

MDL - Assignment-3

NAGA MANOTHAR

2021101128

(2)

Given,

$$\text{Step cost} = R(I, A) = 0.04$$

Probability - direction of action = 0.7

- Perpendicular to direction of action = 0.15

$$\text{discount factor} = \gamma = 0.95$$

$$\text{delta} = \beta = 0.0001$$

Initialise: $U_0(I) = 0$

$$U_0 = \begin{matrix} & \text{col}(j) & & \text{row}(i) \\ & 0 & 1 & 2 \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 0 & 1 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

reward - goal state = 1
Penalty - red state = -1

(~~assumed~~ utility of wall = 0)

Iteration 1:

$$\begin{aligned} U_{t+1}(I) &= \max_A \left(R(I, A) + \gamma \sum_J P(J|I, A) U_t(J) \right) \\ &= R(I, A) + \max_A \left(\gamma \sum_J P(J|I, A) U_t(J) \right) \end{aligned}$$

($\because R(I, A)$ is constant)

$$U_{t+1}(I) = R(I, A) + \gamma \max_A \left(\sum_J P(J|I, A) U_t(J) \right)$$

@ (γ is constant)

Iteration-1:

$$U_1[0][0]$$

=

$$\max(-0.04) + \max \begin{pmatrix} 0.95^0 (0.7^0 0 + 0.15^0 1 + 0.15^0 0), & \text{up} \\ 0.95^0 (0.7^0 1 + 0.15^0 0 + 0.15^0 0), & \text{left} \\ 0.95^0 (0.7^0 0 + 0.15^0 1 + 0.15^0 0), & \text{down} \\ 0.95^0 (0.7^0 0 + 0.15^0 0 + 0.15^0 0), & \text{right} \end{pmatrix}$$

(cycle)

$$= -0.04 + \max(0.1425, 0.665, 0.1425, 0)$$

$$= -0.04 + 0.665 =$$

$$= 0.625$$

$$U_1[0][1] = -0.04 + 0.95^0 \max(0.7^0 0, 0.7^0 1, 0.7^0 0, 0.7^0 0)$$

$$U_1[0][1] = U_0[0][1] \quad (\text{absorption state - no movement})$$

$$U_1[0][2] = U_0[0][2] = -1$$

$$U_1[1][0] = -0.04 + 0.95^0 \max \begin{pmatrix} 0.7^0 0 + 0.15^0 0 + 0.15^0 0, \\ 0.7^0 0 + 0.15^0 0 + 0.15^0 0, \\ 0.7^0 0 + 0.15^0 0 + 0.15^0 0, \\ 0.7^0 0 + 0.15^0 0 + 0.15^0 0 \end{pmatrix}$$

$$U_1[1][0] = -0.04$$

$$U_1[1][1] = -0.04 + 0.95^0 \max \begin{pmatrix} 0.7^0 1 + 0.15^0 0 + 0.15^0 0, \\ 0.7^0 0 + 0.15^0 1 + 0.15^0 0, \\ 0.7^0 0 + 0.15^0 0 + 0.15^0 0, \\ 0.7^0 0 + 0.15^0 1 + 0.15^0 0 \end{pmatrix}$$

$$= -0.04 + 0.95^0 (0.7)$$

$U[i][j]$ - represents a cell/
row col in k^{th} iteration
state
 $0 \leq i < 4$
 $0 \leq j < 3$

$$= 0.625$$

$$u_1[1][2] = -0.04 + 0.95 \cdot \max \begin{pmatrix} 0.7^x(-1) + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x(-1) + 0.15^x0, \\ 0.7^x0 + 0.15^x(0) + 0.15^x0, \\ 0.7^x0 + 0.15^x(-1) + 0.15^x0 \end{pmatrix}$$

$$= -0.04 + 0.95 \cdot \max(-0.7, -0.15, 0, -0.15)$$

$$= -0.04 + 0 =$$

$$= -0.04$$

$$u_1[2][0] = -0.04 + 0.95 \cdot \max \begin{pmatrix} 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0 \end{pmatrix}$$

$$u_1[2][0] = -0.04$$

$$u_1[2][1] = (\text{wall (utility} = 0)) = 0$$

$$u_1[2][2] = -0.04 + 0.95 \cdot \max \begin{pmatrix} 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0 \end{pmatrix}$$

$$u_1[2][2] = -0.04$$

$$u_1[3][0] = -0.04 + 0.95 \cdot \max \begin{pmatrix} 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0 \end{pmatrix}$$

$$u_1[3][0] = -0.04$$

$$u_1[3][1] = -0.04 + 0.95 \cdot \max \begin{pmatrix} 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0 \end{pmatrix}$$

$$u_1[3][1] = -0.04$$

$$u_1[3][2] = -0.04 + 0.95 \cdot \max \begin{pmatrix} 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0, \\ 0.7^x0 + 0.15^x0 + 0.15^x0 \end{pmatrix}$$

$$u_1[3][2] = -0.04$$

$$U = \begin{matrix} & 0 & 1 & 2 \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 0.625 & 1 & -1 \\ -0.04 & 0.625 & -0.04 \\ -0.04 & 0 & -0.04 \\ -0.04 & -0.04 & -0.04 \end{bmatrix} \end{matrix}$$

Here as the agent ~~did~~ ^{did} up but reached same cell $U[0][0]$ thus $U[0][0]$

Iteration-2's

$$U_2[0][0] = -0.04 + 0.95 \times \max \begin{pmatrix} 0.7 \times 0.625 + 0.15 \times 1 + 0.15 \times 0, \\ 0.7 \times 1 + 0.15 \times 0.625 + 0.15 \times (-0.04), \\ 0.7 \times (-0.04) + 0.15 \times 0.625 + 0.15 \times 1, \\ 0.7 \times (0.625) + 0.15 \times 0.625 + 0.15 \times (-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \times \max(0.5875, 0.78775, 0.21575, 0.52525) \\ = -0.04 + 0.95 \times 0.78775 \\ = 0.7083625$$

$$U_2[0][1] = U_1[0][1] = 1 \\ U_2[0][2] = U_1[0][2] = -1$$

$$U_2[1][0] = -0.04 + 0.95 \times \max \begin{pmatrix} 0.7 \times 0.625 + 0.15 \times 0.625 + 0.15 \times (-0.04), \\ 0.7 \times 0.625 + 0.15 \times 0.625 + 0.15 \times (-0.04), \\ 0.7 \times (-0.04) + 0.15 \times 0.625 + 0.15 \times (-0.04), \\ 0.7 \times (-0.04) + 0.15 \times 0.625 + 0.15 \times (-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \times (0.52525)$$

$$U_2[1][0] = 0.4589875$$

$$V_2[1][1] = -0.04 + 0.95 \max \begin{pmatrix} 0.7 \cdot 1 + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (1) + 0.15 \cdot (0) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (1) + 0.15 \cdot (0) \end{pmatrix}$$

$$= -0.04 + 0.95 \cdot 0.688$$

$$= 0.6136$$

$$V_2[1][2] = -0.04 + 0.95 \max \begin{pmatrix} 0.7 \cdot (-1) + 0.15 \cdot (0.625) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (-1) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (0.625) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (0.625) + 0.15 \cdot (-1) + 0.15 \cdot (-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \cdot 0.2815$$

$$= 0.227425$$

$$V_2[2][0] = -0.04 + 0.95 \max \begin{pmatrix} 0.7 \cdot (0.04) + 0.15 \cdot (0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \cdot (-0.04)$$

$$= -0.078$$

$$\left((0.7 + 0.15 + 0.15) \cdot (-0.04) = 1 \cdot (-0.04) = -0.04 \right)$$

$$V_2[2][1] (\text{wall}) = 0$$

$$V_2[2][2] = -0.04 + 0.95 \max \begin{pmatrix} 0.7 \cdot (-0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \\ 0.7 \cdot (-0.04) + 0.15 \cdot (-0.04) + 0.15 \cdot (-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \cdot (-0.04)$$

$$= -0.078$$

$$U_2[3][0] = -0.04 + 0.95 \times \max \begin{pmatrix} 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \times (-0.04)$$

$$= -0.078$$

$$U_2[3][1] = -0.04 + 0.95 \times \max \begin{pmatrix} 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \times (-0.04)$$

$$= -0.078$$

$$U_2[3][2] = -0.04 + 0.95 \times \max \begin{pmatrix} 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \\ 0.7(-0.04) + 0.15(-0.04) + 0.15(-0.04) \end{pmatrix}$$

$$= -0.04 + 0.95 \times (-0.04)$$

$$= -0.078$$

$$U_2 = \begin{matrix} & 0 & 1 & 2 \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 0.7083625 \\ 0.4589875 \\ -0.078 \\ -0.078 \end{bmatrix} & \begin{bmatrix} 1 \\ 0.6136 \\ 0 \\ -0.078 \end{bmatrix} & \begin{bmatrix} -1 \\ 0.227425 \\ -0.078 \\ -0.078 \end{bmatrix} \end{matrix}$$

Thus, the values of U_1, U_2 match with the output of the code.