

Name :- Brajesh Gupta

Roll No. :- 2021101095

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Machine, Data And Learning Assignment 3

Ans Value Iteration Algorithm by Hand is as follows:-

Given :- Grid (Shown Below), Step Cost = -0.04, Discount Factor ( $\gamma$ ) = 0.95  
 $P$  (going in direction of ~~action~~) = 0.7,  $P$  (going in a direction of action) = 0.15  
Reward for reaching reward state, Penalty for reaching penalty state = -1.

	Reward	Penalty
	Wall	
Start		

Initialize Utility Vectors  $dp$  and  $dp\_prev$  (Utility of cell in Previous Iteration) as  $dp = [90, 03, 50, 003, 50, 003, 50, 117]$  and  $dp\_prev = [90, 03, 50, 003, 50, 903, 50, 003]$

Iteration :-

Grid State is :-

		(H)	(L)
3	0	Reward	Penalty
2	0	0	0
1	0	Wall	0
0	0	0	0
	0	1	2

By Bellman Update Formula,

$$dp[i][j] = \text{Step Cost} + \text{Discount Factor} \times \max(dp[i-1][j], dp[i+1][j], dp[i][j-1], dp[i][j+1])$$
$$\begin{aligned} &= -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ &\quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0) \\ &= -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ &\quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0) \\ &= -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ &\quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0) \\ &= -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ &\quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0) \end{aligned}$$

If a state is invalid,  $dp\_prev[i][j]$  is considered in place of  $dp$  per of that state.

$$dp[0][0] = -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ = -0.04 \quad (Up) \quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$$

$$dp[0][1] = -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ = -0.04 \quad (Up) \quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$$

$$dp[0][2] = -0.04 + 0.95 \times \max(0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, \\ = -0.04 \quad (Up) \quad 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0, 0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0)$$

$$\begin{aligned} dp[1,0] &= -0.04 + 0.95 \times \max(0.7x0+0.15x0, 0.7x0+0.15x0+0.15x0) \\ &= -0.04 \quad (\text{Up}) \end{aligned}$$

$dp[1,1] = \text{Wall}$  (Utility Value not Updated) = 0

$$\begin{aligned} dp[1,2] &= -0.04 + 0.95 \times \max(0.7x0+0.15x0, 0.7x0+0.15x0+0.15x0) \\ &= -0.04 \quad (\text{Up}) \end{aligned}$$

$$\begin{aligned} dp[2,0] &= -0.04 + 0.95 \times \max(0.7x0+0.15x0+0.15x0, 0.7x0+0.15x0) \\ &= -0.04 \quad (\text{Up}) \end{aligned}$$

$$\begin{aligned} dp[3,1] &= -0.04 + 0.95 \times \max(0.7x1+0.15x0, 0.7x0+0.15x0+0.15x0) \\ &= -0.04 + 0.95 \times 0.7 \\ &= 0.625 \quad (\text{Up}) \end{aligned}$$

$$\begin{aligned} dp[3,2] &= -0.04 + 0.95 \times \max(0.7x(-1)+0.15x0+0.15x0, 0.7x0+0.15x0+0.15x0) \\ &= -0.04 \quad (\text{Down}) \end{aligned}$$

$$\begin{aligned} dp[3,0] &= -0.04 + 0.95 \times \max(0.7x0+0.15x0+0.15x1, 0.7x0+0.15x0+0.15x1) \\ &= -0.04 + 0.95 \times 0.7 \\ &= 0.625 \quad (\text{Right}) \end{aligned}$$

$dp[3,1] = \text{Reward}$  (Utility Value not Updated) = 1

$dp[3,2] = \text{Penalty}$  (Utility Value not Updated) = -1

Set  $dp\_prev = dp$

$$\Rightarrow dp\_prev = \left[ \begin{array}{c} \{-0.04, -0.04, -0.04\} \\ \{-0.04, 0, -0.04\} \\ \{0.625, 1, -1\} \end{array} \right]$$

Iteration 2:-

Grid State is :-

0.625	Reward (1)	Penalty (-1)
-0.04	0.625	-0.04
0.04	Wall	-0.04
-0.04	-0.04	-0.04

Using the Bellman update formula, we get.

$$\begin{aligned}
 dp[0][0] &= -0.04 + 0.95 \times \max \left( 0.7x(-0.04) + 0.15x(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15x(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15x(-0.04) + 0.15x(-0.04) \right) \\
 &= -0.04 + 0.95 \times (-0.04) \\
 &= -0.04 - \frac{0.95}{0.95} = -0.078 \quad (\text{Up})
 \end{aligned}$$

$$\begin{aligned}
 dp[0][1] &= -0.04 + 0.95 \times \max \left( 0.7x(-0.04) + 0.15(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15x(-0.04) \right) \\
 &= -0.04 + 0.95x(-0.04) \\
 &= -0.078 \quad (\text{Up})
 \end{aligned}$$

$$\begin{aligned}
 dp[0][2] &= -0.04 + 0.95 \times \max \left( 0.7(-0.04) + 0.15(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \right) \\
 &= -0.04 + 0.95(-0.04) \\
 &= -0.078 \quad (\text{Up})
 \end{aligned}$$

$$\begin{aligned}
 dp[1][0] &= -0.04 + 0.95 \times \max \left( 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \right) \\
 &= -0.04 + 0.95(-0.04) \\
 &= -0.078. \quad (\text{Up})
 \end{aligned}$$

$$dp[1][1] = \text{wall} \quad (\text{Utility Value not Updated}) = 0$$

$$\begin{aligned}
 dp[1][2] &= -0.04 + 0.95 \times \max \left( 0.7x(-0.04) + 0.15(-0.04) + 0.15x(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15(-0.04), \right. \\
 &\quad \left. 0.7x(-0.04) + 0.15(-0.04) + 0.15(-0.04) \right) \\
 &= -0.04 + 0.95(-0.04) \\
 &= -0.078 \quad (\text{Up})
 \end{aligned}$$

$$\begin{aligned}
 dp[2][0] &= -0.04 + \max \left( 0.7 \times 0.625 + 0.15 \times 0.625 + 0.15 \times (-0.04), \right. \\
 &\quad 0.7 \times 0 + 0.15 \times 0.625 + 0.15 \times (-0.04), \\
 &\quad 0.7 \times (-0.04) + 0.15 \times 0.625 + 0.15 \times (-0.04), \\
 &\quad 0.7 \times (0.625) + 0.15 \times 0.625 + 0.15 \times (-0.04) \left. \right) \\
 &= -0.04 + 0.95 \times 0.52525 \\
 &= 0.4589875 \approx 0.459. \quad (\text{Right})
 \end{aligned}$$

$$\begin{aligned}
 dp[2][1] &= -0.04 + 0.95 \times \max \left( 0.7 \times 1 + 0.15 \times (-0.04) + 0.15 \times (-0.04), \right. \\
 &\quad 0.7 \times 0 + 0.15 \times (-0.04) + 0.15 \times (-0.04), \\
 &\quad 0.7 \times (-0.04) + 0.15 \times 1 + 0.15 \times 0, \\
 &\quad 0.7 \times (-0.04) + 0.15 \times 1 + 0.15 \times 0 \left. \right) \\
 &= -0.04 + 0.95 \times 0.688 \\
 &= 0.6136 \quad (\text{Up})
 \end{aligned}$$

$$\begin{aligned}
 dp[2][2] &= -0.04 + 0.95 \times \max \left( 0.7 \times 1 + 0.15 \times 0.625 + 0.15 \times (-0.04), \right. \\
 &\quad 0.7 \times (-0.04) + 0.15 \times 0.625 + 0.15 \times (-0.04), \\
 &\quad 0.7 \times 0.625 + 0.15 \times (-1) + 0.15 \times (-0.04), \\
 &\quad 0.7 \times (-0.04) + 0.15 \times (-1) + 0.15 \times (-0.04) \left. \right) \\
 &= -0.04 + 0.95 \times 0.2815 \\
 &= 0.227425 \approx 0.227. \quad (\text{Left})
 \end{aligned}$$

$$\begin{aligned}
 dp[3][0] &= -0.04 + 0.95 \times \max \left( 0.7 \times 0.625 + 0.15 \times 1 + 0.15 \times 0.625, \right. \\
 &\quad 0.7 \times (0.04) + 0.15 \times 1 + 0.15 \times 0.625, \\
 &\quad 0.7 \times (0.625) + 0.15 \times 0.625 + 0.15 \times (-0.04), \\
 &\quad 0.7 \times 1 + 0.15 \times 0.625 + 0.15 \times (-0.04) \left. \right) \\
 &= -0.04 + 0.95 \times 0.2083625 \approx 0.2083625 \quad (\text{Right}) \\
 &= 0.2083625 \approx 0.208
 \end{aligned}$$

$$dp[3][1] = \text{Reward} \text{ (Utility value not updated)} = 1$$

$$dp[3][2] = \text{Penalty} \text{ (Utility value not updated)} = -1$$

On comparing the values obtained in first 2 iterations by running the program, we observe that values obtained are more or less the same. And minor fluctuations are caused by rounding off error in Multiplication in Python.

The Grid State is

0.208	Reward	Penalty
0.459	0.6136	0.227
-0.078	Wall	-0.078
-0.078	-0.078	-0.078