

Grid World

Markov Decision Processes & Reinforced Learning

Noah Prince, Rodney Shaghoulian, David Young

December 2015

Contents

1	Introduction	2
1.1	Submission Overview	2
1.2	Problem Definition.	2
2	Background	4
2.1	Markov Decision Processes	4
2.2	Grid World	4
2.3	Bellman Equations	5
2.4	Value Iteration & Policy Iteration	6
2.5	Reinforcement Learning	7
3	Overview of Source	8
4	Implementation: Solving MDPs	9
4.1	Representing States and Actions	9
4.2	Representing the GridWorld	9
4.3	Getting the Max Utility Action for a Square	10
4.4	Value Iteration	12
4.5	Policy Iteration	13
4.6	Retrieving the Policy from the calculated Utilities	14
5	Implementation: Reinforcement Learning ("Q-Learning")	15
5.1	Modifications to GridSquare	15
5.2	Initializing the GridWorld	16
5.3	Establishing Q Values	16
5.4	Selecting an Action	16
5.5	Getting the Actual Successor State	16
5.6	TD-Update	17
6	Results: Terminal Reward States	18
6.1	Value Iteration	18
6.2	Policy Iteration	20
7	Results: Non-Terminal Reward States	22
7.1	Value Iteration	22
7.2	Policy Iteration	24
8	Results: Q-Learning	26
9	Analysis & Discussion	29
9.1	Potential Improvements	29
10	Individual Contributions	29

1 Introduction

1.1 Submission Overview

This writeup summarizes the procedure and results of two methods for finding policies in MDPs as well as a free model reinforcement learning technique when the transition model and utilities are unknown. It also contains discussion of observed results and attempts to provide some insight and reflection on the behavior of implemented algorithms. This report was produced for course "CS-440: Artificial Intelligence" at University of Illinois Urbana Champaign.

1.2 Problem Definition.

Grid World MDP

Consider the following environment.

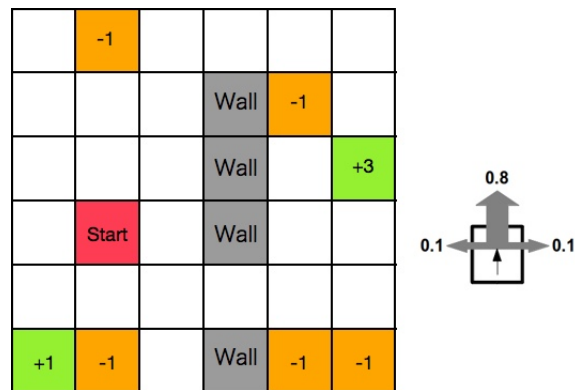


Figure 1: Environment

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 (see the figure above). 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.

Assuming the known transition model and reward function listed above, find the optimal policy and the utilities of all the (non-wall) states using value iteration or policy iteration for two scenarios:

- Reward squares are treated as terminal states: once the agent reaches reward squares, either positive or negative, the agent stops moving.
- Reward squares are treated as non-terminal states: upon reaching reward squares, the agent continues to move. The agent's state sequence is infinite.

For each scenario, display the optimal policy and the utilities of all the states, and plot utility estimates as a function of the number of iterations as in Figure 17.5(a) (for value iteration, you should need no more than 50 iterations to get convergence). In this question and the next one, use a discount factor of 0.99.

Grid World Reinforcement Learning

Consider the reinforcement learning scenario in which the transition model and the reward function are unknown to the agent, but the agent can observe the outcome of its actions and the reward received in any state. (In the implementation, this means that the successor states and rewards will be given to the agent by some black-box functions whose parameters the agent doesn't have access to.)

Use Temporal Difference (TD) Q-Learning to learn an action-utility function only for the terminal scenario described above. Experiment with different parameters for the exploration function and report which choice works the best. For the learning rate function, start with $\alpha(t) = \frac{60}{59+t}$, and play around with the numbers to get the best behavior.

Plot utility estimates and their RMS error (root mean squared error) as a function of the number of trials.

$$RMSE(U', U) = \sqrt{\frac{1}{N} \sum_s (U'(s) - U(s))^2} \quad (1)$$

where $U'(s)$ is the estimated utility of state s , $U(s)$ is the "true" utility as determined by value iteration, and N is the number of states.

2 Background

2.1 Markov Decision Processes

Solving MDPs

Vocabulary

- **States:** s starting with s_0
- **Actions:** a each state s has actions $A(s)$ available to it.
- **Transition Models:** $P(s'|s, a)$ makes use of the Markov assumption: the probability of going to s' from s depends only on s and a and not on any other past actions or states.
- **Reward Functions:** $R(s)$
- **Policy:** π_s is the action that an agent takes in any given state. A policy is the "map" or "solution" to an MDP.

2.2 Grid World

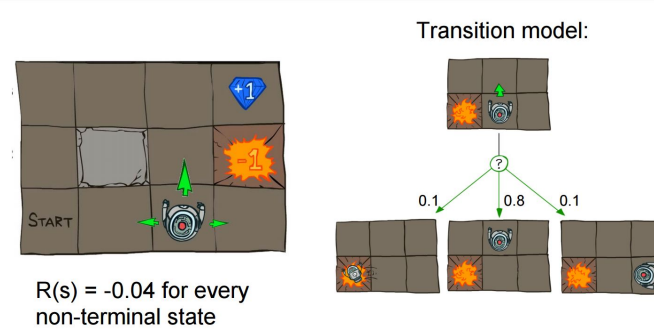


Figure 2: GridWorld Definition.

Goal: Policy

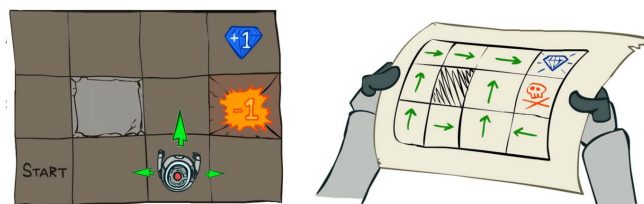


Figure 3: GridWorld Goal: Policy.

2.3 Bellman Equations

The bellman equation describes a recursive relationship between the utilities of successive states given a transition function.

The Bellman equation

- Recursive relationship between the utilities of successive states:

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

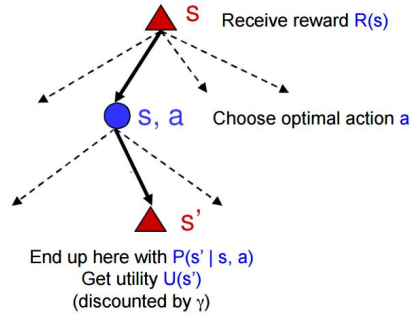


Figure 4: Bellman Equation.

For N states, there are N equations in N unknowns. Solving them will solve the MDP. Unfortunately, trying to solve them through expectimax search might run into trouble with infinite sequences. So instead, solve them algebraically using one of two methods: value iteration and policy iteration.

2.4 Value Iteration & Policy Iteration

Value Iteration:

- Start out with every $U(s) = 0$
- Iterate until convergence
 - During the i th iteration, update the utility of each state according to this rule:

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

- In the limit of infinitely many iterations, guaranteed to find the correct utility values
 - In practice, don't need an infinite number of iterations...

Figure 5: Value Iteration Overview.

Policy Iteration:

- Start with some initial policy π_0 and alternate between the following steps:
 - **Policy evaluation:** calculate $U^{\pi_i}(s)$ for every state s
 - **Policy improvement:** calculate a new policy π_{i+1} based on the updated utilities

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

Figure 6: Policy Iteration Overview.

- Given a fixed policy π , calculate $U^\pi(s)$ for every state s
- The Bellman equation for the optimal policy:

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

- How does it need to change if our policy is fixed?

$$U^\pi(s) = R(s) + \gamma \sum_{s'} P(s' | s, \pi(s)) U^\pi(s')$$

- Can solve a linear system to get all the utilities!
- Alternatively, can apply the following update:

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

Figure 7: Policy Evaluation Overview.

2.5 Reinforcement Learning

Temporal Difference Q-Learning:

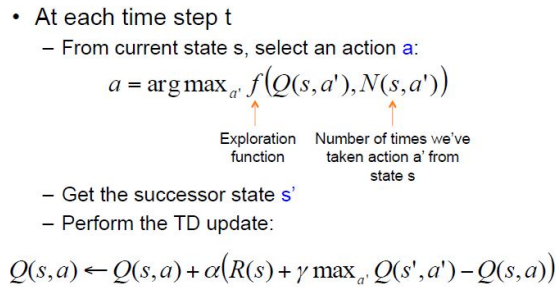


Figure 8: TD Q-Learning Overview.

The main difference in Reinforcement Learning, as compared to the methods we used for Value/Policy iteration, is that 1) The transition model, and 2) the rewards associated with each state, are not known to us. Using a model-free approach, we wish to directly learn the action-utility function $Q(s, a)$ that tells us the value of doing action a in state s .

Exploration vs. Exploitation:

To achieve convergence, we need to balance exploration vs. exploitation. In our initial trials through the maze, we should try to explore to find high-reward states. We try each direction from each state a number of times equal to some threshold value. As our estimates of Q -values gets more accurate, we should exploit more, meaning our actions are decided based off which direction has the highest Q value.

Convergence:

Theoretically, our utility estimates and policy will converge to the same results that value iteration gives us. However, since we are operating in an environment where the transition model and rewards are unknown to us, we will need more iterations and trials through the maze for our results to converge.

3 Overview of Source

Obtaining the source code

The entirety of the code written for this project can be found at the following repository:

<https://github.com/noahprince22/GridWorldMDP/>

Summary of source code

The following source files were written from scratch. All code is well commented with Javadocs; it should be no burden to browse for specific details.

Filename	Description
Direction.java	Enum for movement directions.
DrawingBoard.java	Draws the a simple graphical representation of the gridworld.
GridSquare.java	Holds information for a grid space in the gridworld.
GridWorld.java	Holds the state of the grid world and performs value and policy iteration.
GridWorld_Q.java	Extends GridWorld, and performs Q learning.
Main.java	Entry point to run 1.1.
Main2.java	Entry point to run 1.2.

Additionally, a small third party 2D graphics library was included (StdDraw.java) to provide graphics used in DrawingBoard.java.

4 Implementation: Solving MDPs

4.1 Representing States and Actions

A `gridSquare` class was created to handle the representation of both states and actions. In the grid world, a state is merely a space on the grid and its associated values. Since any action results in the transition to a state, it is possible to represent an action with the state that it leads to. Therefore, both states and actions are represented with the `gridSquare` class.

4.2 Representing the GridWorld

The grid world is represented as an array of `gridSpaces`. The `GridWorld` constructor replicates the grid world environment as outlined in the problem definition, by specifying the wall and reward attributes of appropriate `gridSpaces`.

4.3 Getting the Max Utility Action for a Square

Throughout the policy and value iterations, it is often necessary to get the maximum possibility utility action for a square. This follows the formula given in lecture:

$$\operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s') \quad (2)$$

Key to this implementation is the ability to get the utility from a given action (current square to another gridsquare). This code accomplishes that:

```
1 // Action here is defined as one gridspace to another
2 public double getUtilityForAction(GridSquare currentState, GridSquare successorState) {
3     if (currentState == null || successorState == null)
4         return 0;
5
6     double utility;
7
8     int x = currentState.getxPos();
9     int y = currentState.getyPos();
10
11     GridSquare leftSquare = getGridSquare(x - 1, y);
12     GridSquare rightSquare = getGridSquare(x + 1, y);
13     GridSquare downSquare = getGridSquare(x, y + 1);
14     GridSquare upSquare = getGridSquare(x, y - 1);
15
16     if (isValidLocation(successorState) && ! successorState.isWall())
17         utility = 0.8 * successorState.utility;
18     else
19         utility = 0.8 * currentState.utility; // agent stays in same place as before.
20
21     if (successorState.getxPos() == currentState.getxPos()) { // intended movement is
22         // vertical
23         if (isValidLocation(leftSquare) && ! leftSquare.isWall())
24             utility += 0.1 * leftSquare.utility;
25         else
26             utility += 0.1 * currentState.utility; // agent stays in same place as before.
27
28         if (isValidLocation(rightSquare) && ! rightSquare.isWall())
29             utility += 0.1 * rightSquare.utility;
30         else
31             utility += 0.1 * currentState.utility; // agent stays in same place as before.
32     }
33     else { // intended movement is horizontal
34         if (isValidLocation(upSquare) && ! upSquare.isWall())
35             utility += 0.1 * upSquare.utility;
36         else
37             utility += 0.1 * currentState.utility; // agent stays in same place as before.
38
39         if (isValidLocation(downSquare) && ! downSquare.isWall())
40             utility += 0.1 * downSquare.utility;
41         else
42             utility += 0.1 * currentState.utility; // agent stays in same place as before.
43     }
44
45     return utility;
46 }
```

The following java code calculates the formula using the previous function and returns the Grid-Square which is the state the max action would lead to.

```
1 public GridSquare maxUtilityAction(GridSquare s) {
2     List<GridSquare> validNextSquares = getValidAdjacentSquares(s);
3     if (validNextSquares.size() != 0) {
4         GridSquare maxNode = validNextSquares.get(0);
5         double maxUtility = Double.NEGATIVE_INFINITY;
6
7         // For all valid directions to move, get the node with the expected max utility
8         for (GridSquare s_prime : validNextSquares) {
9             double utility = getUtilityForAction(s, s_prime);
10
11             if (utility > maxUtility) {
12                 maxNode = s_prime;
13                 maxUtility = utility;
14             }
15         }
16         return maxNode;
17     }
18     else
19         return null;
20 }
21 }
```

Along with this function is the helper function maxUtility, which gets the utility from the maxUtilityAction:

```
1 public double maxUtility(GridSquare s) {
2     GridSquare maxAction = maxUtilityAction(s);
3
4     if (maxAction != null)
5         return getUtilityForAction(s, maxAction);
6     else
7         return 0;
8 }
9 }
```

4.4 Value Iteration

Value iteration occurs `numIterations` times. For each value iteration, when the square isn't a wall, we calculate it's utility as:

$$\text{Reward of This Square} + \text{discountFactor} * \text{maxUtility}(\text{square}) \quad (3)$$

After a `GridWorld` is initialized, the following code from `GridWorld.java` performs value iteration to establish utilities for all the grid spaces such that a policy can be determined.

```
1  /* We must not write the utilities to cells until all 36 are calculated */
2  public void establishValueIterationUtilities() {
3      for (int i = 1; i <= numIterations; i++) {
4          double utilities[][] = new double[rows][columns];
5          /* 1st calculate all the utilities before writing to the cells */
6          for (int y = 0; y < rows; y++) {
7              for (int x = 0; x < columns; x++) {
8                  GridSquare currentSquare = getGridSquare(x, y);
9                  if ( ! currentSquare.isWall())
10                     utilities[y][x] = currentSquare.getReward() + discountFactor *
11                         maxUtility(currentSquare);
12             }
13         }
14         /* Now we can copy the utilities */
15         for (int y = 0; y < rows; y++) {
16             for (int x = 0; x < columns; x++) {
17                 grid[y][x].utility = utilities[y][x];
18             }
19         }
20     }
```

4.5 Policy Iteration

A policy iteration consists of two parts: Policy Evaluation and Policy Improvement. In policy evaluation, we update the utilities for all squares according to the following equation. With this equation, we are updating the utilities according to the utility gained by our current policy

$$\text{Utility}(s) = \text{reward}(s) + \text{discountFactor} * \text{Utility for the Action Defined by policy}(s)$$

In policy improvement, we change our policy so that all actions are the actions that produce the most utility under the current model.

After a GridWorld is initialized, the following code from GridWorld.java performs the policy iteration numIterations times.

```
2 public void policyEvaluation() {
3     double utilities [][] = new double[rows][columns];
4     /* 1st calculate all the utilities before writing to the cells */
5     for (int y = 0; y < rows; y++) {
6         for (int x = 0; x < columns; x++) {
7             GridSquare currentSquare = getGridSquare(x, y);
8             if ( ! currentSquare.isWall())
9                 utilities[y][x] = currentSquare.getReward() + discountFactor *
10                    getUtilityForAction(currentSquare, policy[y][x]);
11         }
12     }
13     /* Now we can copy the utilities */
14     for (int y = 0; y < rows; y++) {
15         for (int x = 0; x < columns; x++) {
16             grid[y][x].utility = utilities[y][x];
17         }
18     }
19 }
20 public void policyImprovement() {
21     for (int y = 0; y < rows; y++) {
22         for (int x = 0; x < columns; x++) {
23             policy[y][x] = maxUtilityAction(getGridSquare(x, y));
24         }
25     }
26 }
27
28 public void establishPolicyIterationUtilites() {
29     // Setup initial policy (always go to first in adjacent list)
30     for (int y = 0; y < rows; y++) {
31         for (int x = 0; x < columns; x++) {
32             GridSquare currentSquare = getGridSquare(x, y);
33             List<GridSquare> validAdjacentSquares = getValidAdjacentSquares(currentSquare)
34             ;
35             if (validAdjacentSquares.size() == 0 )
36                 policy[y][x] = null;
37             else
38                 policy[y][x] = validAdjacentSquares.get(0); // policy starts off with all
39                    left-pointing arrows
40         }
41     }
42     for (int i = 0; i < numIterations; i++) {
43         policyEvaluation();
44         policyImprovement();
45     }
46 }
```

4.6 Retrieving the Policy from the calculated Utilities

The policy is continuously updated throughout the iterations. After the iterations have completed, the instance policy variable contains the optimal policy. This policy, in terms of the optimal grid square to move to from the current grid square, can then be analysed.

Analyzing this policy means translating the policy into directions to move at each grid. This is done via the generateDirectionPolicy function which uses the getDirection helper function for every node in the policy. Getting the direction from our definition of actions just means figuring out which direction the action square is relative to the current square. This is done as follows:

```
2 public Direction getDirection(int row, int column, GridSquare destination){  
    if (destination == null)  
        return null;  
4    if (column - 1 == destination.getxPos())  
        return Direction.LEFT;  
6    else if (column + 1 == destination.getxPos())  
        return Direction.RIGHT;  
8    else if (row + 1 == destination.getyPos())  
        return Direction.DOWN;  
10    else if (row - 1 == destination.getyPos())  
        return Direction.UP;  
12    return null; // should never execute  
}
```

5 Implementation: Reinforcement Learning ("Q-Learning")

5.1 Modifications to GridSquare

For the reinforcement Q-Learning, the following fields were added to GridSquare.

```
1 //q values for the 4 possible actions
2 public double qValueLeft;
3 public double qValueRight;
4 public double qValueUp;
5 public double qValueDown;
6
7 //counters for the number of times a given action has been taken
8 public int actionCounterLeft;
9 public int actionCounterRight;
10 public int actionCounterUp;
11 public int actionCounterDown;
```

The following methods were added to GridSquare as well.

```
1 public void updateUtility(){
2     utility = Math.max(qValueLeft, Math.max(qValueRight, Math.max(qValueUp, qValueDown)));
3 }
4
5 public Direction highestUtilityDirection(){
6     if (qValueLeft > Math.max(qValueRight, Math.max(qValueUp, qValueDown)))
7         return Direction.LEFT;
8     else if (qValueRight > Math.max(qValueUp, qValueDown))
9         return Direction.RIGHT;
10    else if (qValueUp > qValueDown)
11        return Direction.UP;
12    else
13        return Direction.DOWN;
14 }
15
16 public Direction leastTriedDirection(){
17     if (actionCounterLeft < Math.min(actionCounterRight, Math.min(actionCounterUp,
18         actionCounterDown)))
19         return Direction.LEFT;
20     else if (actionCounterRight < Math.min(actionCounterUp, actionCounterDown))
21         return Direction.RIGHT;
22     else if (actionCounterUp < actionCounterDown)
23         return Direction.UP;
24     else
25         return Direction.DOWN;
26 }
```


5.2 Initializing the GridWorld

The new fields and constructor for GridWorld_Q are shown below. The new GridWorld_Q class extends the previous GridWorld.

```
1  int threshold;
   GridSquare [][] solutionGrid;
3
   public GridWorld_Q(boolean rewardsTerminal, int iterations, double discountFactor, int
       threshold, GridSquare [][] solutionGrid) {
5       super(rewardsTerminal, iterations, discountFactor);
       this.threshold = threshold;
7       this.solutionGrid = solutionGrid;
   }
```

5.3 Establishing Q Values

```
public void establish_Q_Uilities() {
2     //while(notConverged()){ // can try this instead of the for loop line below
   for (int i = 1; i <= numIterations; i++){
4         GridSquare currentState = start;
         GridSquare intendedSuccessorState = null;
6         GridSquare actualSuccessorState = null;
         while (true){
8             if(currentState.isTerminal()){
                 currentState.utility = currentState.getReward();
10                break;
            }
12            intendedSuccessorState = selectAction(currentState); //this is an action
            Direction intendedDirection = getDirection(currentState,
                intendedSuccessorState);
14            actualSuccessorState = getSuccessorState(currentState, intendedSuccessorState,
                intendedDirection);
            TD.Update(currentState, actualSuccessorState, intendedDirection);
16            updateOtherVariables(currentState, intendedDirection);
            currentState = actualSuccessorState;
18        }
   }
20   policyImprovement();
   generateDirectionPolicy();
22 }
```

5.4 Selecting an Action

When selecting an action, we explore before exploiting. For each state, we first try to explore. If any of the actions from the specific state have been tried less than threshold=500 number of times, we try that action. If multiple actions have been tried less than 500 times, we pick the action with the fewest previous attempts. If the previous trial count of all actions exceeds threshold (500), we select the action that yields the highest $Q(s,a)$, which corresponds to "exploiting".

5.5 Getting the Actual Successor State

Because the agent does not know the transition function it must take actions based off the intended outcome. In reality, an action results in one of three possible movements according to a probability function. To simulate this real probability function, each time the agent makes a move the code uses a virtual dice roll to select the actual state resulting from an action. A random number generator is used such that the intended successor state is chosen with probability 0.8, and the left and right states are chosen with probabilities 0.1 each.

5.6 TD-Update

After doing the action, we need to update $Q(s,a)$. We use the following formula to update $Q(s,a)$, while taking into account the learning rate. Nothing tricky here.

$$Q_{new}(s, a) \leftarrow Q(s, a) + \alpha(R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a)) \quad (4)$$

6 Results: Terminal Reward States

For the Value Iteration and Policy Iteration with Terminal Reward States, 50 iterations were performed. Looking at the utility estimate plots as a function of the number of iterations, one can see this converges.

6.1 Value Iteration

Optimal Policy

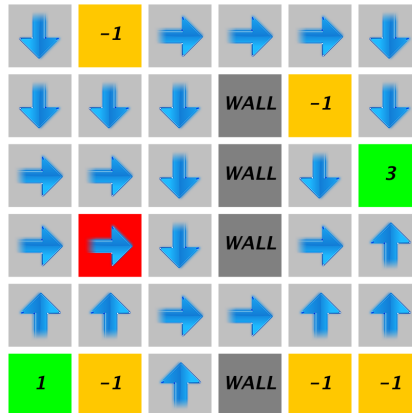


Figure 9: Optimal Policy for terminal case, post value iteration.

Utilities of All States

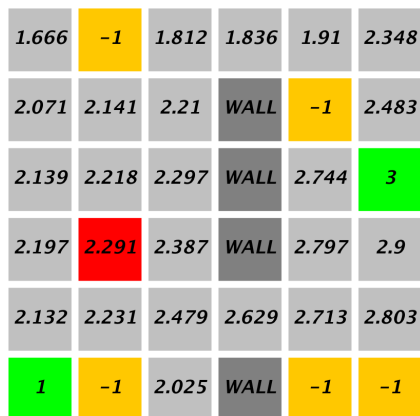


Figure 10: Utilities for terminal case, post value iteration.

Plot of utility estimates as a function of the number of iterations

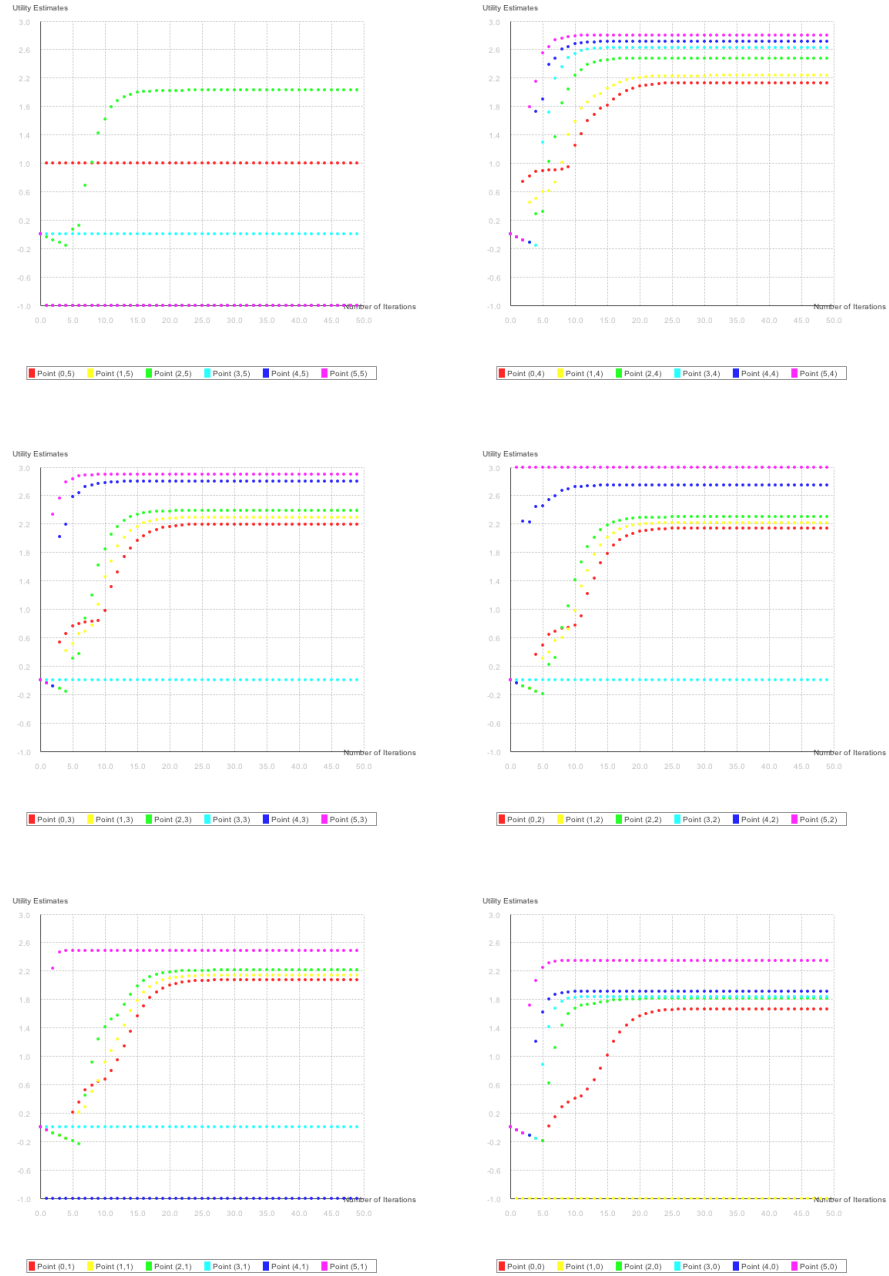


Figure 11: Utility Estimates vs Number of Iterations (Terminal Case)

6.2 Policy Iteration

Optimal Policy

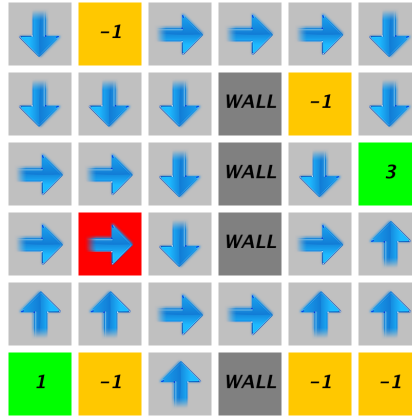


Figure 12: Optimal Policy in the terminal Case, post policy iteration.

Utilities of All States

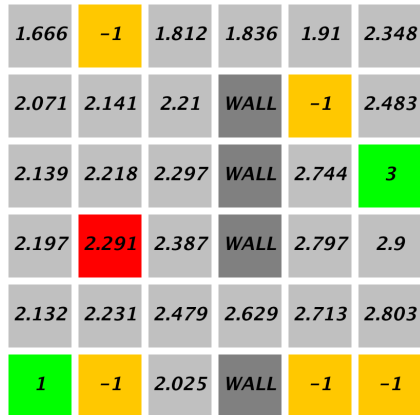


Figure 13: Utilities in terminal case, post policy iteration.

Plot of utility estimates as a function of the number of iterations

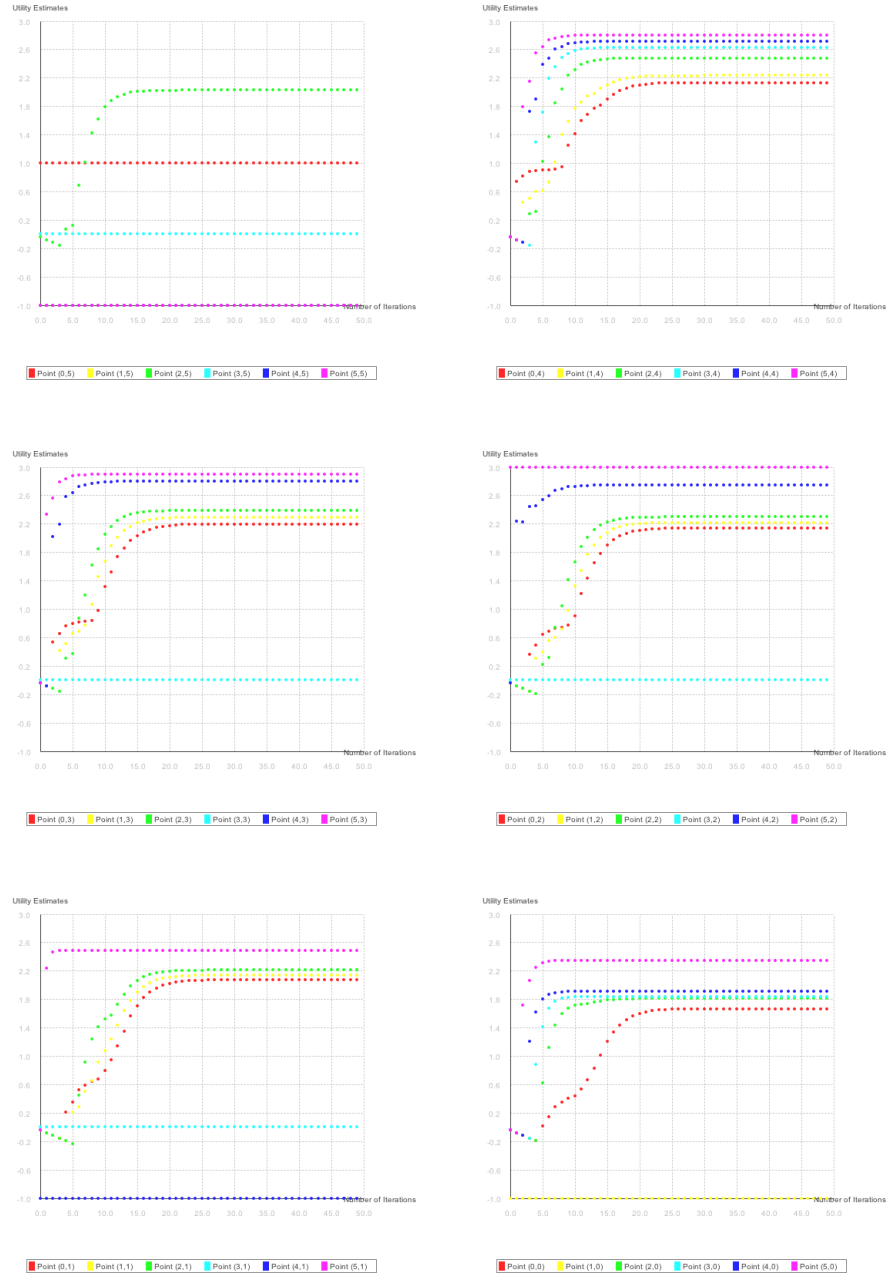


Figure 14: Utility Estimates vs Number of Iterations (Terminal Case)

7 Results: Non-Terminal Reward States

For the Value Iteration and Policy Iteration with Non-Terminal Reward States, 1000 iterations were performed. Looking at the utility estimate plots as a function of the number of iterations, one can see this converges.

7.1 Value Iteration

Optimal Policy

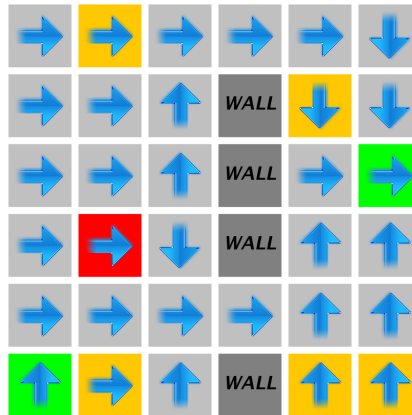


Figure 15: Optimal Policy for non-terminal case, post value iteration.

Utilities of All States

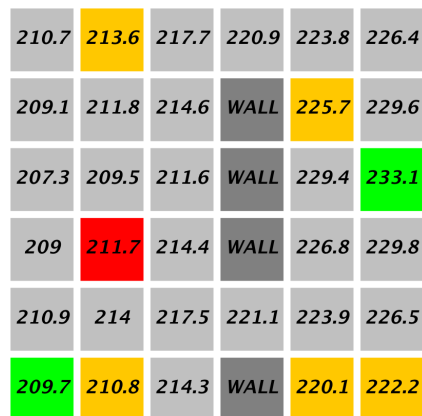


Figure 16: Utilities for non-terminal case, post value iteration.

Plot of utility estimates as a function of the number of iterations

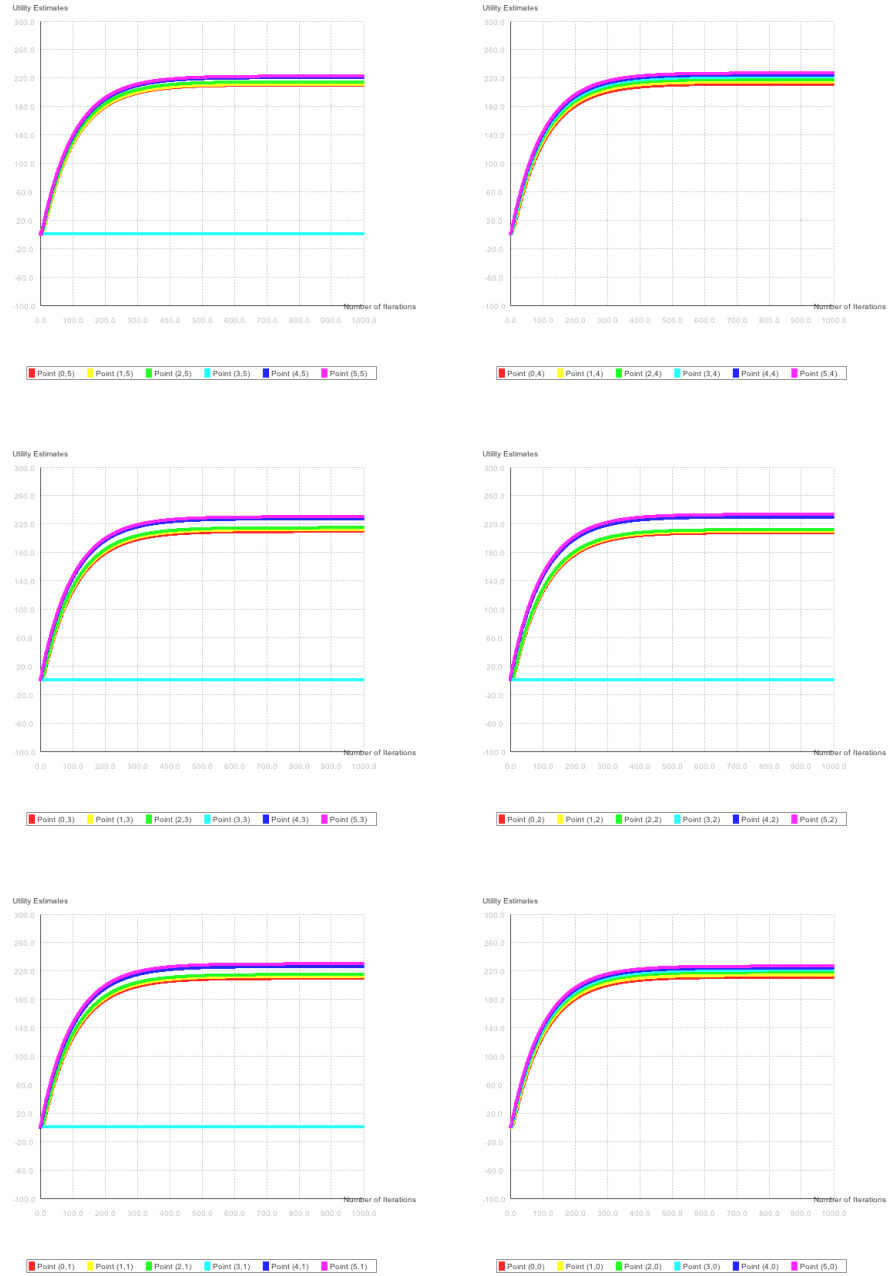


Figure 17: Utility Estimates vs Number of Iterations (Non-Terminal Case)

7.2 Policy Iteration

Optimal Policy

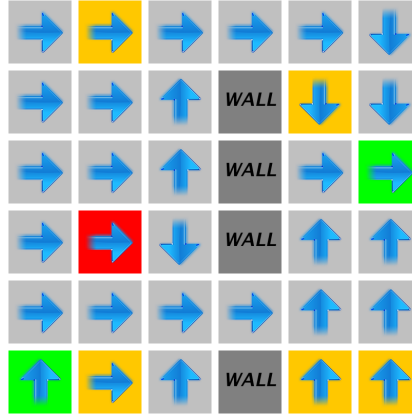


Figure 18: Optimal Policy in the non-terminal Case, post policy iteration.

Utilities of All States

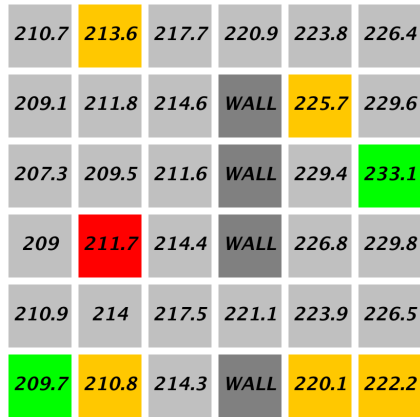


Figure 19: Utilities in the non-terminal case, post policy iteration.

Plot of utility estimates as a function of the number of iterations

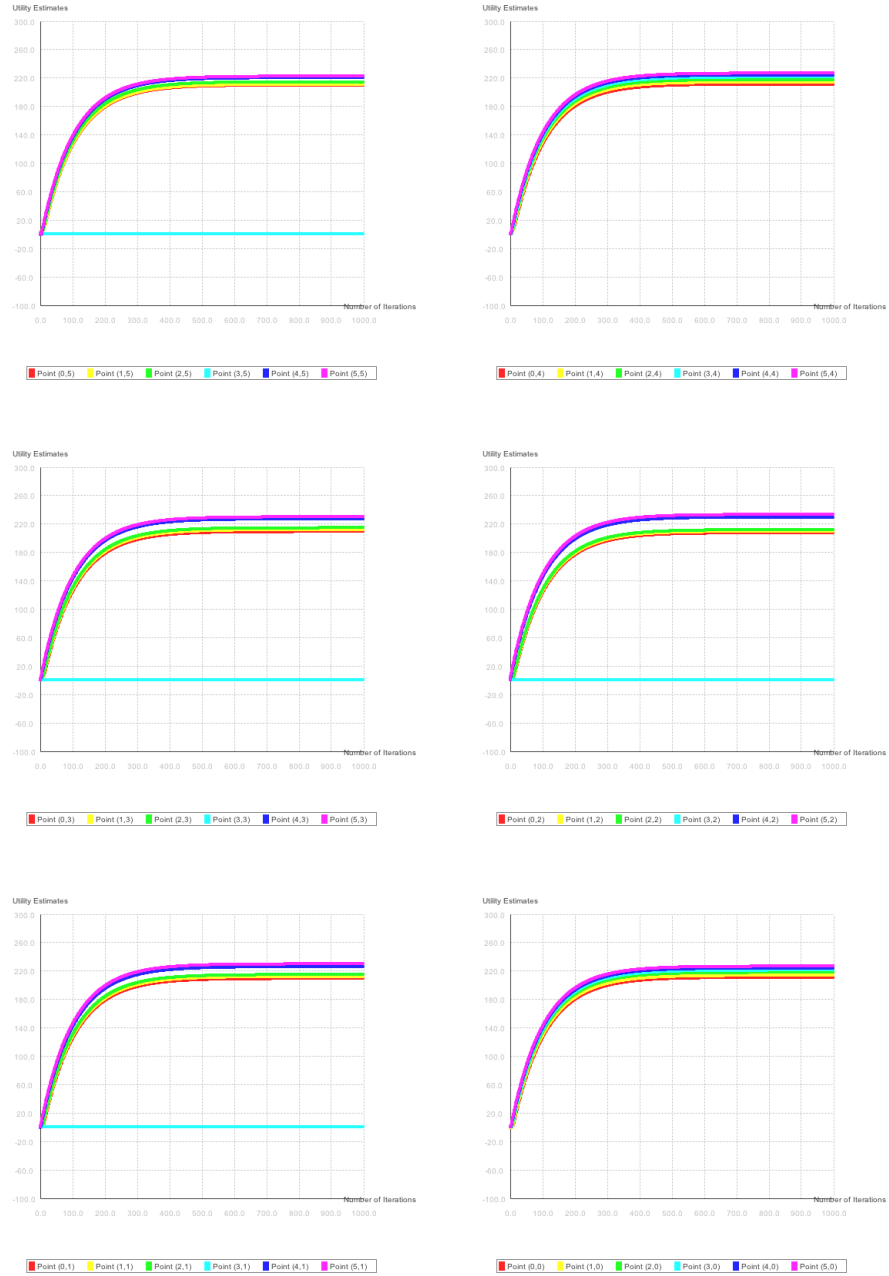


Figure 20: Utility Estimates vs Number of Iterations (Non-Terminal Case)

8 Results: Q-Learning

Utility Estimates



Figure 21: Final Policy (left) and Utility Estimates (right) after Q-Learning

RMS error as a function of the number of trials

Looking at the RMSE plot, it is clear that the error decreases extremely quickly as the number of trials increases. This was expected as major fluctuations occur early in the process.

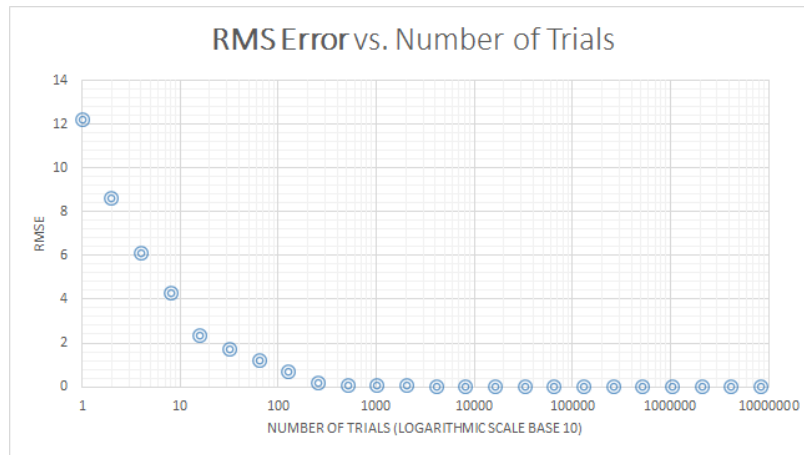


Figure 22: RMS error as a function of the number of trials

Utility plots

Each plot corresponds to the states in a specific row. The general trend is that the utilities converge. For row 0, it is more difficult to see the convergence since we don't try actions from that row as much as from the middle rows. The policy gives us the path to the highest reward state has more accurate utilities (which converge faster) since we try those actions a large number of times. We don't try row 0 actions quite as many times so our plot for row 0 has a slower convergence than the other rows.

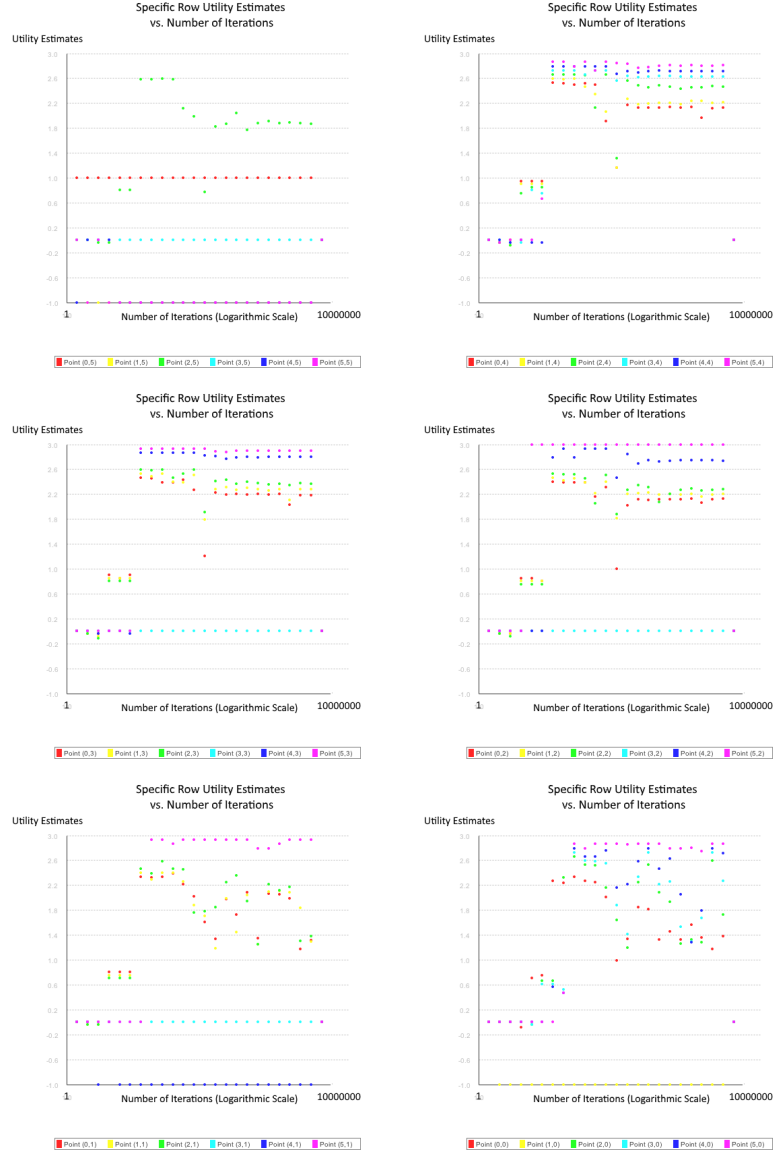


Figure 23: Utility Estimates vs Number of Iterations

Experimenting with Different Parameters

There are two main parameters that we tune to get different results:

1. Exploration Function - Specifically, by changing the threshold value of the number of times we try each direction from a state, we can specify how much exploration we want our agent to do. We currently use threshold = 500, which is a fair amount of exploration before we try exploiting.
2. Learning Rate - Our "t" value for alpha is done slightly differently than the method presented in lecture. We tie our "t" value to each action from a state (i.e. there are 4 "t" values for each state). "t" will range from 1 to 500, so for learning rate, we used a modified formula of $\frac{600}{599+t}$. This learning rate value decays more steadily since t can be as large as 500. Other learning rates such as $\frac{60}{59+t}$ and $\frac{6000}{5999+t}$ were attempted but resulted in slower convergence.