.5      | MATLAB Implementation . . . . .                   | 28        |
| <b>6</b> | <b>Additional Scenarios</b>                       | <b>31</b> |
| 6.1      | Trajectory with error . . . . .                   | 31        |
| 6.2      | Trajectory with altered reward . . . . .          | 31        |
| 6.3      | Conclusions . . . . .                             | 31        |
| 6.4      | MATLAB Implementation . . . . .                   | 32        |
| <b>7</b> | <b>Appendix and Demonstration</b>                 | <b>34</b> |
| <b>8</b> | <b>References</b>                                 | <b>35</b> |

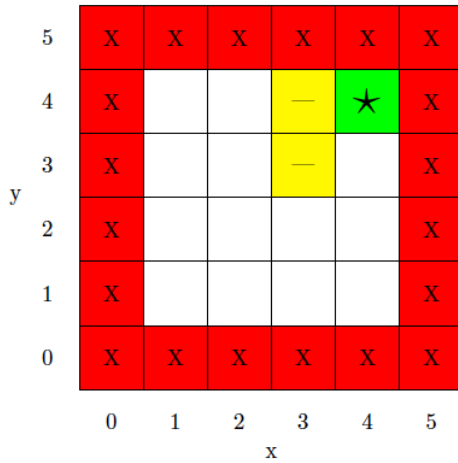
# 1 Introduction and Lab Overview

Markov decision process (MDP) is a model widely used for machine learning. MDP allows the robot to determine the optimal action to take for the given state it is currently in. In order to achieve this, the general MDP model has to be given four elements: state  $s$ , action  $a$ , transition function  $T(s, a, s')$ , and reward  $R(s, a, s')$ . The transition function and the reward are dependent on the current state, the action being taken, and the next state. The solution of the MDP model is to obtain the optimal policy that prescribes the action the robot should take at a given state,  $s$ . This is achieved by finding a policy that prescribes the actions that maximizes the cumulative reward in the following equation:

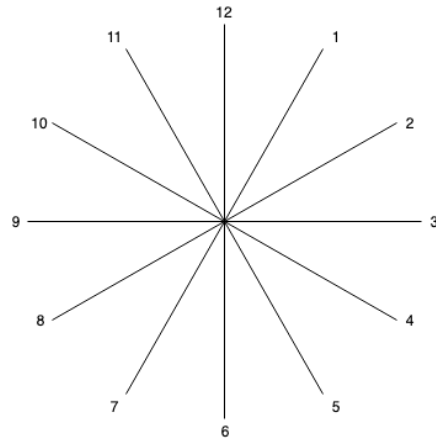
$$\sum_{t=0}^{\infty} \gamma^t R_{a_t}(s_t, s_{t+1}) \quad (1)$$

Where  $a_t = \pi(s_t)$ , i.e. actions given by policy. Note that  $\gamma$  is the discount factor.

The state is the positional configuration plus the heading of the robot. In an MDP, all possible states are known, and the robot will always exist in exactly one of these states at any given time. For this lab, the robot runs in a grid world with finite cells, and its current state is represented by its  $x$  and  $y$  coordinates (each ranging from 1 to 6 for the 6 by 6 grid, shown in fig. 1a), and its heading (ranging from 1 to 12, laid out like the numbers on a clock, shown in fig. 1b).



(a) Grid World given in Lab 3 instruction



(b) 12 possible headings of the robot

The actions are a set of movements that the MPD model can dictate the robot to take in order to transfer it from the current state to the next state. For this lab, the robot is only allowed to take 7 actions: Forward, and turn left, no turn, or turn right (current heading -1, 0, and +1, respectively), backward, and turn left, no turn, or turn right (current heading -1, 0, and +1, respectively), and no forward or backward movement and no rotation.

The transition matrix describes the transition from state  $s$  to the state  $s'$  by taking action  $a$ . For an ideal system, the next state,  $s'$ , is always the result of being at the current state and taking the given action. For a system that contains noise, the next state is determined by being at the state,  $s$ , and taking the given action with a probabilistic model. For this lab, the transition function follows a probabilistic model with error called “pre-rotation”. “Pre-rotation” occurs before the robot performs the linear motion (forward/backward) part of a prescribed action. It will change the heading of the robot, resulting in the robot taking action while in a different state than the current state input. This “pre-rotation” occurs with the probability  $pe$  for both pre-rotating right and left; the probability of the robot not “pre-rotating” and taking the given action at the inputted state is  $1 - 2pe$ .

The reward describes the reward that the robot will get when it moves into state,  $s'$ . The reward in the MDP model serves as an indicator to the robot of which state is preferred and which states are to be avoided. For this lab, the reward is set to specific values, defined only by the state, indicated from the lab manual (shown in fig. 1a ), with white being 0, red being  $-100$ , yellow being  $-10$ , and green being  $+1$ , the goal square.

## 2 MDP system

### 2.1 State space

The state space is  $N_S \in \mathbb{R}^{6 \times 6 \times 12}$ . It consists of all possible states and is represented by a 6 by 6 by 12 three dimensional matrix. The first two dimensions give the x and y values, each with possibilities 1 through 6. The third dimension gives the heading of the state. The headings are numbered 1 through 12 and are arranged the same way as numbers on a clock. For example, moving forward with right rotation while in state  $x = 1$ ,  $y = 1$ , and heading  $= 12$ , results in moving to state  $x = 1$ ,  $y = 2$ , and heading  $= 1$ .

The state space  $N_S$  is implemented in MATLAB as *stateSpace* in the function *init.m* (listing 1).

### 2.2 Action space

The action space is  $N_A \in \mathbb{R}^{7 \times 2}$ . It consists of 7 actions. The first three are move forward, with left, no, right rotation respectively. The next three are move backward, with left, no, right rotation respectively. The last action is no movement and no rotation.

The state space  $N_S$  is implemented in MATLAB as *action* in the function *init.m* (listing 1).

### 2.3 Probability function

$p_{sa}(s')$  is the function that returns the probability, given  $p_e, s, a, s'$ . It first determines if  $s'$  is a possible next state given the current state ( $s$ ), the action ( $a$ ), and error probability ( $p_e$ ). There will only be at most 3 possible next states for  $s'$ . This is due to the fact that the action dictate whether the robot will move forward or backward and its rotation, so the only thing that is not determined is the pre-rotation, which has 3 possibilities; left with probability  $p_e$ , right with probability  $p_e$ , or no rotation with probability  $1 - 2p_e$ . Accordingly, the function calculates the probability that each of these pre-rotations happen, meaning that the probability of all possible states is known given the current state, action, and  $p_e$ . The probability of the requested state  $s'$  can then be returned by the function, as required.

The probability function is implemented in MATLAB as *probability.m* (listing 2).

### 2.4 Next state function

The function that returns the next state  $s'$  given the current state  $s$ , action  $a$ , and error probability  $p_e$  uses the probability function defined in *probability.m* (listing 2). Since the probability of each of the possible next states is known, the MATLAB code generates a random number over a sufficient large range of values, (the size of which is given by the variable *pe\_m*, calculated by 10 to the power of the number of decimal places in  $p_e$  ( $10^{pe\_power}$ )). This way

it is assured that  $pe\_m$  is a sufficiently large number. We then assign the first  $pe \cdot 10^{pe\_power}$  values to signify the state given in the event of a pre-rotation to right, the next  $pe \cdot 10^{pe\_power}$  values to signify the state given in the event of a pre-rotation to the left, and the remaining  $10^{pe\_power} - 2 \cdot pe \cdot 10^{pe\_power}$  values to signify the state given in the event of no pre-rotation. Using this method, it can be guaranteed that the right, left, and no pre-rotations happen for each movement of the robot with the probability given by any  $p_e$  values.

The next state function is implemented in MATLAB as *next\_state.m* (listing 3).

## 2.5 Matlab Implementation

Listing 1: MATLAB code for Initialization function

```

1 function [stateSpace,pe,lambda,s,action] = init()
2 % Initialization of value
3 value_int= zeros(6);
4
5 % Expend the value matrix to 12 headings' 3D matrix space
6 [rows,cols] = size(value_int);
7 for x = 1:rows
8     for y = 1:cols
9         for h = 1:12
10             stateSpace(x,y,h) = value_int(x,y);
11         end
12     end
13 end
14
15 % Initialization of error percentage (pe)
16 pe = 0.1;
17
18 % Initialization of discount factor
19 lambda = 0.9;
20
21 % Initialization of initial state
22 s = [2;5;6];
23
24 % Initialization of actions set
25 % column 1: forward 1 / no movement 0 / backward -1
26 % column 2: leftturn -1 / no rotation 0 / rightturn 1
27 action = [1  -1;...
28           1   0;...
29           1   1;...
30          -1 -1;...
31          -1  0;...
32          -1  1;...
33           0  0];
34 end

```

Listing 2: MATLAB code for Probability function

```

1 function [prob,state_p] = probability(pe,state_c,action,state_n)
2     % ===== Heading ===== %
3     % circular loop to account for the flaw in mod 12
4     heading = [12 ,1 ,2, 3, 4, 5, 6, 7, 8, 9, 10, 11 ,12, 1];
5
6     % ===== Possible States ===== %
7     % state_p is a R2
8     % cause there are three possible prerotation outcome
9     state_p = [state_c state_c state_c];
10    if (pe ~= 0 && action(1) ~=0)
11        % prerotation -1 (left)
12        state_p(3,1) = heading(state_p(3,1));
13        % prerotation 0 (no prerotation)
14        state_p(3,2) = heading(state_p(3,2)+1);
15        % prerotation 1 (right)
16        state_p(3,3) = heading(state_p(3,3)+2);
17    end
18
19    % ===== Update Possible Next State ===== %
20    % State depedning on heading
21    for i =1:3
22        % (x,y)
23        % If statements account for when robot is at the border
24        % when @ the border, robot can still rotate but cannot
        % perform movement
25        switch state_p(3,i)
26            case {11,12,1}
27                state_p(2,i) = state_p(2,i) + action(1);
28                if (state_p(2,i) == 7 || state_p(2,i) == 0)
29                    state_p(2,i) = state_c(2);
30                end
31            case {2,3,4}
32                state_p(1,i) = state_p(1,i) + action(1);
33                if (state_p(1,i) == 7 || state_p(1,i) == 0)
34                    state_p(1,i) = state_c(1);
35                end
36            case {5,6,7}
37                state_p(2,i) = state_p(2,i) - action(1);
38                if (state_p(2,i) == 0 || state_p(2,i) == 7)
39                    state_p(2,i) = state_c(2);
40                end
41            case {8,9,10}
42                state_p(1,i) = state_p(1,i) - action(1);
43                if (state_p(1,i) == 0 || state_p(1,i) == 7)
44                    state_p(1,i) = state_c(1);
45                end

```

```

46     end
47     % Heading
48     state_p(3,i) = heading((state_p(3,i)+1) + action(2));
49
50 end
51 % Probability
52 % Attach the probability corresponding to each possible states
53 % For prerotated (-1/1) states: pe
54 % For non-prerotated (0) states: 1-2*pe
55 state_p(4,:) = [pe, 1-2*pe, pe];
56
57 % ===== Finding probability for the selected state ===== %
58 prob = 0;
59 for i = 1:3
60     % Check if the possible next state == next state
61     % If so save the probability attached to that state
62     if state_p(1:3,i) == state_n
63         prob = state_p(4,i);
64         break;
65     end
66 end

```

Listing 3: MATLAB code for Next State function

```

1 function [state_n,prerotation] = next_state(pe,state_c,action)
2     % ===== pe && prerotation ===== %
3     if (pe == 0)
4         prerotation = 0;
5     else
6         % initialization of pe array
7         pe_power = 0;
8         pe_m = [];
9         % calculate array size needed to account for pe's # of
          digits after decimal point
10        while (floor(pe*10^pe_power) ~= pe*10^pe_power)
11            pe_power = pe_power + 1;
12        end
13        % create an array to store movements based on pe
14        for i = 1:pe*power(10,pe_power)
15            pe_m = [pe_m 1 -1];
16        end
17        pe_m = [pe_m zeros(1,power(10,pe_power) - 2*pe*power(10,
          pe_power))];
18        % prerotation based on pe (defined in init.m)
19        prerotation = pe_m(randi(power(10,pe_power)));
20    end
21
22    % ===== Initialize next state ===== %

```



```

23     state_n = state_c;
24
25     % ===== Update Next State ===== %
26     % circular loop to account for the flaw in mod 12
27     heading = [12 ,1 ,2, 3, 4, 5, 6, 7, 8, 9, 10, 11 ,12, 1];
28     % Heading (prerotation)
29     % account for pe = 0 (no error percentage), will not incur an
        error rotation
30     % account for choosing to stay still, will not incur an error
        rotation
31     if (pe ~= 0 && action(1) ~=0)
32         state_n(3) = heading((state_n(3)+1) + prerotation);
33     end
34     % (x,y)
35     % If statements account for when robot is at the border
36     % when @ the border, robot can still rotate but cannot perform
        movement
37     switch state_n(3)
38         case {11,12,1}
39             state_n(2) = state_c(2) + action(1);
40             if (state_n(2) == 7 || state_n(2) == 0)
41                 state_n(2) = state_c(2);
42             end
43         case {2,3,4}
44             state_n(1) = state_c(1) + action(1);
45             if (state_n(1) == 7 || state_n(1) == 0)
46                 state_n(1) = state_c(1);
47             end
48         case {5,6,7}
49             state_n(2) = state_c(2) - action(1);
50             if (state_n(2) == 0 || state_n(2) == 7)
51                 state_n(2) = state_c(2);
52             end
53         case {8,9,10}
54             state_n(1) = state_c(1) - action(1);
55             if (state_n(1) == 0 || state_n(1) == 7)
56                 state_n(1) = state_c(1);
57             end
58     end
59     % Heading
60     state_n(3) = heading((state_n(3)+1) + action(2));
61 end

```

## 3 Planning Problem

### 3.1 Reward function

For the reward function  $R(s)$ , a matrix was created that contained the value of the reward at that state in the appropriate state. For example, when cell (5,5) was referenced, it gave the reward of the goal state (5,5) of +1. When a state  $s$  was inputted, the function simply returned the value of the cell of matrix given by the  $x$  and  $y$  values in  $s$ . This simple implementation was possible because reward in a state is independent of heading, so the heading in state  $s$  could simply be ignored, with only the  $x$  and  $y$  values taken into account. In order to check if the reward function returns correct reward for any given state, we used boolean statement to test for all cases in MATLAB *main.m*.

The next state function is implemented in MATLAB as *reward.m* (listing 4).

### 3.2 Matlab Implementation

Listing 4: MATLAB code for Reward function

```

1 function reward_disp = reward(state_c)
2     % ===== Grid World Reward ===== %
3     % REFER to the grid world shown in the lab instruction
4     gw_reward = [-100 -100 -100 -100 -100 -100;...
5                  -100 0      0      -10  1      -100;...
6                  -100 0      0      -10  0      -100;...
7                  -100 0      0      0      0      -100;...
8                  -100 0      0      0      0      -100;...
9                  -100 -100 -100 -100 -100 -100];
10    [rows,cols] = size(gw_reward);
11
12    % ===== Map the reward value ===== %
13    reward_disp = gw_reward((rows+1) - state_c(2),state_c(1));
14 end

```

## 4 Policy Iteration

### 4.1 Overview

Policy iteration changes the policy of the robot's actions after each movement and then recompute the value of the state. After repeated iterations, the policy will converge to a single policy, and this policy is guaranteed to be the optimal policy. Usually, the optimal policy is converged upon before the optimal value function is converged upon.

The optimal policy is found using the Bellman equation eq. (3) to perform policy iteration, which improves  $\pi$  until it converges on the optimal policy function. The equations for optimal value and optimal policy are eq. (4) and eq. (5), respectively.

$$V^\pi(s) = E \left[ \sum_{i=1}^T \gamma^{i-1} r_i \right] \quad \forall s \in \mathbb{S} \quad (2)$$

$$V^*(S) = \max_a \left[ R(s, a) + \gamma \sum_{s' \in \mathbb{S}} p(s'|s, a) V^*(s') \right] \quad (3)$$

$$V^*(s) = \max_{\pi} V^\pi(s) \quad \forall s \in \mathbb{S} \quad (4)$$

$$\pi^* = \underset{\pi}{\operatorname{argmax}} V^\pi(s) \quad \forall s \in \mathbb{S} \quad (5)$$

The optimal Q-function (eq. (6)) will give the total expected reward gained by the robot given that the robot starts in state  $s$ , takes action  $a$ , and then continues to take the best action for the rest of its movement.

$$V^*(s) = \max_a Q^*(s, a) \quad \forall s \in \mathbb{S} \quad (6)$$

### 4.2 Initial policy matrix

To create a matrix that stored the action  $a = \pi_0(s)$ , the action given by the initial policy, a function was written to assign the appropriate action to each state. This function took into account position (x,y) and heading and assigned the action prescribed the initial policy of moving forward or backward, whichever took the robot closest to the goal square, or, if the goal square was to the right of left of the robot, moving forward and turning toward the goal square.

The action  $a = \pi_0(s)$  is implemented in MATLAB as *init\_policy.m* (listing 5).

### 4.3 Robot trajectory plotting function

The function that plotted the trajectory given policy matrix  $\pi$ , initial state  $s_0$  and probability of error  $p_e$ . We first implemented the Grid World shown in Lab 3 instructions (fig. 1a) into MATLAB figure, then live update robot trajectory onto that existing figure, with the heading indicated. The action and prerotation for next state are also shown in the Grid World for the ease to verify and debug the code, a sample figure when the robot reaches the goal square is shown in (fig. 2).

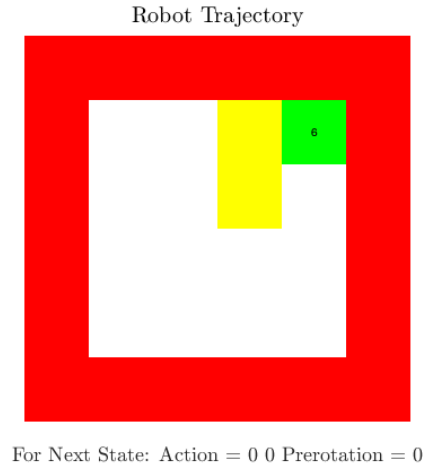


Figure 2: Grid World implemented on MATLAB and robot trajectory plot

The robot trajectory plotting function is implemented in MATLAB as *plot\_route.m* (listing 6).

### 4.4 Trajectory plot using policy $\pi_0$

The plot of the trajectory using policy  $\pi_0$  from initial state  $x = 2, y = 5, h = 6$ , assuming  $p_e = 0$  is plotted using (listing 5) and (listing 6).

The trajectory plotted is shown in video in the Appendix and Demonstration section.

### 4.5 Policy evaluation computation

This function returns a matrix of value  $v = V^\pi(s)$  when indexed by state  $s$ .

The function is implemented in MATLAB as *get\_value.m* (listing 7).

### 4.6 Value of the trajectory in section 4.3

The value matrix is shown in *val* in MATLAB *policy\_init.m* (listing 8).

## 4.7 One-step lookahead optimal policy $\pi$

This function returns a matrix  $\pi$  given a one-step lookahead on value

The function is implemented in MATLAB as *get\_policy.m* (listing 9).

## 4.8 Policy Iteration

Combining the *get\_value.m* and *get\_policy.m* functions, the policy iteration *policy.m* code returns optimal policy  $\pi^*$  and optimal value  $V$ .

### Pseudo Code for Policy Iteration

```

Start with an arbitrary policy  $\pi_0$ 
Rerun
     $\pi$  corresponds to  $\pi_0$ 
    Compute values using  $\pi$  by solving the following equation:
        
$$V^\pi(S) = E[r|s, \pi(s)] + \gamma \sum_{s' \in S} P(s'|s, \pi(s)) V^\pi(s')$$

    At each state, improve policy
        
$$\pi'(s) = \operatorname{argmax}_a (E[r|s, a] + \gamma \sum_{s' \in S} P(s'|s, a) V^\pi(s'))$$

Continue until  $\pi = \pi'$ 

```

Figure 3: Pseudo code for Policy Iteration

The policy iteration function is implemented in MATLAB as *policy.m* (listing 10).

## 4.9 Optimal policy trajectory plot

The plot of the trajectory under optimal policy  $\pi^*$  and  $p_e = 0$ .

The trajectory plotted is shown in video in the Appendix and Demonstration section.

## 4.10 Computation time

We use the MATLAB built in *tic – toc* function to determine the compute time of Policy Iteration, as shown in *main.m* (listing 11), the compute time is then printed to the commend window as shown in fig. 4, it is determined to be 4.6858 [s].

```

===== 2.3(a) =====
REFER to 'init_policy.m' for detail
===== 2.3(b) =====
REFER to 'plot_route.m' for detail
===== 2.3(c) =====
REFER to 'policy_init.m' for detail
REFER to figure plotted or video linked in the lab report reference
===== 2.3(d) =====
REFER to 'get_value.m' for detail
===== 2.3(e) =====
REFER to 'policy_init.m' for detail
===== 2.3(f) =====
REFER to 'get_policy.m' for detail
===== 2.3(g) =====
REFER to 'policy.m' for detail
===== 2.3(h) =====
REFER to figure plotted or video linked in the lab report reference
===== 2.3(i) =====
Policy Iteration Compute Time = 4.6858 [s]>>

```

Figure 4: Compute time of Policy Iteration

## 4.11 Matlab Implementation

Listing 5: MATLAB code for  $a = \pi_0(s)$ 

```

1 function [policy] = init_policy(rows)
2     % ===== Initialization of initial policy \pi_0 ===== %
3     % REFER to lab 3 instructions for the initialization process
4     % If the goal is in front of you, move forward
5     % If it is behind you, move backward
6     % Then turn the amount that aligns your next direction of
7     % travel closer towards the goal (if necessary).
8     % If the goal is directly to your left or right, move
9     % forward then turn appropriately.
10    for x = 1:6
11        for y = 1:6
12            % Iterate through all headings
13            for h = 1:12
14                % Case 1: (x,y) = (<5,<5)
15                if (x < 5 && (rows+1) - y < 5)
16                    switch h
17                        case {11,12,1}
18                            policy(y,x,h,1) = 1;
19                            policy(y,x,h,2) = 1;
20                        case {2,3,4}
21                            policy(y,x,h,1) = 1;
22                            policy(y,x,h,2) = -1;
23                        case {5,6,7}
24                            policy(y,x,h,1) = -1;
25                            policy(y,x,h,2) = -1;
26                        case {8,9,10}
27                            policy(y,x,h,1) = -1;
28                            policy(y,x,h,2) = 1;
29                    end
30                end
31            end
32        end
33    end

```

```

29 % Case 2: (x,y) = (<5,5)
30 if (x < 5 && (rows+1) - y == 5)
31     switch h
32     case {11,12,1}
33         policy(y,x,h,1) = 1;
34         policy(y,x,h,2) = 1;
35     case {2,3,4}
36         policy(y,x,h,1) = 1;
37         policy(y,x,h,2) = 0;
38     case {5,6,7}
39         policy(y,x,h,1) = 1;
40         policy(y,x,h,2) = -1;
41     case {8,9,10}
42         policy(y,x,h,1) = -1;
43         policy(y,x,h,2) = 0;
44     end
45 end
46 % Case 3: (x,y) = (<5,6)
47 if (x < 5 && (rows+1) - y == 6)
48     switch h
49     case {11,12,1}
50         policy(y,x,h,1) = -1;
51         policy(y,x,h,2) = 1;
52     case {2,3,4}
53         policy(y,x,h,1) = 1;
54         policy(y,x,h,2) = 1;
55     case {5,6,7}
56         policy(y,x,h,1) = 1;
57         policy(y,x,h,2) = -1;
58     case {8,9,10}
59         policy(y,x,h,1) = -1;
60         policy(y,x,h,2) = -1;
61     end
62 end
63 % Case 4: (x,y) = (5,<5)
64 if (x == 5 && (rows+1) - y < 5)
65     switch h
66     case {11,12,1}
67         policy(y,x,h,1) = 1;
68         policy(y,x,h,2) = 0;
69     case {2,3,4}
70         policy(y,x,h,1) = 1;
71         policy(y,x,h,2) = -1;
72     case {5,6,7}
73         policy(y,x,h,1) = -1;
74         policy(y,x,h,2) = 0;
75     case {8,9,10}

```

```

76         policy(y,x,h,1) = 1;
77         policy(y,x,h,2) = 1;
78     end
79 end
80 % Case 5: (x,y) = (5,5) REACHED GOAL SQUARE
81 if (x == 5 && (rows+1) - y == 5)
82     policy(y,x,h,1) = 0;
83     policy(y,x,h,2) = 0;
84 end
85 % Case 6: (x,y) = (5,6)
86 if (x == 5 && (rows+1) - y == 6)
87     switch h
88     case {11,12,1}
89         policy(y,x,h,1) = -1;
90         policy(y,x,h,2) = 0;
91     case {2,3,4}
92         policy(y,x,h,1) = 1;
93         policy(y,x,h,2) = 1;
94     case {5,6,7}
95         policy(y,x,h,1) = 1;
96         policy(y,x,h,2) = 0;
97     case {8,9,10}
98         policy(y,x,h,1) = 1;
99         policy(y,x,h,2) = -1;
100    end
101 end
102 % Case 7: (x,y) = (6,<5)
103 if (x == 6 && (rows+1) - y < 5)
104     switch h
105     case {11,12,1}
106         policy(y,x,h,1) = 1;
107         policy(y,x,h,2) = -1;
108     case {2,3,4}
109         policy(y,x,h,1) = -1;
110         policy(y,x,h,2) = -1;
111     case {5,6,7}
112         policy(y,x,h,1) = -1;
113         policy(y,x,h,2) = 1;
114     case {8,9,10}
115         policy(y,x,h,1) = 1;
116         policy(y,x,h,2) = 1;
117     end
118 end
119 % Case 8: (x,y) = (6,5)
120 if (x == 6 && (rows+1) - y == 5)
121     switch h
122     case {11,12,1}

```



```

123         policy(y,x,h,1) = 1;
124         policy(y,x,h,2) = -1;
125         case {2,3,4}
126             policy(y,x,h,1) = -1;
127             policy(y,x,h,2) = 0;
128         case {5,6,7}
129             policy(y,x,h,1) = 1;
130             policy(y,x,h,2) = 1;
131         case {8,9,10}
132             policy(y,x,h,1) = 1;
133             policy(y,x,h,2) = 0;
134         end
135     end
136     % Case 9: (x,y) = (6,6)
137     if (x == 6 && (rows+1) - y ==6)
138         switch h
139             case {11,12,1}
140                 policy(y,x,h,1) = -1;
141                 policy(y,x,h,2) = -1;
142             case {2,3,4}
143                 policy(y,x,h,1) = -1;
144                 policy(y,x,h,2) = 1;
145             case {5,6,7}
146                 policy(y,x,h,1) = 1;
147                 policy(y,x,h,2) = 1;
148             case {8,9,10}
149                 policy(y,x,h,1) = 1;
150                 policy(y,x,h,2) = -1;
151         end
152     end
153 end
154 end
155 end
156 end

```

Listing 6: MATLAB code for robot trajectory plotting

```

1 function plot_route(robot_state,prerotation,action)
2 % Using [RGB] to plot the corresponding colored grid world
3 rgb_1 = [1 1 1 1 1 1;...
4         1 1 1 1 0 1;...
5         1 1 1 1 1 1;...
6         1 1 1 1 1 1;...
7         1 1 1 1 1 1;...
8         1 1 1 1 1 1];
9 rgb_2 = [0 0 0 0 0 0;...
10         0 1 1 1 1 0;...
11         0 1 1 1 1 0;...

```

```

12         0 1 1 1 1 0;...
13         0 1 1 1 1 0;...
14         0 0 0 0 0 0];
15 rgb_3 = [0 0 0 0 0 0;...
16         0 1 1 0 0 0;...
17         0 1 1 0 1 0;...
18         0 1 1 1 1 0;...
19         0 1 1 1 1 0;...
20         0 0 0 0 0 0];
21 rgb = cat(3,rgb_1,rgb_2,rgb_3);
22 [rows,cols,deps] = size(rgb);
23 GW = figure(1);
24 gridWorld = imshow(rgb,'InitialMagnification','fit');
25 set(gca,'ticklabelinterpreter','latex');
26 title('Robot Trajectory','fontsize',20,'interpreter','latex');
27
28 % Plot the robot state on grid world
29 % Number shown on the grid is the heading
30 for gw_x = 1:rows
31     for gw_y = 1:cols
32         if (robot_state(gw_x,gw_y) ~= 0)
33             text(gw_y,gw_x,num2str(robot_state(gw_x,gw_y)),...
34                 'HorizontalAlignment','center','VerticalAlignment','
35                     middle');
36             xlabelStr = ['For Next State: Action = ',num2str(action)
37                 , ' Prerotation = ',num2str(prerotation)];
38             xlabel(xlabelStr,'fontsize',18,'interpreter','latex');
39             drawnow;      refreshdata;
40         end
41     end
42 end
43
44 % 'tightfig.m' is an open source function
45 % Trim off the extra region in the figure
46 tightfig;
47 end

```

Listing 7: MATLAB code for computing policy evaluation

```

1 function [val] = get_value(command,lambda,pe,rows)
2     val = init();
3     value_difference = 10e5;
4     % Policy Evaluation
5     while value_difference > 10e-5
6         % Store previous value to temp
7         temp = val;
8         % Compute and update new value and store
9         for x = 1:6

```

```

10         for y = 1:6
11             for h = 1:12
12                 % Find all possible movements due to prerotation
13                 % Note: here the next state and doesn't matter
14                 [garbage,s] = probability(pe,[x;y;h],command((
15                     rows+1)-y,x,h,:),[1,1,1]'));
16                 % Update the value calculated for policy before
17                 % update
18                 val((rows+1)-y,x,h) = s(4,1) * (reward([s(1,1);s
19                     (2,1)]) + lambda*val((rows+1)-s(2,1),s(1,1),s
20                     (3,1))) + ...
21                     s(4,2) * (reward([s(1,2);s
22                     (2,2)]) + lambda*val((
23                     rows+1)-s(2,2),s(1,2),s
24                     (3,2))) + ...
25                     s(4,3) * (reward([s(1,3);s
26                     (2,3)]) + lambda*val((
27                     rows+1)-s(2,3),s(1,3),s
28                     (3,3))));
29             end
30         end
31     end
32     % Compare new value with previous value
33     value_difference = sum(sum(sum(abs(temp-val))));
34 end
35 end

```

Listing 8: MATLAB code for running the  $\pi_0$  case

```

1 %% EE 183DA Lab 3
2 % Team Buffalo
3     % Iou-Sheng (Danny) Chang    UID: 804-743-003
4     % William                    Argus    UID: 004-610-455
5     % XianXing (Gray) Jiang      UID: 604-958-018
6     % Ho                        (Bobby) Dong    UID: 604-954-176
7 % Markov Decision Processes (MDPs)
8 clc; clear all; close all;
9
10 %% For initial policy \pi_0
11 robot = zeros(6);
12 [rows,cols] = size(robot);
13
14 % Initialization
15 [optimal_value_policy,pe,lambda,s,action] = init();
16 policy_p = init_policy(rows);
17 prer = 0;
18 pe = 0;
19

```

```

20 val = get_value(policy_p,lambda,pe,rows);
21
22 % Robot trajectory plot using initial policy
23 % NOTE: comment out rest of the code in section 2.3 before running
    this part
24 % Initialization
25
26 robot((rows+1)-s(2),s(1)) = s(3);
27 for j = 1:1000
28     act = policy_p((rows+1)-s(2),s(1),s(3),:);
29     plot_route(robot,prer,act);
30     optimal_policy_policy(j,1) = act(1,1,1,1);
31     optimal_policy_policy(j,2) = act(1,1,1,2);
32     if (s(1) == 5 && s(2) ==5)
33         break;
34     end
35     robot((rows+1)-s(2),s(1)) = 0;
36     [s prer]= next_state(pe,s,act);
37     robot((rows+1)-s(2),s(1)) = s(3);
38 end
39
40 % Command Window
41 % 2.3(c)
42 fprintf('===== 2.3(c) =====\n');
43 fprintf('REFER to figure plotted or video linked in the lab report
    reference\n')
44
45 % 2.3(e)
46 fprintf('===== 2.3(e) =====\n');
47 fprintf('REFER to ''val'' matrix\n');

```

Listing 9: MATLAB code for optimal policy

```

1 function [command] = get_policy(command,val,lambda,pe,rows,action)
2
3 for x = 1:6
4     for y = 1:6
5         for h = 1:12
6             set = [];
7             for j = 1:6
8                 % Find all possible movements due to prerotation
9                 % Note: here the next state and doesn't matter
10                [garbage,s] = probability(pe,[x;y;h],action(j,:),
                    ,[1,1,1]');
11                % Store the value calculated for each actions
12                % Stored in a temp array
13                set(j) = s(4,1) * (reward([s(1,1);s(2,1)]) + lambda*val
                    ((rows+1)-s(2,1),s(1,1),s(3,1))) +...

```

```

14         s(4,2) * (reward([s(1,2);s(2,2)]) + lambda*val
15             ((rows+1)-s(2,2),s(1,2),s(3,2))) +...
16         s(4,3) * (reward([s(1,3);s(2,3)]) + lambda*val
17             ((rows+1)-s(2,3),s(1,3),s(3,3)));
18     end
19     % Special case
20     % When we reach the goal (green) square
21     % Robot can perform action (0,0) --> stay still (no
22     % movement, no rotation)
23     if (x == 5 && y == 5)
24         [garbage,s] = probability(pe,[x;y;h],action(7,:),
25             [1,1,1]');
26         set(7) = s(4,1) * (reward([s(1,1);s(2,1)]) + lambda*
27             val((rows+1)-s(2,1),s(1,1),s(3,1))) +...
28             s(4,2) * (reward([s(1,2);s(2,2)]) + lambda*
29             val((rows+1)-s(2,2),s(1,2),s(3,2))) +...
30             s(4,3) * (reward([s(1,3);s(2,3)]) + lambda*
31             val((rows+1)-s(2,3),s(1,3),s(3,3)));
32     end
33     % Find the best action (maximum value) and store it
34     [M,I] = max(set);
35     % Find value of original action
36     [garbage,orig] = probability(pe,[x;y;h],command((rows+1)
37         -y,x,h,:),[1,1,1]');
38     temp = orig(4,1) * (reward([orig(1,1);orig(2,1)]) +
39         lambda*val((rows+1)-orig(2,1),orig(1,1),orig(3,1))) +
40         ...
41         orig(4,2) * (reward([orig(1,2);orig(2,2)]) +
42         lambda*val((rows+1)-orig(2,2),orig(1,2),orig
43         (3,2))) +...
44         orig(4,3) * (reward([orig(1,3);orig(2,3)]) +
45         lambda*val((rows+1)-orig(2,3),orig(1,3),orig
46         (3,3)));
47     % Store new policy if find better
48     if temp < M
49         % Store movement (forward/backward) into command
50         command((rows+1)-y,x,h,1) = action(I,1);
51         % Store rotation (left/none/right) into command
52         command((rows+1)-y,x,h,2) = action(I,2);
53     end
54 end
55 end
56 end
57 end

```

Listing 10: MATLAB code for Policy Iteration

```

1 function [command,val] = policy(policy_p,lambda,pe,rows,action)
2     % Initialization of initial policy \pi_0

```

```

3     command = policy_p;
4     exit = 0;
5
6     while exit == 0
7         command_copy = command;
8         % Initialization
9         val = get_value(command,lambda,pe,rows);
10        command = get_policy(command_copy,val,lambda,pe,rows,action)
11        ;
12        % Policy Improvement
13
14        % If policy converges, exit the while loop and set to true
15        (1)
16        if isequal(command_copy,command) == 0
17            exit = 0;
18        else
19            exit = 1;
20        end
21    end
22    val = get_value(command,lambda,pe,rows);
23 end

```

Listing 11: MATLAB code for running policy iteration and value iteration

```

1 %% EE 183DA Lab 3
2 % Team Buffalo
3 % Iou-Sheng (Danny) Chang    UID: 804-743-003
4 % William                    Argus    UID: 004-610-455
5 % XianXing (Gray) Jiang      UID: 604-958-018
6 % Ho (Bobby) Dong           UID: 604-954-176
7 % Markov Decision Processes (MDPs)
8 clc; clear all; close all;
9
10 %% Robot State Initialization
11 robot = zeros(6);
12 [rows,cols] = size(robot);
13
14 %% 2.1 MDP System
15 % 2.1(a)
16 % REFER to stateSpace R^{3} matrix
17 fprintf('===== 2.1(a) =====\n');
18 fprintf('REFER to stateSpace R^{3} matrix in 'init.m'\n');
19 % 2.1(b)
20 % REFER to action R^{2} matrix
21 fprintf('===== 2.1(b) =====\n');
22 fprintf('REFER to action R^{2} matrix in 'init.m'\n');
23 % 2.1(c)

```

```

24 % REFER to 'probability.m' function for detail
25 fprintf('===== 2.1(c) =====\n');
26 fprintf('REFER to 'probability.m' function for detail\n');
27
28 % 2.1(d)
29 % REFER to 'next_state.m' function for detail
30 fprintf('===== 2.1(d) =====\n');
31 fprintf('REFER to 'next_state.m' function for detail\n');
32
33 %% 2.2 Planning problem
34 % 2.2(a)
35 % REFER to 'reward.m' function for detail
36 % Use boolean to check given state if return correct reward value
37 fprintf('===== 2.2(a) =====\n');
38 fprintf('REFER to 'reward.m' function for detail\n');
39 % border states (reward = -100)
40 fprintf('border states reward value = -100 (yes(1) / no(0)): %d\n',
    ...
41     (reward([1,1]) == -100) && (reward([1,2]) == -100) && (
        reward([1,3]) == -100) && (reward([1,4]) == -100) && (
        reward([1,5]) == -100) &&...
42     (reward([6,2]) == -100) && (reward([6,3]) == -100) && (
        reward([6,4]) == -100) && (reward([6,5]) == -100) && (
        reward([6,6]) == -100) &&...
43     (reward([2,1]) == -100) && (reward([3,1]) == -100) && (
        reward([4,1]) == -100) && (reward([5,1]) == -100) && (
        reward([6,1]) == -100) &&...
44     (reward([1,6]) == -100) && (reward([2,6]) == -100) && (
        reward([3,6]) == -100) && (reward([4,6]) == -100) && (
        reward([5,6]) == -100));
45 % lane markers (reward = -10)
46 fprintf('lane markers reward value = -10 (yes(1) / no(0)): %d\n',...
47     (reward([4,4]) == -10) && (reward([4,5]) == -10));
48 % goal square (reward = 1)
49 fprintf('goal square reward value = 1 (yes(1) / no(0)): %d\n',...
50     (reward([5,5]) == 1));
51 % other states (reward = 0)
52 fprintf('other states reward value = 0 (yes(1) / no(0)): %d\n',...
53     (reward([2,2]) == 0) && (reward([2,3]) == 0) && (reward
        ([2,4]) == 0) && (reward([2,5]) == 0) &&...
54     (reward([3,2]) == 0) && (reward([3,3]) == 0) && (reward
        ([3,4]) == 0) && (reward([3,5]) == 0) &&...
55     (reward([4,2]) == 0) && (reward([4,3]) == 0) &&...
56     (reward([5,2]) == 0) && (reward([5,3]) == 0) && (reward
        ([5,4]) == 0));
57
58 %% 2.3 Policy Iteration

```

```

59 % Use 'Run Section' to only run the policy iteration
60 clc; clear all; close all;
61
62 % Robot State Initialization
63 robot = zeros(6);
64 [rows,cols] = size(robot);
65
66 % Initialization
67 [optimal_value_policy,pe,lambda,s,action] = init();
68 policy_p = init_policy(rows);
69 prer = 0;
70
71 % Using MATLAB built in 'tic toc' command to obtain compute time
72 tic
73 [policy_p,optimal_value_policy] = policy(policy_p,lambda,pe,rows,
    action);
74 computeTime_policyIterationStr = toc;
75
76 % Display initial state
77 robot((rows+1)-s(2),s(1)) = s(3);
78
79 % Plot robot state
80 % Grid World has the same configuration as the figure shown in lab3
    instruction
81 % Value shown on the graph is the heading
82 % Next action and prerotation is shown in the x axis
83 % when we reach the goal block, break out of the loop
84 for j = 1:1000
85     act = policy_p((rows+1)-s(2),s(1),s(3),:);
86     plot_route(robot,prer,act);
87     optimal_policy_policy(j,1) = act(1,1,1,1);
88     optimal_policy_policy(j,2) = act(1,1,1,2);
89     if (s(1) == 5 && s(2) ==5)
90         break;
91     end
92     robot((rows+1)-s(2),s(1)) = 0;
93     [s prer]= next_state(pe,s,act);
94     robot((rows+1)-s(2),s(1)) = s(3);
95 end
96
97 % Command Window
98 % 2.3(a)
99 fprintf('===== 2.3(a) =====\n');
100 fprintf('REFER to ''init_policy.m'' for detail\n');
101 % 2.3(b)
102 fprintf('===== 2.3(b) =====\n');
103 fprintf('REFER to ''plot_route.m'' for detail\n');

```



```

104 % 2.3(c)
105 fprintf('===== 2.3(c) =====\n');
106 fprintf('REFER to ''policy_init.m'' for detail\n');
107 fprintf('REFER to figure plotted or video linked in the lab report
    reference\n');
108 % 2.3(d)
109 fprintf('===== 2.3(d) =====\n');
110 fprintf('REFER to ''get_value.m'' for detail\n');
111 % 2.3(e)
112 fprintf('===== 2.3(e) =====\n');
113 fprintf('REFER to ''policy_init.m'' for detail\n');
114 % 2.3(f)
115 fprintf('===== 2.3(f) =====\n');
116 fprintf('REFER to ''get_policy.m'' for detail\n');
117 % 2.3(g)
118 fprintf('===== 2.3(g) =====\n');
119 fprintf('REFER to ''policy.m'' for detail\n');
120 % 2.3(h)
121 fprintf('===== 2.3(h) =====\n');
122 fprintf('REFER to figure plotted or video linked in the lab report
    reference\n');
123 % 2.3(i)
124 fprintf('===== 2.3(i) =====\n');
125 computeTime_policyIterationStr = ['Policy Iteration Compute Time =
    4.6858 [s]'];
126 %computeTime_policyIterationStr = ['Policy Iteration Compute Time =
    ',num2str(computeTime_policyIterationStr),' [s]'];
127 fprintf(computeTime_policyIterationStr);
128
129 %% 2.4 Value iteration
130 % Use 'Run Section' to only run the value iteration
131 clc; clear all; close all;
132
133 % Robot State Initialization
134 robot = zeros(6);
135 [rows,cols] = size(robot);
136
137 % Initialization
138 [value_prev,pe,lambda,s,action] = init();
139 optimal_value_value = value_prev;
140 value_difference = 10e5;
141 prer = 0;
142
143 % Using MATLAB built in 'tic toc' command to obtain compute time
144 tic
145 while value_difference > 10e-5
146     [value_prev,policy_v] = value(optimal_value_value,action,lambda,

```

```

        pe,rows);
147     value_difference = sum(sum(sum(abs(value_prev -
        optimal_value_value))));
148     optimal_value_value = value_prev;
149 end
150 computeTime_valueIteration = toc;
151
152 % Display initial state
153 robot((rows+1)-s(2),s(1)) = s(3);
154
155 % Plot robot state
156 % Grid World has the same configuration as the figure shown in lab3
    instruction
157 % Value shown on the graph is the heading
158 % Next action and prerotation is shown in the x axis
159 % when we reach the goal block, break out of the loop
160 for j = 1:1000
161     act = policy_v((rows+1)-s(2),s(1),s(3),:);
162     plot_route(robot,prer,act);
163     optimal_policy_value(j,1) = act(1,1,1,1);
164     optimal_policy_value(j,2) = act(1,1,1,2);
165     if (s(1) == 5 && s(2) ==5)
166         break;
167     end
168     robot((rows+1)-s(2),s(1)) = 0;
169     [s prer]= next_state(pe,s,act);
170     robot((rows+1)-s(2),s(1)) = s(3);
171 end
172
173 % Command Window
174 % 2.4(a)
175 fprintf('===== 2.4(a) =====\n');
176 fprintf('REFER to ''optimal_policy_value'' matrix\n');
177 fprintf('REFER to ''optimal_value_value'' matrix\n');
178 % 2.4(b)
179 fprintf('===== 2.4(b) =====\n');
180 fprintf('REFER to figure plotted or video linked in the lab report
    reference\n');
181 % 2.4(c)
182 fprintf('===== 2.4(c) =====\n');
183 computeTime_valueIterationStr = ['Value Iteration Compute Time = ',
    num2str(computeTime_valueIteration),' [s]'];
184 fprintf(computeTime_valueIterationStr);

```

## 5 Value Iteration

### 5.1 Overview

Value iteration uses a value function to assign a value to each state. The higher the value of a given state, the better it is for the robot to be in that state. The value function eq. (7) gives the state a value equal to the total expected rewards the robot will gain from starting in that state and moving according to its prescribed policy.

$$V^\pi(s) = E \left[ \sum_{i=1}^T \gamma^{i-1} r_i \right] \forall s \in \mathbb{S} \quad (7)$$

The optimal policy is found using the Bellman equation eq. (8) to perform value iteration, which improves  $V(s)$  until it converges on the optimal value function eq. (9). The optimal policy eq. (10) is the policy that is used in the optimal value function.

$$V^*(s) = \max_a \left[ R(s, a) + \gamma \sum_{s' \in \mathbb{S}} p(s'|s, a) V^*(s') \right] \quad (8)$$

$$V^*(s) = \max_{\pi} V^\pi(s) \quad \forall s \in \mathbb{S} \quad (9)$$

$$\pi^* = \underset{\pi}{\operatorname{argmax}} V^\pi(s) \quad \forall s \in \mathbb{S} \quad (10)$$

The optimal Q-function eq. (11) will give the total expected reward gained by the robot given that the robot starts in state  $s$ , takes action  $a$ , and then continues to take the best action for the rest of its movement.

$$V^*(s) = \max_a Q^*(s, a) \quad \forall s \in \mathbb{S} \quad (11)$$

### 5.2 Value Iteration

The function returns optimal policy  $\pi^*$  and optimal value  $V$ , assuming  $V(s) = 0 \quad \forall s \in \mathbb{S}$ .

#### Pseudo Code for Value Iteration

```

Assign arbitrary values to V(s)
Rerun
  For all states in the state space
    For all actions in the action space
       $Q(s,a)$ , where  $Q(s,a) = E[r|s,a] + \gamma \sum_{s' \in S} P(s'|s,a)V(s')$ 
     $V(s)$ , where  $V(s)$  is  $\max_a Q(s,a)$ 
  Continue until  $V(s)$  converges

```

Figure 5: Pseudo code for Value Iteration

The value iteration function is implemented in MATLAB as *value.m* (listing 12).

### 5.3 Optimal policy trajectory plot

The plot of the trajectory using optimal policy  $\pi^*$  and  $P_e = 0$ .

The trajectory plotted is shown in video in the Appendix and Demonstration section.

### 5.4 Computation time

We use the MATLAB built in *tic – toc* function to determine the compute time of Value Iteration, as shown in *main.m* (listing 11), the compute time is then printed to the commend window as shown in fig. 6, it is determined to be 8.6225 [s].

```

===== 2.4(a) =====
REFER to 'optimal_policy_value' matrix
REFER to 'optimal_value_value' matrix
===== 2.4(b) =====
REFER to figure plotted or video linked in the lab report reference
===== 2.4(c) =====
Value Iteration Compute Time = 8.6225 [s]>>

```

Figure 6: Compute time of Value Iteration

### 5.5 Matlab Implementation

Listing 12: MATLAB code for Value Iteration

```

1 function [val,command] = value(V,action,lambda,pe,rows)
2   % ===== Initialization ===== %
3   val = V;
4   command = [];
5
6   % ===== ?ompute the value ===== %
7   for x = 1:6

```

```

8     for y = 1:6
9         % Iterate through 12 possible headings
10        for h = 1:12
11            % Iterate through 7 possible actions
12            value_allActions = [];
13            for j = 1:6
14                % Find all possible movements due to prerotation
15                % Note: here the next state and doesn't matter
16                [garbage,s] = probability(pe,[x;y;h],action(j,:),
17                    ,[1,1,1]');
18                % Store the value calculated for each actions
19                % Stored in a temp array
20                value_allActions(j) = s(4,1) * (reward([s(1,1);s
21                    (2,1)]) + lambda*val((rows+1)-s(2,1),s(1,1),s
22                    (3,1))) +...
23                    s(4,2) * (reward([s(1,2);s
24                    (2,2)]) + lambda*val((
25                    rows+1)-s(2,2),s(1,2),s
26                    (3,2))) +...
27                    s(4,3) * (reward([s(1,3);s
28                    (2,3)]) + lambda*val((
29                    rows+1)-s(2,3),s(1,3),s
30                    (3,3)));
31            end
32            % Special case
33            % When we reach the goal (green) square
34            % Robot can perform action (0,0) --> stay still (no
35            movement, no rotation)
36            if (x == 5 && y == 5)
37                [garbage,s] = probability(pe,[x;y;h],action(7,:),
38                    ,[1,1,1]');
39                value_allActions(7) = s(4,1) * (reward([s(1,1);s
40                    (2,1)]) + lambda*val((rows+1)-s(2,1),s(1,1),s
41                    (3,1))) +...
42                    s(4,2) * (reward([s(1,2);s
43                    (2,2)]) + lambda*val((
44                    rows+1)-s(2,2),s(1,2),s
45                    (3,2))) +...
46                    s(4,3) * (reward([s(1,3);s
47                    (2,3)]) + lambda*val((
48                    rows+1)-s(2,3),s(1,3),s
49                    (3,3)));
50            end
51            % Find the best action (maximum value) and store it
52            [M,I] = max(value_allActions);
53            val((rows+1)-y,x,h)= M;
54            % Store movement (forward/backward) into command

```

```
36         command((rows+1)-y,x,h,1) = action(I,1);
37         % Store rotation (left/none/right) into command
38         command((rows+1)-y,x,h,2) = action(I,2);
39     end
40 end
41 end
42 end
```

## 6 Additional Scenarios

For the additional scenarios, we run the Policy Iteration with  $p_e = 0.1$  using the original code, only changing the  $p_e$  value in *init.m* (listing 1). However for the altered reward section (the reward of goal square +1 only applies when the robot is pointing straight down  $h = 6$ , it is 0 otherwise), we duplicated the original code and further changed the *reward.m* function to account for the scenarios.

The modified *reward.m* is combined with the *init.m* into the MATLAB function *add\_init.m* (listing 13).

### 6.1 Trajectory with error

The plot of the trajectory using Policy Iteration under optimal policy  $\pi^*$  and  $p_e = 0.1$ .

The trajectory plotted is shown in video in the Appendix and Demonstration section, we can also see from the video that the compute time is now 6.1113 [s].

### 6.2 Trajectory with altered reward

The plot of the trajectory using Policy Iteration under optimal policy  $\pi^*$  and  $p_e = 0$  &  $p_e = 0.25$ .

The trajectory plotted is shown in video in the Appendix and Demonstration section, we can also see from the video that the compute time for  $p_e = 0$  is 3.8029 [s] (fig. 7) and the compute time for  $p_e = 0.25$  is 4.1386 [s] (fig. 8).

```
===== 2.5(b) =====
REFER to figure plotted or video linked in the lab report reference
Policy Iteration Compute Time = 3.8029 [s]>>
```

Figure 7: Compute time under altered reward  $p_e = 0$

```
===== 2.5(b) =====
REFER to figure plotted or video linked in the lab report reference
Policy Iteration Compute Time = 4.1386 [s]>>
```

Figure 8: Compute time under altered reward  $p_e = 0.25$

### 6.3 Conclusions

Policy and Value iteration achieve the same goal using different algorithms. The compute time of policy iteration was much smaller than the compute time of value iteration. This is

seen in sections 4.10 and 5.4, and is expected, as policy iteration's compute time is usually faster than value iteration compute time, as the policy iteration algorithm converges to the optimal policy faster than the value iteration algorithm converges to the optimal value policy.

The altered reward scenario (section 6.2) resulted in the robot only receiving the value given by the final state when it had a specific heading within the final state, also known as the goal square. The computation time of the policy iteration algorithm was less for this scenario than it was for the general scenario policy iteration computation. This is believed to be caused by the fact that because the final state was specific (there was an exact coordinate and heading required), compared to the general final state (only an exact coordinate required), there were less possible movements to achieve the final state that resulted in gaining reward.

## 6.4 Matlab Implementation

Listing 13: MATLAB code for the altered reward

```

1 function [value,reward,pe,lambda,s_init,action] = init()
2 % Initialization of value
3 value_int= zeros(6);
4 reward_temp = [-100 -100 -100 -100 -100 -100;...
5               -100  0    0    -10  0    -100;...
6               -100  0    0    -10  0    -100;...
7               -100  0    0     0    0    -100;...
8               -100  0    0     0    0    -100;...
9               -100 -100 -100 -100 -100 -100];
10
11 % Expend the value matrix to 12 headings' 3D matrix space
12 [rows,cols] = size(value_int);
13 for x = 1:rows
14     for y = 1:cols
15         for h = 1:12
16             value(x,y,h) = value_int(x,y);
17             reward(x,y,h) = reward_temp(x,y);
18         end
19     end
20 end
21
22 reward(2,5,6) = 1;
23
24 % Initialization of error percentage (pe)
25 pe = 0.25;
26
27 % Initialization of discount factor
28 lambda = 0.9;
29
30 % Initialization of initial state
31 s_init = [2;5;6];

```



```
32
33 % Initialization of actions set
34 % column 1: forward 1 / no movement 0 / backward -1
35 % column 2: leftturn -1 / no rotation 0 / rightturn 1
36 action = [1  -1;...
37           1   0;...
38           1   1;...
39          -1  -1;...
40          -1   0;...
41          -1   1;...
42           0   0];
43 end
```

## 7 Appendix and Demonstration

All class materials and documents can be found on [GitHub](#).

All codes, videos, figures and related documents for Lab 3 can be found on [GitHub](#).

For specific experiment videos and demonstrations please visit the links attached to the topics:

- [Video: MATLAB live robot trajectory using  \$\pi\_0\$ ,  \$p\_e = 0\$](#)
- [Video: MATLAB live robot trajectory for Policy Iteration under  \$\pi^\*\$ ,  \$p\_e = 0\$](#)
- [Video: MATLAB live robot trajectory for Value Iteration under  \$\pi^\*\$ ,  \$p\_e = 0\$](#)
- [Video: MATLAB live robot trajectory for Policy Iteration under  \$\pi^\*\$ ,  \$p\_e = 0.1\$](#)
- [Video: MATLAB live robot trajectory for Policy Iteration under altered reward &  \$\pi^\*\$ ,  \$p\_e = 0\$](#)
- [Video: MATLAB live robot trajectory for Policy Iteration under altered reward &  \$\pi^\*\$ ,  \$p\_e = 0.25\$](#)
- [MATLAB code - original](#)
- [MATLAB code - altered reward](#)
- [Images and Figures](#)
- [Images for compute time](#)

## 8 References

1. Mehta, Ankur. "Design of Robotic Systems I." EE183DA. University of California Los Angeles, Los Angeles, California. Lectures.
2. Mehta, Ankur. "Design of Robotic Systems I." EE183DA. University of California Los Angeles, Los Angeles, California. Lab 3 Instructions.
3. Feinberg, Eugene A., and Adam Schwartz. Handbook of Markov Decision Processes: Methods and Applications. Kluwer Academic Publishers, 2002.
4. Ross, Sheldon M. Introduction to Stochastic Dynamic Programming. Acad. Press, 1995.
5. Alzantot, Moustafa. "Deep Reinforcement Learning Demystified - Policy Iteration, Value Iteration and..." Medium.com, Medium, 9 July 2017.