

2019101056

Aryan Jain

Page:

Date: / /

MDL Assgn 2 Part 1

Ques 1 At either of the state (A/B/C/R) we can have 4 possible moves which are Up (U), Down (D), Left (L), Right (R).

To(s') \Rightarrow From(s) \Downarrow	A				B				C				R			
	L	R	U	D	L	R	U	D	L	R	U	D	L	R	U	D
A	0	0.2	0.2	0	0	0.8	0	0	0	0.8	0	0	0	0	0	0
B	0.8	0	0	0	0.2	0	0.2	0	0	0	0	0	0	0.8	0	0
C	0	0	0	0.8	0	0	0	0	0	0.75	0.2	0	0.25	0	0	0

Here each value represents $P(s'|s, a)$ i.e. probability of reaching a state s' from ~~initial~~ given initial state s and given action a .

Ques 2 Going from $A \rightarrow R$ has different possibilities like

$$A \rightarrow B \rightarrow R$$

$$A \rightarrow C \rightarrow R$$

$$A \rightarrow B \rightarrow C \rightarrow R$$

$$A \rightarrow C \rightarrow B \rightarrow R$$

} these will add on some extra cost than above two.

\therefore Best path will be chosen from $A \rightarrow B \rightarrow R$ and $A \rightarrow C \rightarrow R$

Given $P(B \rightarrow R) = 0.8$ Step Cost $(B \rightarrow R) = -4$
 $P(C \rightarrow R) = 0.25$ Step Cost $(C \rightarrow R) = -3$

$$\therefore \text{Expected utility } (B \rightarrow R) = 0.8 \times (-4) = -3.2$$

$$\text{Expected utility } (C \rightarrow R) = 0.25 \times (-3) = -0.75$$

\therefore ~~Optimal solⁿs follows MEU principle~~

By MEU principle (Maximum Expected Utility)
we can say $A \rightarrow C \rightarrow R$ is the optimal path.

Ques 3 roll-number = 2019101056

$$R = \text{Arr}[2019101056 \cdot 15] = \text{Arr}[1] = 9$$

$$\rightarrow \text{Reward Value } (R) = 9$$

To perform the value iteration algorithm we need to apply Bellman equation viz.

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

Also for $t=0$ $U_0(A)=0$ $U_0(B)=0$
 $U_0(C)=0$ $U_0(D)=9$

0	9
0	0

Let's call this as Destination D

$t=0$ 1st iteration a.) For state 'A': R, U are only possible actions

$$U_1(A) = \max [R(A, R) + \gamma [P(B|A, R) \cdot U_0(B) + P(A|A, R) \cdot U_0(A)]$$

$$R(A, U) + \gamma [P(C|A, U) \cdot U_0(C) + P(D|A, U) \cdot U_0(D)]$$

$$U_1(A) = \max [-1 + 0.2 [0.8 \times 0 + 0.2 \times 0]]$$

$$-1 + 0.2 [0.8 \times 0 + 0.2 \times 0]$$

$$\therefore U_1(A) = -1$$

b.) For state 'B': L, U are 2 possible actions.

$$U_1(B) = \max [R(B, L) + \gamma [P(A|B, L) \cdot U_0(A) + P(B|B, L) \cdot U_0(B)]]$$

$$R(B, U) + \gamma [P(D|B, U) \cdot U_0(D) + P(B|B, U) \cdot U_0(B)]$$

↓ Destination
 To calculate $R(B, U)$: The reward associated with Action = 'U' and initial state = 'B' is also dependent on the destination state.

$$\therefore R(B, U) = \sum_J R(B, U, J) \cdot P(J|B, U)$$

$$\therefore R(B, U) = -4 \times 0.8 + (-1) \times 0.2 = -3.4$$

$$\therefore U_1(B) = \max \begin{cases} -1 + 0.2[0.8 \times 0 + 0.2 \times 0] = -1 \\ -3.4 + 0.2[0.8 \times 9 + 0.2 \times 0] = -1.96 \end{cases}$$

$$\therefore \boxed{U_1(B) = -1}$$

c) For state 'C': R, D are 2 possible actions.

$$U_1(C) = \max \begin{cases} R(C, R) + \gamma [P(D|C, R) \cdot U_0(D) + P(C, C, R) \cdot U_0(C)] \\ R(C, D) + \gamma [P(A|C, R) \cdot U_0(A) + P(C, C, D) \cdot U_0(C)] \end{cases}$$

$$R(C, R) = -3 \times 0.25 + (-1) \times 0.75 = -1.5$$

$$U_1(C) = \max \begin{cases} -1.5 + 0.2[0.25 \times 9 + 0.75 \times 0] = -1.05 \\ -1 + 0.2[0.8 \times 0 + 0.2 \times 0] = -1 \end{cases}$$

$$\therefore \boxed{U_1(C) = -1}$$

\therefore New utilities values are

-1	10
-1	-1

$$\max_x [U_{t+1}(x) - U_t(x)] \quad \text{Max. difference} = 1 > 0.01 \quad (\delta = 0.01)$$

\Rightarrow ~~Not~~ Not converged

$t=1$
2nd iteration

By taking the formula from last iteration and substituting $U_0(x)$ with $U_1(x)$, we have :-

a) For state 'A'

$$U_2(A) = \max \begin{cases} -1 + 0.2[0.8 \times (-1) + 0.2 \times (-1)] = -1.2 \\ -1 + 0.2[0.8 \times (-1) + 0.2 \times (-1)] = -1.2 \end{cases}$$

$$\Rightarrow \boxed{U_2(A) = -1.2} \quad \boxed{U_2(A) = -1.2}$$

b) For state 'B'

$$U_2(B) \begin{matrix} \text{max} \\ \text{min} \end{matrix} \begin{matrix} -1 + 0.2[0.8 \times (-1) + 0.2 \times (-1)] = \cancel{-1.02} \\ -3.4 + 0.2[0.8 \times 9 + 0.2 \times (-1)] = -2.00 \end{matrix}$$

$\therefore \cancel{U_2(B) = -1.02} \quad \boxed{U_2(B) = -1.2}$

c) For state 'C'

$$U_2(C) \begin{matrix} \text{max} \\ \text{min} \end{matrix} \begin{matrix} -1.5 + 0.2[0.25 \times 9 + 0.75 \times (-1)] = \cancel{0} \\ -1 + 0.2[0.8 \times (-1) + 0.2 \times (-1)] = \cancel{-1.02} \end{matrix}$$

$\therefore \boxed{\cancel{U_2(C) = 0}} \quad \boxed{\cancel{U_2(C) = -1.02}} \quad \boxed{U_2(C) = -1.2}$

\therefore New utilities values are

1.02	9	-1.2	9
-1.02	-1.02	-1.2	-1.2

Max. difference = $\cancel{0.02}^{0.2} > 0.01 \quad (\delta = 0.01)$

\Rightarrow Not converged

∴ New utilities values are

-1.203	9
-1.204	-1.204

∴ Max. difference =

$t=2$

3rd iteration

a) For state 'A'

$$U_2(A) \begin{cases} -1 + 0.2 [0.8 \times (-1.2) + 0.2 \times (-1.2)] \\ \text{max.} \quad -1 + 0.2 [0.8 \times (-1.2) + 0.2 \times (-1.2)] \end{cases}$$

$$\therefore \boxed{U_2(A) = -1.24}$$

b) For state 'B'

$$U_2(B) \begin{cases} -1 + 0.2 [0.8 \times (-1.2) + 0.2 \times (-1.2)] \\ \text{max.} \quad -3.4 + 0.2 [0.8 \times 9 + 0.2 \times (-1.2)] \end{cases}$$

$$\therefore \boxed{U_2(B) = -1.24}$$

c) For state 'C'

$$U_2(C) \begin{cases} -1.5 + 0.2 [0.25 \times 9 + 0.75 \times (-1.2)] \\ \text{max.} \quad -1 + 0.2 [0.8 \times (-1.2) + 0.2 \times (-1.2)] \end{cases}$$

$$\therefore \boxed{U_2(C) = -1.23}$$

∴ New utilities values are

-1.23	9
-1.24	-1.24

∴ Max. difference = $0.04 > 0.01$ ($\delta = 0.01$)
⇒ Not converged

$t=3$

4th iteration

a) For state 'A'

$$U_4(A) \begin{cases} -1 + 0.2 [0.8 \times (-1.24) + 0.2 \times (-1.24)] \\ \text{max.} \quad -1 + 0.2 [0.8 \times (-1.23) + 0.2 \times (-1.24)] \end{cases}$$

$$\therefore \boxed{U_4(A) = -1.2464}$$

b) For state 'B'

$$U_4(B) \begin{cases} -1 + 0.2 [0.25 \times 9 + 0.75 \times (-1.24)] \\ \text{max. } -3.4 + 0.2 [0.8 \times 9 + 0.2 \times (-1.24)] \end{cases}$$

$$\therefore U_4(B) = -0.736$$

b) For state 'B'

$$U_4(B) \begin{cases} -1 + 0.2 [0.8 \times (-1.24) + 0.2 \times (-1.24)] \\ \text{max. } -3.4 + 0.2 [0.8 \times 9 + 0.2 \times (-1.24)] \end{cases}$$

$$U_4(B) = -1.248$$

c) For state 'c'

$$U_4(c) \begin{cases} -1.5 + 0.2 [0.25 \times 9 + 0.75 \times (-1.23)] \\ \text{max. } -1 + 0.2 [0.8 \times (-1.24) + 0.2 \times (-1.23)] \end{cases}$$

$$\therefore U_4(c) = -1.2345$$

\therefore New utility values are

-1.2345	9
-1.2464	-1.248

\therefore Max. difference = 0.0064 < 0.01
 \Rightarrow Converged (At 4th iteration)

Ans 4 Now, we run the algorithm, one last time to determine the policy.

Utilities were:

$$U_5(A) \begin{cases} -1 + 0.2 [0.8 \times (-1.248) + 0.2 \times (-1.2464)] \\ \text{max. } -1 + 0.2 [0.8 \times (-1.2345) + 0.2 \times (-1.2464)] \end{cases}$$

Action

① \leftarrow Right

② \leftarrow UP

Clearly $U_5(A) > U_5(B)$

\Rightarrow From A you should take an UP action

② > ①

(By comparison)

$$① = -1.249536 \quad ② = -1.247376$$

$U_5(B) \rightarrow$ calculating it will be useless.

$$U_5(c) \begin{cases} -1.5 + 0.2 [0.25 \times 9 + 0.75 \times (-1.2345)] \\ -1 + 0.2 [0.8 \times (-1.2464) + 0.2 \times (-1.2345)] \end{cases}$$

$$U_5(c) \begin{cases} -1.235175 \leftarrow \text{Right} \\ -1.248804 \leftarrow \text{Down} \end{cases}$$

$$U_5(c) = -1.235175$$

\Rightarrow From c u should move Right

\Rightarrow A \rightarrow c \rightarrow R is optimal path
 (which we also guessed in Ques. 2)

Ans. 5 As the reward value increases the optimal path $A \rightarrow C \rightarrow R$ starts changing to $A \rightarrow B \rightarrow R$. This is because now we have a high reward so the algorithm will tend to more risk seeking methodology.
(Means more risk can be taken).