

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/326357651>

Finding Optimal Policy in Grid World Problem Using TDLearning Reinforcement learning project 2

Technical Report · July 2018

CITATIONS

0

READS

24

1 author:



[Sadra Hemmati](#)

Michigan Technological University

17 PUBLICATIONS 1 CITATION

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



NextCar: Next generation Technologies for Connected Vehicles [View project](#)



Fire Protection Analytics [View project](#)

“Finding Optimal Policy in Grid World Problem Using TD-Learning with Eligibility Traces”

Project 2 of the Reinforcement Learning course

Sadra Hemmati

Fall 2015

University of Tennessee, Knoxville

Abstract:

In this project, we implement Temporal Difference Learning ($TD(\lambda)$) to enable an agent to learn the shortest path from any start position to the goal in a maze. Specifically, we implement $SARSA(\lambda)$ algorithm that utilizes eligibility traces to improve the rate of learning for the agent. In our result section, we show how the agent finds the optimal path and the values of the eligibility traces as the agent gets closer to the path. The promising performance of this algorithm gives a measure of strength of this learning paradigm in the specific problem of grid world.

Introduction and Background:

Introduction:

Eligibility traces are one of the basic mechanisms of reinforcement learning. There are two ways to view eligibility traces: One is that these methods are filling the gap from TD to Monte Carlo methods(forward view) and the other is seeing it as a temporary record of an event such as visiting a state or taking an action.

Eligibility traces can be combined with any TD learning method to build more efficient algorithms than other TD-based 01-step method.They are methods that help relate events with state-action pairs and assign rewards to the state-actions depending on the distance they have with the goal state.

Problem Definition:

The problem in this project is defined as the agent in the environment tries to maximize its rewards by reaching the goal state. Here a rectangular grid is used to represent the value functions for a simple finite MDP. Each cell of the grid corresponds to a state and set of possible actions for the agent at each state are defined as: North, South, West and East. After choosing an specific action, the agent deterministically moves one cell in the respective direction on the grid unless the action transfers the agent to outside of the grid or into a wall in which case the location is unchanged. If the agent reaches the goal state, a positive reward of 1 is given to the agent.

We chose to implement Tabular Sarsa(λ) with the addition of 10% exploration as a simple and robust approach to the problem.

In this project we build a reasonably complex grid world of size 20×20 in the excel software and transfer it to MATLAB software. A fixed goal is set up for the mazes and the starts are chosen in a random fashion. An episode is defined as the time in which the agents reaches from the start state to the terminal state. Reaching the goal state has a +1 reward while other states have a zero reward.

At first we made some mazes which depending on the starting point, could form various optimal paths hence making the problem more challenging. Therefore we come up with another set of mazes which have a specific solution ensuring that an optimal policy exists. We added blind alleys to the maze in order to increase the complexity of the maze.

Design:

Design Objective:

Here the high level goal is to implement the Sarsa learning algorithm for the Grid World problem. In order to achieve this, some requirements should be fulfilled which are the design parameters. We summarise these parameters as following:

States:

The state space is represented by a 20 by 20 array. This array includes the wall states but they are never visited. A function was created to determine the next state given a state action pair.

Actions:

As mentioned previously, we have defined four actions for the agent at each step keeping in mind that actions that lead to a wall are not allowed execution. In the table below, the current and next states change in array value after taking different actions is shown below:

State at time t	State at time t+1 after taking action left	State at time t+1 after taking action right	State at time t+1 after taking action up	State at time t+1 after taking action down
(row,col)	(row,col-1)	(row,col+1)	(row+1,col)	(row-1,col)

State at time t	State at time t+1 if action leads outside the maze or into a wall
(row,col)	(row,col)

Rewards:

As explained before, the agent receives a reward of 1 upon reaching the goal and zero for other actions that do not result in reaching the goal.

Implementation:

We implemented SARSA(λ) algorithm based on the Sutton's book and recreated the example following the algorithm description in the book (figure 1). The results of the evolution for the agent with a gamma of 0.9 is shown in figure 2.

```

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
Repeat (for each episode):
   $E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
  Initialize  $S, A$ 
  Repeat (for each step of episode):
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$ 
     $E(S, A) \leftarrow E(S, A) + \delta$  (accumulating traces)
    or  $E(S, A) \leftarrow (1 - \alpha)E(S, A) + \delta$  (dutch traces)
    or  $E(S, A) \leftarrow \delta$  (replacing traces)
    For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :
       $Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$ 
       $E(s, a) \leftarrow \gamma \lambda E(s, a)$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal

```

Figure 7.13: Tabular Sarsa(λ).

Example 7.2: Traces in Gridworld The use of eligibility traces can substantially increase the efficiency of control algorithms. The reason for this is illustrated by the gridworld example in Figure 7.14. The first panel shows the path taken by an agent in a single episode, ending at a location of high reward, marked by the *. In this example the values were all initially 0, and all rewards were zero except for a positive

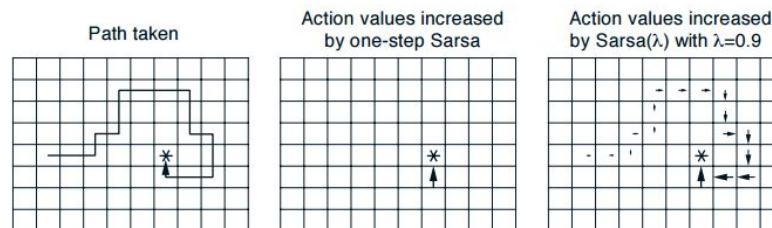


Figure 7.14: Gridworld example of the speedup of policy learning due to the use of eligibility traces.

[illegible]

Initial Trace Example

This example depicts the evolution of Q values in an empty maze as the agent takes different actions and learns by TD method with accumulating eligibility traces. For example in the picture with UP identifier, the state-action values of cells that the agent performs the action of going up has a non-zero value while the other values are zero. As the agent gets closer to the goal, the values increases which is in line with the observation that the 4th row has the highest value since it is closer to the goal state vertically. Similar explanations can be considered for the other actions(down, left, right).

[illegible]

[illegible][illegible][illegible]

AVERAGE*4											
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0.517139	1.900352	0.163968	0	0	0	0	0	0	0	0
0	0.313811	2.089316	0.864085	0.81	0.9	1	0	0	0	0	0
0	0	0.279218	0.401391	0	0	0	0	0	0	0	0
0	0	0.018248	0.392939	0	0	0	0	0	0	0	0
0.00707	0	0.016423	0.185302	0	0	0	0	0	0	0	0
0.019944	0.008728	0.026753	0	0	0	0	0	0	0	0	0
0.005154	0.009698	0.024078	0	0	0	0	0	0	0	0	0
0.008813	0	0	0	0	0	0	0	0	0	0	0

The above average effectively shows the full trace / path to the goal. This path can be seen by following the progression of the values. It is noteworthy that in this figure, the average of cell (4,3) is more than the possible reward of 1. This is the effect of accumulation of traces such that bigger values are assigned to some states. Accumulation trace effects are addressed in challenge 1.

Design Challenges:

There were 3 main challenges when implementing Tabular Sarsa.

1. Accumulating traces leading to perpetually banging into walls.
2. The agent getting lost in endless spirals seeking greed gains
3. Defining lambda for optimal trace length

The first of these was the biggest challenge in running our algorithm because it led to long run times. In instances where the agent took stationary actions right before reaching the end state, the eligibility of the stationary state action pair accumulated disproportionately. This led to future runs only taking this action and running in an infinite loop. This was solved by changing our eligibility tracing method to replacing traces instead of accumulation. We found this worked well for our simple problem statement.

The second problem happened when the agent starts to effectively chase its own tail. This only happened on rare occasions with certain maze designs. This was solved by implementing a 1/10 random chance to take a random action instead of the calculated optimal. However, this problem is maze dependant, and was not an issue with our main maze shown in our result. In that maze, random actions were not used to optimize results.

The last problem was the most open ended. By setting lambda, you effectively determine the length of the traces tail. The main design goal is for it to be long enough to reach out of / into dead ends and optimally reach from any starting point accross the optimal path to the goal. However, this value is set specifically to the maze's unique properties. We decided on a value of .9 after trying to find a value that works robustly with the different kinds of mazes we designed even if it is not 100% optimal for each one.

Technical Approach:

We implemented the following algorithm:

```
Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
Repeat (for each episode):
     $E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
    Initialize  $S, A$ 
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$  or rand  $A$  with 10% chance

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$ 
         $E(S, A) \leftarrow 1$  (replacing traces)
        For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :
             $Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$ 
             $E(s, a) \leftarrow \gamma \lambda E(s, a)$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal
```

Four scripts were used to implement this. These scripts can be found in the Appendix.

1. **TDLambda.m**

This script does everything in the above algorithm. It is the master script that runs through all the iterations and stores all the values.

2. **rewardSmall.m**

This script is a function that returns a reward value when passed a state action pair. It calculates this based on the current maze.

3. **nextStateSmall.m**

This function returns the next state given a state action pair based on the current maze

4. **getDigiEgg.m**

This script automatically generates a maze based on a dimension (such as 4 by 4 or 20 by 20). It would be more useful for future potential application discussed in the summary.

As discussed earlier in challenge 3, lambda was determined to be 0.9. γ was set to allow lambda full control over the eligibility decay and to fully account for the future state action value because if the value isn't 0 then it for sure is on the path to the exit. Alpha was also set to 1 to fully incorporate delta and the eligibility. These values were found experimentally to solve the maze. For some mazes, adjusting alpha and gamma had an effect because of alternative paths but in all cases the trace length was found to be the most important factor.

The algorithm runs by constantly calculating delta based on the seen reward and the value of the next state. For the majority of the first run, these values remain zero until it hits the single reward for the first time. After that, the values of each station action are updated and the trace value is discounted.

These values could of been chosen with more thought, however the given Maze 1 below was solved after a single iteration so they are considered optimal.

One additional note is that for the results in Maze 1, running with 10% random exploration was not optimal and was not used. However, it was used in the additional results of Maze 2.

Experiments and Results:

Below shows our first full sized maze. The green represents the movable path and the red is the reward state. The cell values are arbitrarily set for matlab interpretation.

1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	-1
0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

Maze 1

This maze is very effectively solved by our algorithm. After the first completed run, it will optimally run for future starting values. This can be mainly attributed to the eligibility trace reaching the reward and propagating its value outward. This maze was also designed to be scalable and created automatically with an algorithm. In early testing, we ran smaller versions of the maze seen below:

1	1	0	-1
0	1	1	1
0	1	0	0
1	1	1	1

Maze 1 (4x4)

1	1	1	1	1	0	0	0	0	-1
0	0	0	0	1	1	1	1	1	1
0	0	0	0	1	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0
0	0	0	0	1	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1

Maze 1 (10 x 10)

These smaller versions are useful because of how long the first initial run is.

Eligibility Traces for the moving up

Eligibility Traces for the moving right

Eligibility Traces for the moving down

State-action values for the moving Up

[illegible]

State-action values for the moving right

[illegible]

State-action values for the moving down

5.77E-133	1.34E-132	4.67E-133	3.06E-133	6.51E-132	1.87E-131	6.71E-130	1.92E-129	4.40E-130	1.47E-126	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	3.31E-128	6.08E-06	0.004638398	0.047101287	0.030903154	0.018248D04	0.000176964	0.228767925	0.59049	0.43046721	0
0	0	0	0	0	0	0	0	0	1.21E-128	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	3.99E-06	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.37E-07	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	2.94E-06	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.01E-06	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.72E-06	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	4.17E-10	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	3.37E-10	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	6.98E-11	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	5.06E-11	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.16E-11	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.84E-12	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.61E-17	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	8.56E-18	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	2.95E-12	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	1.14E-12	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	4.42E-19	0	0	0	0	0	0	0	0	0	0
2.04E-36	3.11E-36	1.49E-36	7.89E-37	6.30E-38	5.67E-38	4.59E-38	2.67E-39	2.89E-20	1.01E-13	7.37E-14	6.64E-14	5.37E-14	1.52E-14	1.23E-14	8.07E-15	4.29E-15	8.05E-22	3.07E-23	1.81E-23

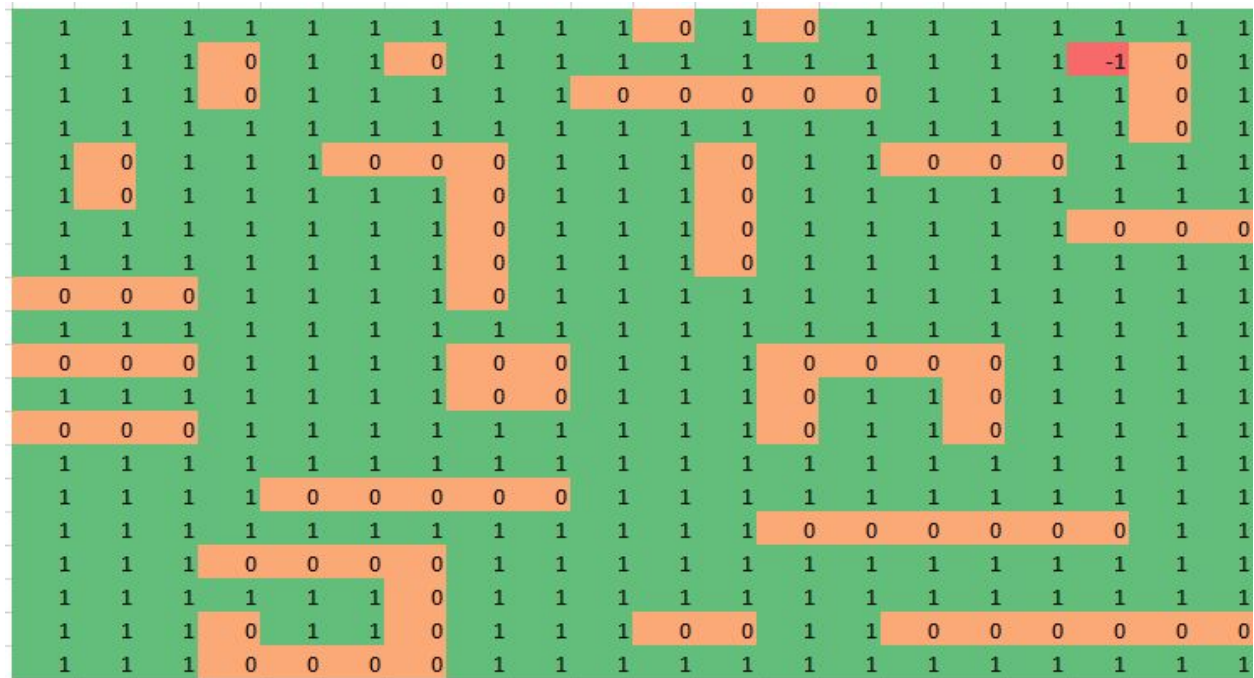
State-action values for the moving left

2	2	2	2	2	2	2	2	2	2	3	0	0	0	0	0	0	0	0	0	
0	0	0	0	0	0	0	0	0	0	2	2	2	2	2	2	2	2	2	2	1
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
2	2	2	2	2	2	2	2	2	2	1	4	4	4	4	4	4	4	4	4	4

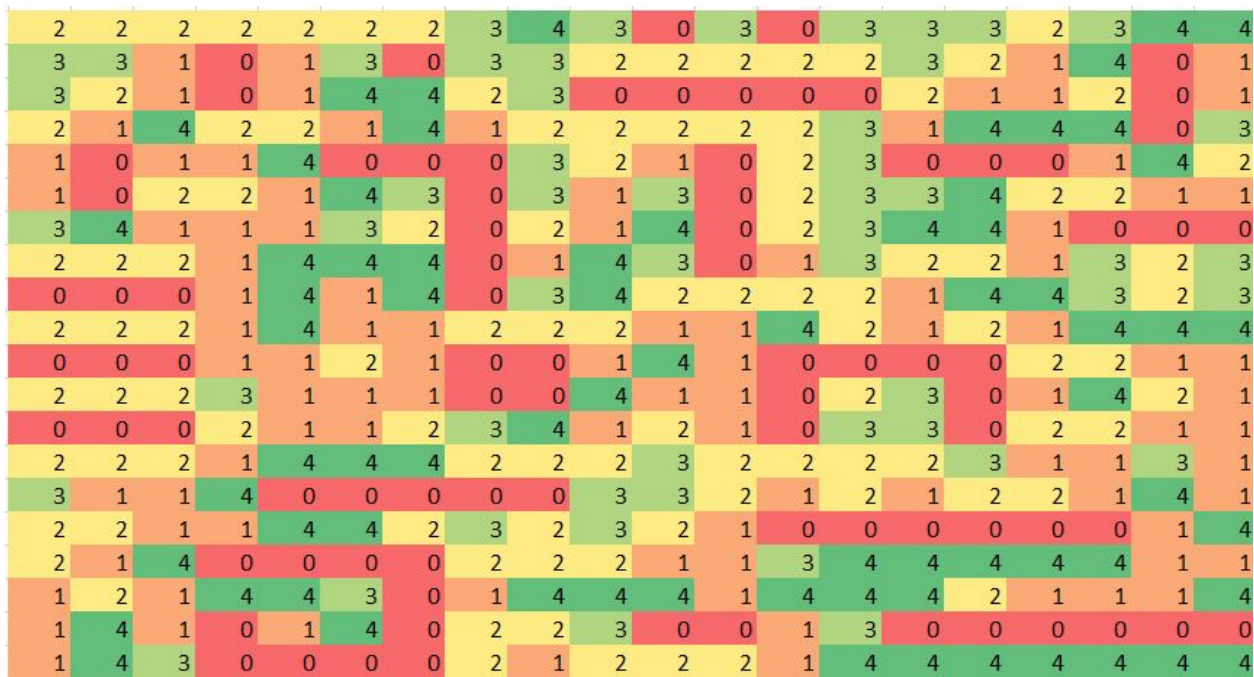
Maze 1 (20 x 20) results after 1 run

Above is a table of the optimal actions calculated after the first iteration of the maze. The actions correspond perfectly to the optimal path that should be taken (1 is up, 2 is right, 3 is down, 4 is left). Every further iteration with random location optimally goes to the goal.

Additionally, we ran our algorithm on other kinds of mazes. However, they do not run as optimally as the above results. They can be seen below:



Maze 2



Maze 2 results after 108 runs

Above it can be seen that the program solves many features of the maze. This can clearly be seen in the bottom right corner. Although we did not calculate it explicitly, the maze runs

functionally well (compared to random navigation). Each run is around 8 to 20 steps which is reasonably fast and moves clearly to the goal. This shows that our solution is robust enough to solve different types of mazes.

Summary

In this project, the formulation of the problem was the first step. In this stage the states and action-state value pairs were mapped to the cells in the maze. The rewards, actions, and factors that affect the learning were initialized in the code. The problem was first addressed by simply implementing eligibility trace TD in an open state space and then iteratively applying it to more complex problems until it fully fulfilled the design scope. Through this process, we learned about many aspects of applied temporal difference learning.

Ultimately, in this project we learned how to implement TD learning with eligibility traces and how to interpret the impact of different parameters (learning, discount, eligibility, etc) on the outcome of the simulation. This knowledge let us overcome a the difficult challenge of not having negative rewards.

Future work could be done in further optimizing the hyper variables to make the most optimal program. With the scalable nature of Maze 1, it would be interesting to scale up the optimal policy with the size of the maze to save on compute time.

Appendix A:

Matlab Code:

TDLambda.m -

```
%Initialize matrixes
clc
clear all
close all
global mazesize
%2 1 2 1 best, avg 9.2
% mazesize = 20;
% maze = getDigiEgg(mazesize);
mazesize = 20;
global mazeOverride
mazeOverride = 1;
global override
override = -2;
%maze = getDigiEgg(mazesize);
maze = getDigiEgg(-2);
exploration = .9;

optSeed = [optA,optG,optL];
oldmaxavgResultVars = zeros(1,3);
maxavgResultVars = ones(1,3);
startRow = mazesize;
startCol = 1;
runs = 1000;
global run;

%Set goalRow and goalCol
for mazerow = 1:mazesize(1)
    for mazecol = 1:mazesize(1)
        if(maze(mazerow,mazecol) == -1)
            goalRow = mazerow;
            goalCol = mazecol;
        end
    end
end

%might need this up here maxavgResultVars = ones(1,3);

Q = zeros(mazesize(1),mazesize(1),4);
E = zeros(mazesize(1),mazesize(1),4);

row = startRow;
col = startCol;
nextRow= startRow;
nextCol= startCol;
maxLimit = 1000;

%New starting position
rowStarts = zeros(1,runs);
colStarts = zeros(1,runs);

for(ij = 1:runs)
```

```

newRow = randi(mazesize(1),1);
newCol = randi(mazesize(1),1);
while (maze(newRow,newCol) == 0 || maze(newRow,newCol) == -1)
    newRow = randi(mazesize(1),1);
    newCol = randi(mazesize(1),1);
end
rowStarts(ij) = newRow;
colStarts(ij) = newCol;
end

moves = zeros(1,runs);
limit = 0
%Run Temporal code
for run = 1:runs
    limit = limit + 1
    limit = 0;
    E = zeros(mazesize(1),mazesize(1),4); %reset E
    moves(run) = 0;
    while (~(row == goalRow && col == goalCol)) && (limit < maxLimit))
        limit = limit + 1;
        moves(run) = moves(run) + 1;
        %Finds optimal action, picks random if tie
        temp = Q(row,col,:);
        temp2 = find(temp==(max(max(temp))));

        optimalA = temp2(randi([1 length(temp2)],1));
        %if(rand(1) > exploration)
        %    optimalA = (randi([1 4],1))
        %end
        nextReward = rewardSmall(row,col,optimalA);
        [nextRow, nextCol] = nextStateSmall( row,col,optimalA);

        temp = Q(nextRow,nextCol,:);
        temp2 = find(temp==(max(max(temp))));
        optimalAprime = temp2(randi([1 length(temp2)],1));

        delta = nextReward + gamma * Q(nextRow,nextCol,optimalAprime) - Q(row,col,optimalA);
        %E(row,col,optimalA) = E(row,col,optimalA) + 1;
        %E(row,col,optimalA) = (1-alpha)*E(row,col,optimalA) + 1
        E(row,col,optimalA) = 1;
        avgE = E(:, :, 1) + E(:, :, 2) + E(:, :, 3) + E(:, :, 4)
        if(limit > maxLimit*.9)
            pause = 1;
            disp(row)
            disp(col)
            disp(optimalA)
            disp(Q(row,col,optimalA))
            disp(optimalAprime)
            disp(Q(nextRow,nextCol,optimalAprime))
        end

        for staterow = 1:mazesize(1)
            for statecol = 1:mazesize(1)
                for action = 1:4
                    Q(staterow,statecol,action) = Q(staterow,statecol,action) +
alpha*delta*E(staterow,statecol,action);
                    E(staterow,statecol,action) = gamma * lambda * E(staterow,statecol,action);
                end
            end
        end

        row = nextRow;
        col = nextCol;
    end
end

```

```
end

if(limit == maxLimit)
    for iii = run:runs
        moves(iii) = maxLimit*100;
    end

end

row = rowStarts(run);
col = colStarts(run);
end
```


rewardSmall.m -

```
function [ myReward ] = rewardSmall( row,col,action)

nextRow = row;
nextCol = col;

standStillReward = 0;
goalReward = 1;
% maze = [1 1 0 -1;
% 0 1 1 1;
% 0 1 0 0;
% 1 1 1 1];
global mazesize
global mazeOverride
global override
if(mazeOverride ~= 1)
maze = getDigiEgg(mazesize);
else
maze = getDigiEgg(override);
end
%maze = getDigiEgg(-1); %temp

%1 up, 2 right, 3 down, 4 left
if(action == 1)
nextRow = nextRow - 1;
end
if(action == 2)
nextCol = nextCol + 1;
end
if(action == 3)
nextRow = nextRow + 1;
end
if(action == 4)
nextCol = nextCol - 1;
end

if(nextRow > mazesize(1) || nextRow < 1 || nextCol > mazesize(1) || nextCol < 1)
%Illegal move, stand still
myReward = standStillReward;
elseif(maze(nextRow, nextCol) == 0)
myReward = standStillReward;
elseif(maze(nextRow, nextCol) == -1)
myReward = goalReward;
else
myReward = 0;
end

end
```

nextStateSmall.m -

```
function [ finalNextRow, finalNextCol ] = nextStateSmall( row,col,action)
    nextRow = row;
    nextCol = col;

    % maze = [1 1 0 -1;
    %         0 1 1 1;
    %         0 1 0 0;
    %         1 1 1 1];
    global mazesize
    global mazeOverride
    global override
    if(mazeOverride ~= 1)
        maze = getDigiEgg(mazesize);
    else
        maze = getDigiEgg(override);
    end
    %maze = getDigiEgg(-1); %temp
    %1 up, 2 right, 3 down, 4 left
    if(action == 1)
        nextRow = nextRow - 1;
    end
    if(action == 2)
        nextCol = nextCol + 1;
    end
    if(action == 3)
        nextRow = nextRow + 1;
    end
    if(action == 4)
        nextCol = nextCol - 1;
    end

    if(nextRow > mazesize(1) || nextRow < 1 || nextCol > mazesize(1) || nextCol < 1)
        %Illegal move, stand still
        finalNextRow = row;
        finalNextCol = col;
    elseif(maze(nextRow, nextCol) == 0)
        finalNextRow = row;
        finalNextCol = col;
    else
        finalNextRow = nextRow;
        finalNextCol = nextCol;
    end

end
```

getDigiEgg.m -

```
function [ matrix ] = getDigiEgg(evolution)
%     if(evolution == 1)
%         maze = [1 1 0 -1;
%                 0 1 1 1;
%                 0 1 0 0;
%                 1 1 1 1];
%     end
maze = zeros(evolution, evolution);
%floor(evolution)
%bottom slice
if(evolution>0)
for(filler = 1:evolution)
    maze(evolution, filler) = 1;
end
%column
for(filler = 1:evolution)
    maze(filler, floor(evolution/2)) = 1;
end
%top
for(filler = 1:floor(evolution/2))
    maze(1, filler) = 1;
end
%end branch
for(filler = 1:ceil(evolution/2))
    maze(2, evolution + 1 - filler) = 1;
end

%goal
maze(1, evolution) = -1;

end
if(evolution==1)
maze = [1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 -1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1];
end
if(evolution==2)
maze = [1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1;
1 1 1 0 1 1 0 1 1 1 1 1 1 1 1 1 -1 0 1;
1 1 1 0 1 1 1 1 1 1 0 0 0 0 1 1 1 1 0 1;
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1;
1 0 1 1 1 0 0 0 1 1 1 0 1 1 0 0 0 1 1 1;
1 0 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1 0 0;
1 1 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1 1 1;
0 0 0 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1;
0 0 0 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 1 1;
0 0 0 1 1 1 1 0 0 1 1 1 0 0 0 0 1 1 1 1;
1 1 1 1 1 1 1 0 0 1 1 1 0 1 1 0 1 1 1 1;
0 0 0 1 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1;
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1;
1 1 1 1 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 1 1;
1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1;
1 1 1 0 1 1 0 1 1 1 0 0 1 1 0 0 0 0 0 0;
1 1 1 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1];
end
```

```

if(evolution==3)
maze=[1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1;
1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 -1 1;
1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 1 1;
1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 0 1 1;
1 1 1 1 0 0 0 0 0 0 0 1 1 1 1 1 0 1 1;
1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 0 1 1;
1 1 0 0 0 1 0 0 0 0 0 1 1 1 1 0 0 1 1;
1 1 1 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1 1;
1 0 1 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1 1;
1 0 1 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1 1;
1 0 1 0 0 0 0 1 1 1 1 1 0 1 1 0 1 1 1;
1 0 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1;
1 0 1 1 0 0 0 1 1 1 1 1 0 1 1 0 0 0 0;
1 0 1 1 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1;
1 0 1 1 0 1 1 1 1 1 1 0 1 0 0 0 0 1 0;
1 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1;
1 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1;
1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1];
end
if(evolution==4)
maze = [1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 -1 0 0;
1 0 0 1 1 1 1 1 1 1 0 1 1 1 1 0 1 0 0;
1 0 0 1 0 0 1 0 0 1 0 0 1 0 0 0 1 0 0;
1 1 1 1 0 0 1 0 0 0 0 0 1 0 0 0 1 1 1;
0 0 0 1 0 0 1 1 1 1 0 1 0 0 1 1 0 0 0;
0 0 0 1 0 0 0 0 0 0 1 0 1 1 0 1 0 0 0;
0 0 0 1 0 1 0 0 0 0 1 0 0 1 0 1 1 1 1;
0 1 1 1 0 1 1 1 0 1 0 0 1 0 1 0 0 0 0;
0 1 0 1 0 0 0 1 0 1 0 0 1 0 1 1 1 1 1;
0 1 0 1 1 0 0 1 0 1 1 1 1 0 0 0 0 0 0;
0 1 0 0 1 1 0 1 0 0 0 0 0 0 0 1 1 0 0;
0 0 0 0 0 1 1 1 1 0 0 1 0 1 1 1 1 0 0;
0 0 1 1 0 0 1 1 0 0 1 0 0 0 0 0 0 1 1;
0 0 1 0 0 0 1 0 0 1 1 1 1 0 0 0 0 1 0;
0 0 1 1 1 1 1 1 0 0 0 0 0 1 1 1 1 0;
0 0 0 0 0 0 1 1 1 0 0 0 1 0 0 0 1 0 0;
0 1 1 1 1 0 0 0 1 0 0 0 0 0 0 1 1 1 0;
0 1 0 0 1 0 1 0 1 1 1 1 1 1 0 1 0 0 0;
0 1 0 0 1 0 1 0 0 0 0 0 0 1 1 1 0 0 0;
1 1 0 0 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0];
end

matrix = maze;
end

```