

MDL - Assignment-3

ROMICA RAISINGHANI

Implementing Value Iteration for MDP 2021101053

According to value iteration algorithm,

Step 1:

We initialize the utilities of each state to zero.

$$i.e., U_0(I) = 0$$

$U_0(I) \rightarrow$ utility of state I at time $t=0$.

Step 2: Iterate:

According to Bellman Update Equation:
$$U_{t+1}(I) = \max_A \left[R(I, A) + \gamma \sum_J P(J|I, A) \times U_t(J) \right]$$

Prob. of going in the direction of action } = 0.7

where,

$R(I, A) \rightarrow$ immediate reward

OR

Reward on taking action A on state I (step cost)

Prob. of going perpendicular to direction of action } = 0.15

Prob. of going opposite to the direction of action } = 0

$P(J|I, A) \rightarrow$ transition probability on reaching state J from state I upon action A .

$U_t(I) \rightarrow$ utility value of state I in the previous iteration

$\gamma \rightarrow$ discount factor

We will continue iterating until the values $\{U_t, U_{t+1}\}$ converge.

At the end of iteration, we calculate optimal policy:

Policy $(I) =$

$$\arg \max_A \left[R(I, A) + \gamma \sum_J P(J|I, A) \times U_{t+1}(J) \right]$$

$\rightarrow (0,0)$

	0	1	2
0		Goal	Penalty
1			
2			
3			

According to algorithm:

Initial utilities of all states
 $= u[i][j] = 0$

For (0,0):

North \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

West/East \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0]$$

max of all actions = -0.04

$$u[0][0] = -0.04$$

For (0,1): Goal state

$$u[0][1] = 1$$

For (0,2): Penalty state

$$u[0][2] = -1$$

For (1,0):

North \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

West/East \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

max of all directions = -0.04

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

For (1,1):

North \rightarrow

$$-0.04 + 0.95[0.7 \times 1 + 0.15 \times 0 + 0.15 \times 0] = -0.04 + 0.665 = 0.625$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

West/East \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 1 + 0.15 \times 0] = -0.04 + 0.95 \times 0.15 = 0.1025$$

$$\Rightarrow u[1][1] = 0.625$$

For (1,2):

North:

$$-0.04 + 0.95[0.7 \times -1 + 0.15 \times 0 + 0.15 \times 0] = -0.04 - 0.95 \times 0.7 = -0.705$$

South:

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] = -0.04$$

West/East:

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times (-1) + 0.15 \times 0]$$

$$= -0.04 - 0.95 \times 0.15$$

$$= -0.1825$$

max of all directions = -0.04

$$u[1,2] = -0.04$$

For (2,0):

North →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

South →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

West/East →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

max of all directions = -0.04

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

For (2,k): wall

$$u[2][k] = 0$$

For (2,2):

North →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

South →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

West/East →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

max of all directions = -0.04

$$u[2][2] = -0.04$$

For (3,0):

North →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

South →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

East/West →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

max of all directions = -0.04

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

For (3,1):

North →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

South →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

West/East →

$$-0.04 + 0.95 \left[\begin{matrix} 0.7 \times 0 + 0.15 \times 0 \\ + 0.15 \times 0 \end{matrix} \right]$$

$$= -0.04$$

max of all dir's = -0.04

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

For (3,2):

North →

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] \\ = -0.04$$

South →

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] \\ = -0.04$$

West/East →

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 0] \\ = -0.04$$

max of all dirⁿs = -0.04

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

ITERATION 2:

For (0,0):

North →

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times 1] \\ = -0.04 + 0.95 \times 0.15 \\ = 0.1025$$

South →

$$-0.04 + 0.95[0.7 \times (-0.04) + 0.15 \times 0 + 0.15 \times 1] \\ = -0.04 - 0.0666$$

West →

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15 \times (-0.04)] \\ = -0.0457$$

East →

$$-0.04 + 0.95[0.7 \times 1 + 0.15 \times 0 + 0.15 \times (-0.04)] \\ = -0.04 + 0.95[0.7 - 0.006] \\ = 0.6193$$

max of all directions = 0.6193
 $U[0][0] = 0.6193$

For (0,1): Goal state
 $U[0][1] = 1$

For (0,2): Penalty state
 $U[0][2] = -1$

For (1,0):

North →

$$-0.04 + 0.95[0.7 \times (-0.04) + 0.15 \times 0 + 0.15 \times 0.625] \\ = -0.04 + 0.95[-0.028 + 0.09375] \\ = -0.04 + 0.95[0.06575] \\ = 0.0224625$$

South →

$$-0.04 + 0.95[0.7 \times (-0.04) + 0.15 \times 0 + 0.15 \times 0.625] \\ = 0.0224625$$

West →

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times (-0.04) + 0.15 \times (-0.04)] \\ = -0.04 - 0.0114 \\ = -0.0514$$

East \rightarrow

$$-0.04 + 0.95[0.7 \times 0.625 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= -0.04 + 0.95[0.4375 - 0.012]$$

$$= 0.364225$$

max of all directions = 0.364225

$$u[1][0] = 0.364225$$

For (1,1):

North \rightarrow

$$-0.04 + 0.95[0.7 \times 1 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= 0.62386$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= -0.0514$$

West \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0.15 \times 1 + 0.15 \times 0]$$

$$= 0.0759$$

East \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0.15 \times 1 + 0.15 \times 0]$$

$$= 0.0759$$

max of all directions =

$$u[1][1] = 0.62386$$

For (1,2):

North \rightarrow

$$-0.04 + 0.95[0.7(-1) + 0.15 \times 0 + 0.15(0.625)]$$

$$= -0.6159375$$

South \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0.15 \times 0 + 0.15 \times 0.625]$$

$$= 0.0224625$$

West \rightarrow

$$-0.04 + 0.95[0.7(0.625) + 0.15(-1) + 0.15(-0.04)]$$

$$= 0.227425$$

East \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15(-1) + 0.15(-0.04)]$$

$$= -0.1882$$

max of all directions = 0.227425

$$u[1][2] = 0.227425$$

For (2,0):

North \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0]$$

$$= -0.0666$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times (-0.04) + 0]$$

$$= -0.0666$$

West/East \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= -0.0514$$

max of all directions = -0.0514

$$u[2][0] = -0.0514$$

For (2,1): wall

$$u[2][1] = 0$$

For (2,2):

North \rightarrow

$$-0.04 + 0.95[0.7 \times (-0.04) + 0] = -0.0666$$

South \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0] = -0.0666$$

West/East \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= -0.0514$$

max of all directions = -0.0514

$$u[2][2] = -0.0514$$

For (3,0):

North \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0 + 0.15(-0.04)]$$

$$= -0.0723$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15(-0.04)]$$

$$= -0.0457$$

West \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15 \times 0 + 0.15(-0.04)]$$

$$= -0.0457$$

East \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0.15(-0.04) + 0.15 \times 0]$$

$$= -0.0723$$

max of all directions = -0.0457

$$u[3][0] = -0.0457$$

For (3,1):

North \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= -0.0514$$

South \rightarrow

$$-0.04 + 0.95[0.7 \times 0 + 0.15(-0.04) + 0.15(-0.04)]$$

$$= -0.0514$$

West/East \rightarrow

$$-0.04 + 0.95[0.7(-0.04) + 0 + 0]$$

$$= -0.0666$$

max of all directions = -0.0514

$$u[3][1] = -0.0514$$

For (3,2):

North \rightarrow

$$-0.04 + 0.95 [0.7(-0.04) + 0.15(-0.04) + 0] = -0.0723$$

South \rightarrow

$$-0.04 + 0.95 [0.7 \times 0 + 0.15(-0.04) + 0] = -0.0457$$

Fast West \rightarrow

$$-0.04 + 0.95 [0.7(-0.04) + 0.15(-0.04) + 0] = -0.0457$$

$$= -0.0723$$

Fast \rightarrow

$$-0.04 + 0.95 [0.7 \times 0 + 0.15(-0.04) + 0] = -0.0457$$

max of all directions = -0.0457

$$u[3][2] = -0.0457$$

Matrix after Iteration 1:

$$\begin{bmatrix} 0.62499 & 1 & -1 \\ -0.04 & 0.62499 & -0.04 \\ -0.04 & 0 & -0.04 \\ -0.04 & -0.04 & -0.04 \end{bmatrix}$$

Matrix after Iteration 2:

$$\begin{bmatrix} 0.10836 & 1 & -1 \\ 0.4589 & 0.6135 & 0.2274 \\ -0.0779 & 0 & -0.0779 \\ -0.0779 & -0.0779 & -0.0779 \end{bmatrix}$$

Both iteration values matches with the code in python.

Hence, the algorithm works by manual calculation as well as on code.

There are 22 iterations taken for the algorithm to converge.