Ketan Karnakota
3 October 2019

# MDP Techniques

## Introduction:

Our central research pursuit is to compare the speed of different techniques for Markov Decision Processes (MDPs) in the context of solving mazes. Two basic algorithms for solving MDPs are value iteration and policy iteration.

The MDP selected is to solve mazes. There are 4 choices the agent in the maze can take viz. Up, Down, Left, Right. And the states are the positions in the maze. We are considering deterministic actions. The Reward is in the first block that is the (0,0) which is denoted by red. This has a reward of 10. The agent or the mouse gets in from the bottom right and needs to go up to the first block that has the reward or cheese.

## Setup

There is a grid that has '#'s in them as walls and the other parts are the paths to the goal. We have policies that point to the next state we take. We need to come-up/improve a policy such that we reach the target reward state in the least number of steps that is without wandering around and repeating the states in the grid.

To visualize the algorithm there is a LOGO kind of a plugin turtle that is used to visualize actions that are taken at each state to understand what the policy looks like.

## Policy Iteration

Starting with a random policy $\pi 0$ , this approach consists of two steps:

I. Policy Evaluation

$$V(s) = \sum_{s' \in S} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V(s')], \forall s \in S$$

II. Policy Improvement

$$\pi(s) \leftarrow \underset{a}{argmax} \sum_{s' \in S} T(s, a, s')[R(s, a, s') + \gamma V(s')], \forall s \in S$$

These steps are repeated until π converges to π*

## Value Iteration

The value iteration approach finds the optimal policy π* by calculating the optimal value function, V*. Starting with V(s) = 0 for all states s, the values of each state are iteratively updated to get the next value function V, which converges towards V*.
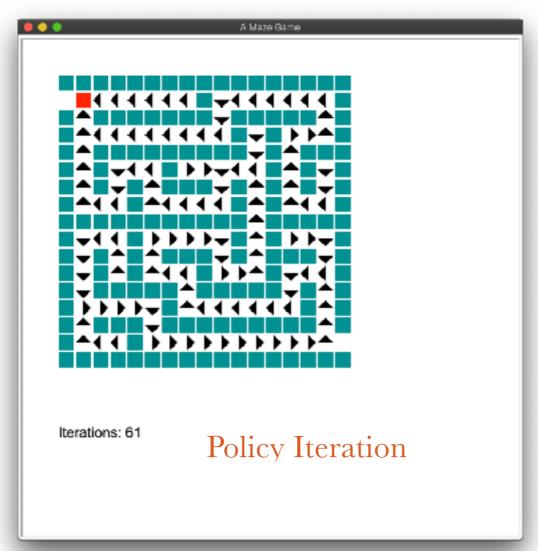
$$V_{i+1}^* \leftarrow \underset{a}{argmax} \sum_{s' \in S} T(s, a, s')[R(s, a, s') + V_i^*(s')]$$
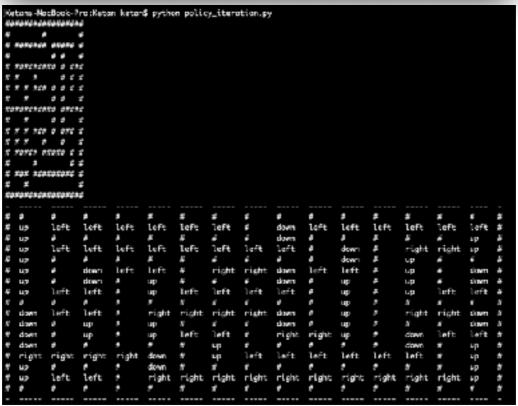
## Files List:

AnimationFile.py - Has the turtle Animation to visualize the policy and it's updates

maze.py - This creates a solvable maze for the size given in the calling classes

maze2mdp.py - Converts the maze matrix to an MDP problem with actions and states

mdpstate.py - The definition class to define the variables/parameters for MDP

policy_iteration.py - Actual code for PI having animations for the optimal policy.

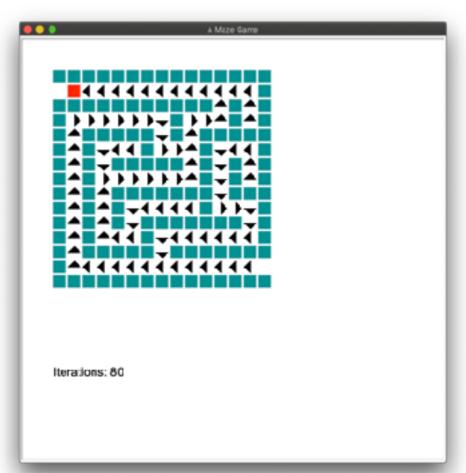value_iteration.py - Code for VI without visualization that computes is less time.

vi_visualization.py - Code for VI with visualization and animation this takes way longer than the other version because it has to check the actions for all the stages to show it in the animation. This can be used to check the policy/(actions for each state) for every VI update for V(s) for the states in the matrix.
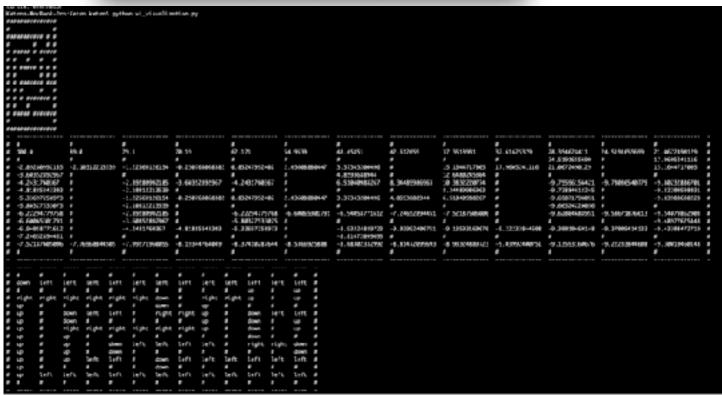
Policy Iteration

A Maze Game

Iterations: 80