

# CS489 HW1 Report

Sheng Cen

March 2021

## 1 Introduction

In known Markov reward/decision process, it's necessary for us to evaluate certain policy and try to improve it, so that we could achieve optimal value given each state. Usually, we need to solve Bellman Expectation Equation and Bellman Optimality Equation, which cannot be tackled in acceptable time periods when the number of states is large. Dynamic programming is designed to use iterative methods to find desired value/policy for each state, which can usually reach convergence quickly. This report experiments on a simple grid world game, using dynamic programming to evaluate the random policy, and try to improve it using both policy iteration and value iteration. Comparison of these two methods is analyzed in the final part, and we conclude the report by assessing dynamic programming used in this setting.

## 2 Setting of Grid World

In the 6\*6 grid world we have set up (see Figure 1), each grid represents a state. Grid 1 and 35 are two terminated states, in which we expect to reach as quickly as possible. In each state, four actions (towards east, west, south and north) could be taken, and in order to reach terminated states quickly, we penalize and assign a reward of -1 for each action. Also, actions are bounded in the 6\*6 grid world, meaning that agent could not get out of the world. The initial policy we use is random policy, i.e., each direction we take in any non-terminated state is assigned to a probability of 1/4.

## 3 Dynamic Programming in Markov Reward and Decision Process

### 3.1 Policy Evaluation

For a known Markov reward/decision process, it's necessary to evaluate a given policy an agent plans to take. For each state, a value would be assigned to reflect its expected return, under the given policy. The value functions for all

0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
18	19	20	21	22	23
24	25	26	27	28	29
30	31	32	33	34	35

Figure 1: The Grid World

states could be obtained through dynamic programming, when convergence is reached. The pseudocode for the policy evaluation is shown in Figure 2.

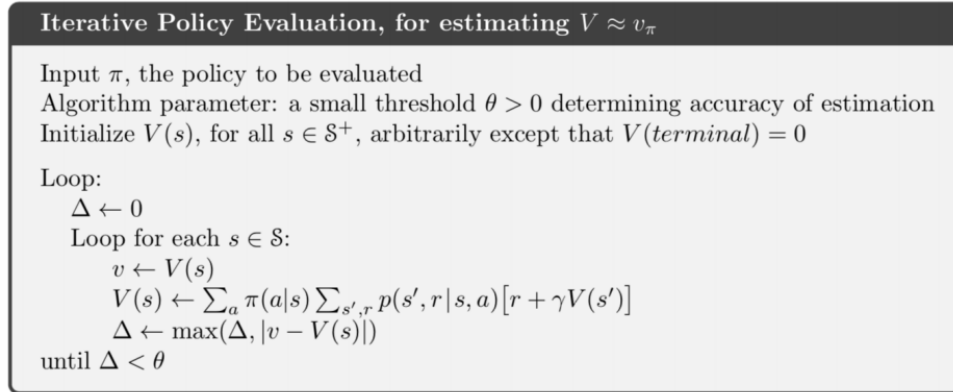


Figure 2: Policy Evaluation

### 3.2 Policy Iteration

It's meaningless to only evaluate a given policy, so improvement should be executed after we evaluate the policy. Policy would be updated only if there exists a state that would gain a higher value when taken another action. Policy im-

provement and evaluation need to be done repetitively until we reach an optimal policy. The pseudocode for the policy iteration is shown in Figure 3.

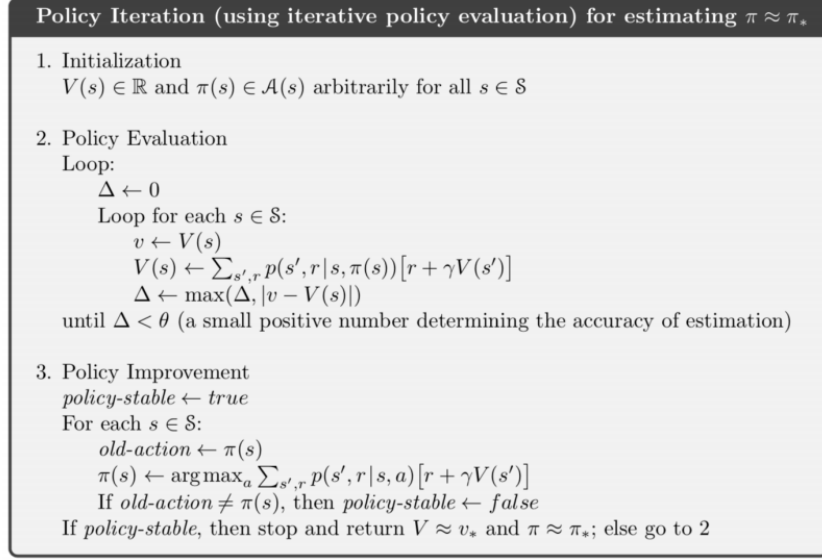


Figure 3: Policy Iteration

### 3.3 Value Iteration

We could also improve the policy by finding the optimum value function for each state first in iterative manner, then choose the optimum policy given these value functions. Therefore, we do not need to iteratively evaluate intermediate policy, saving lots of time. The pseudocode for the value iteration is shown in Figure 4.

## 4 Experiment

We use Python 3 to simulate the grid world, and use dynamic programming to effectively evaluate and improve the random policy.

### 4.1 Policy Evaluation

We set the discounting factor  $\gamma$  to be 1, and get the resulting value functions for each state. The accuracy threshold  $\theta$  is set to be 0.001. Value function for each state is shown in Figure 5, and we need to iterate 216 times to reach convergence.

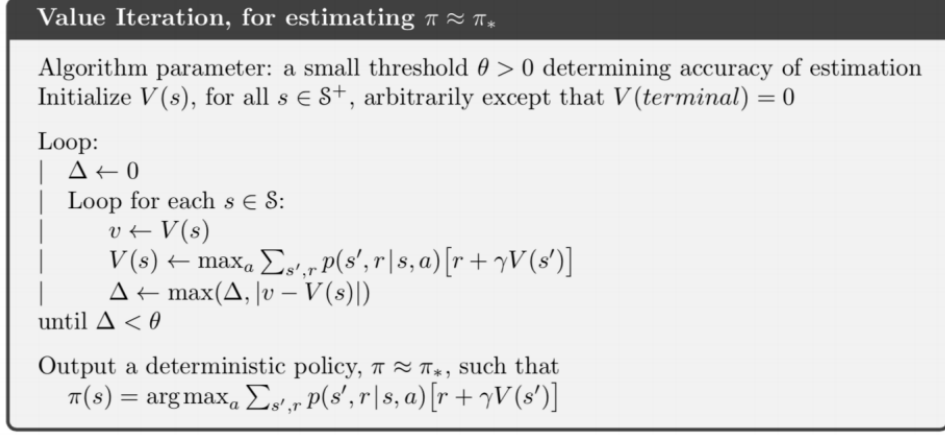


Figure 4: Value Iteration

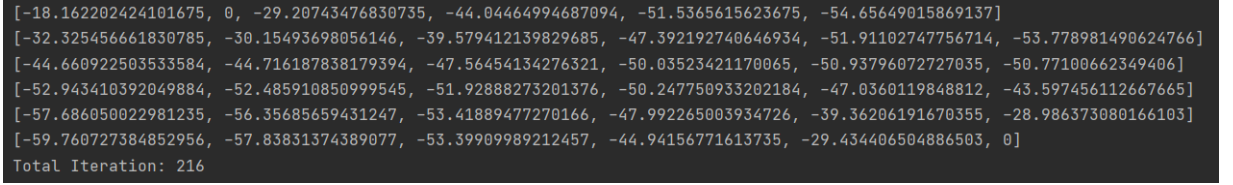


Figure 5: Result for Policy Evaluation

## 4.2 Policy Iteration

We set the discounting factor  $\gamma$  to be 0.9, and the accuracy threshold  $\theta$  is set to be 0.001. Probabilistic policy matrix for each state is shown in Figure 6, while value function is shown in Figure 7. Policy gets stable after three rounds of improvement.

## 4.3 Value Iteration

We set the discounting factor  $\gamma$  to be 0.9, and the accuracy threshold  $\theta$  is set to be 0.001. Value function for each state is shown in Figure 8, while probabilistic policy matrix is shown in Figure 9. Optimum value functions are obtained after six rounds of improvement.

# 5 Analysis for Experiment

For the policy evaluation part, we notice that the final stable value function is intuitive. Expected return is small in the upper-right and bottom-left part of

```

Up probabilistic matrix:
[0, 0, 0, 0, 0, 0]
[0.5, 1.0, 0.5, 0.5, 0.5, 0]
[0.5, 1.0, 0.5, 0.5, 0, 0]
[0.5, 1.0, 0.5, 0, 0, 0]
[0.5, 1.0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 0]
Down probabilistic matrix:
[0, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 1.0]
[0, 0, 0, 0, 0.5, 1.0]
[0, 0, 0, 0.5, 0.5, 1.0]
[0, 0, 0.5, 0.5, 0.5, 1.0]
[0, 0, 0, 0, 0, 0]
Left probabilistic matrix:
[0, 0, 1.0, 1.0, 1.0, 1.0]
[0, 0, 0.5, 0.5, 0.5, 0]
[0, 0, 0.5, 0.5, 0, 0]
[0, 0, 0.5, 0, 0, 0]
[0, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 0]
Right probabilistic matrix:
[1.0, 0, 0, 0, 0, 0]
[0.5, 0, 0, 0, 0, 0]
[0.5, 0, 0, 0, 0.5, 0]
[0.5, 0, 0, 0.5, 0.5, 0]
[0.5, 0, 0.5, 0.5, 0.5, 0]
[1.0, 1.0, 1.0, 1.0, 1.0, 0]

```

Figure 6: Probabilistic policy matrix for Policy Iteration

```

Value matrix:
[-1.0, 0, -1.0, -1.9, -2.71, -3.439]
[-1.9, -1.0, -1.9, -2.71, -3.439, -3.439]
[-2.71, -1.9, -2.71, -3.439, -3.439, -2.71]
[-3.439, -2.71, -3.439, -3.439, -2.71, -1.9]
[-4.0951, -3.439, -3.439, -2.71, -1.9, -1.0]
[-4.0951, -3.439, -2.71, -1.9, -1.0, 0]
Improvement Count:
3

```

Figure 7: Value Function for Policy Iteration

```

Value matrix:
[-1.0, 0, -1.0, -1.9, -2.71, -3.439]
[-1.9, -1.0, -1.9, -2.71, -3.439, -3.439]
[-2.71, -1.9, -2.71, -3.439, -3.439, -2.71]
[-3.439, -2.71, -3.439, -3.439, -2.71, -1.9]
[-4.0951, -3.439, -3.439, -2.71, -1.9, -1.0]
[-4.0951, -3.439, -2.71, -1.9, -1.0, 0]

```

Figure 8: Value Function for Value Iteration

```

Up probabilistic matrix:
[0, 0, 0, 0, 0, 0]
[0.5, 1.0, 0.5, 0.5, 0.5, 0]
[0.5, 1.0, 0.5, 0.5, 0, 0]
[0.5, 1.0, 0.5, 0, 0, 0]
[0.5, 1.0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 0]
Down probabilistic matrix:
[0, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 1.0]
[0, 0, 0, 0, 0.5, 1.0]
[0, 0, 0, 0.5, 0.5, 1.0]
[0, 0, 0.5, 0.5, 0.5, 1.0]
[0, 0, 0, 0, 0, 0]
Left probabilistic matrix:
[0, 0, 1.0, 1.0, 1.0, 1.0]
[0, 0, 0.5, 0.5, 0.5, 0]
[0, 0, 0.5, 0.5, 0, 0]
[0, 0, 0.5, 0, 0, 0]
[0, 0, 0, 0, 0, 0]
[0, 0, 0, 0, 0, 0]
Right probabilistic matrix:
[1.0, 0, 0, 0, 0, 0]
[0.5, 0, 0, 0, 0, 0]
[0.5, 0, 0, 0, 0.5, 0]
[0.5, 0, 0, 0.5, 0.5, 0]
[0.5, 0, 0.5, 0.5, 0.5, 0]
[1.0, 1.0, 1.0, 1.0, 1.0, 0]
Iteration to get optimal state value:
6

```

Figure 9: Probabilistic policy matrix for Value Iteration

the grid world, since these states are far away from the terminated states. And the expected return is large when close to grid 1 and 35.

For the last two experiments meant to improve the random policy and get optimum value function and policy, we find that both methods (policy iteration and value iteration) get the optimum policy and value function. Policy iteration needs fewer rounds of improvement, but needs to evaluate the intermediate policies after each improvement; while value iteration needs more rounds of improvement, but only needs to derive the optimum policy. More time needs to be consumed when using policy iteration.

The final optimum policy is shown in Figure 10.

## 6 Conclusion

Dynamic programming avoids the problem of solving Bellman Equation directly, and in the simple grid world setting, it can achieve optimum policy in few rounds of improvement.

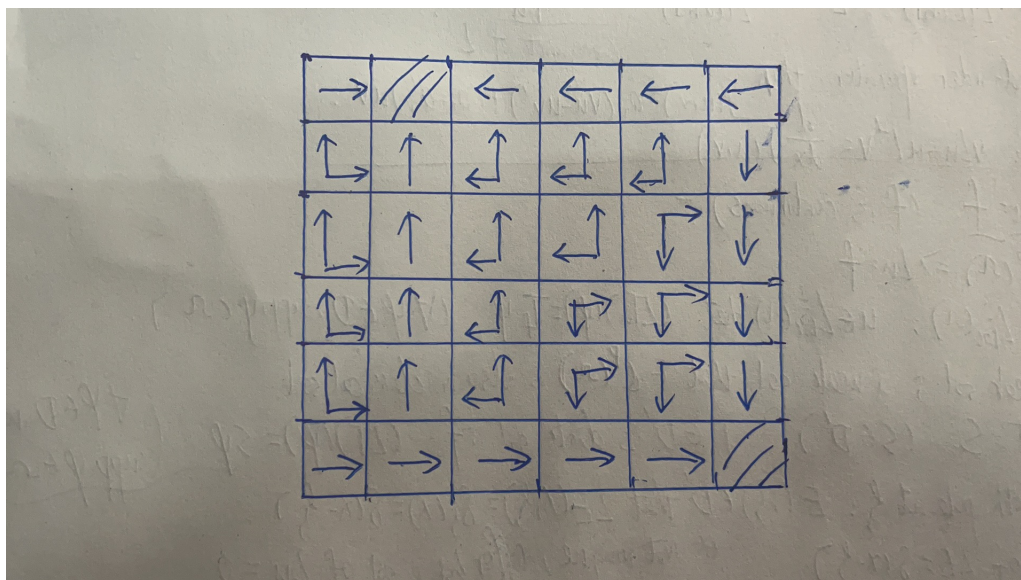


Figure 10: Optimum Policy Derived from the Experiments

## References

- [1] Prof. Junni Zou. *Reinforcement Learning, Lecture 3*.