

MIHIR BANI , 20189113003

Transition table

Let	$S \rightarrow$ State (current)	$R \rightarrow$ Right
	$S' \rightarrow$ next state	$U \rightarrow$ up
	$A \rightarrow$ Action	$D \rightarrow$ down
	$P \rightarrow$ Probability	
	$T \rightarrow$ terminal state	

Table : (without reward R)

<u>S</u>	<u>A</u>	<u>S'</u>	<u>P</u>
A	R	B	0.8
A	R	A	0.2
A	U	C	0.8
A	U	A	0.2
B	L	A	0.8
B	L	B	0.2
B	U	T	0.8
B	U	B	0.2
C	R	T	0.25
C	R	C	0.75
C	D	A	0.8
C	D	C	0.2

Q2. I think the best path can be $A \rightarrow B \rightarrow T$ as it is equally probable to reach B or C from A. But reaching T from C is less probable than reaching T from B. ($0.25 < 0.8$)
 $C \rightarrow T$ $B \rightarrow T$

Q.3 Reward = 16, $\gamma = 0.2$, $\delta = 0.01$

C	0	16	T
A	0	0	B

Initial state

I 1st iteration

$$\rightarrow A: R \rightarrow (0.8)(-1 + 0.2 \times 0) + 0.2(-1 + 0.2 \times 0) = -1$$

$$U \rightarrow 0.8(-1 + 0.2 \times 0) + 0.2(-1 + 0.2 \times 0) = -1$$

$$A = \max(-1, -1) = -1$$

$$\rightarrow B: L \rightarrow 0.8(-1 + 0.2 \times 0) + 0.2(-1 + 0.2 \times 0) = -1$$

$$U \rightarrow 0.2(-1 + 0.2 \times 0) + 0.8(-4 + 0.2 \times 16) = -0.84$$

$$B = \max(-1, -0.84) = -0.84$$

$$\rightarrow C: R \rightarrow 0.75(-1 + 0.2 \times 0) + 0.25(-3 + 0.2 \times 16) = -0.7$$

$$D \rightarrow 0.8(-1 + 0.2 \times 0) + 0.2(-1 + 0.2 \times 0) = -1$$

$$C = \max(-0.7, -1) = -0.7$$

-0.84	16
-1	-0.7

I 2nd iteration

$$\rightarrow A: U \rightarrow 0.2(-1 + 0.2(-1)) + 0.8(-1 + 0.2(-0.7)) \\ = -1.152$$

$$R \rightarrow 0.2(-1 + 0.2(-1)) + 0.8(-1 + 0.2(-0.84)) \\ = -1.1744$$

$$A = \max(-1.152, -1.1744) \\ = -1.152$$

$$\rightarrow B: L \rightarrow 0.8(-1 + 0.2(-1)) + 0.2(-1 + 0.2(-0.84)) \\ = -1.1936$$

$$U \rightarrow 0.2(-1 + 0.2(-0.84)) + 0.8(-1 + 0.2(16)) \\ = -0.8736$$

$$B = -0.8736$$

$$\rightarrow C: R \rightarrow 0.75(-1 + 0.2(-0.7)) + 0.25(-3 + 0.2(16)) \\ = -0.805$$

$$D \rightarrow 0.8(-1 + 0.2(-1)) + 0.2(-1 + 0.2(-0.7)) \\ = -1.188$$

$$\Rightarrow C = -0.805$$

-0.805	16
-1.152	-0.8736

III 3rd iteration

$$\rightarrow A: U \rightarrow 0.2(-1 + 0.2(-1.152)) + 0.8(-1 + 0.2(-0.805)) \\ = -1.17488$$

$$R \rightarrow 0.2(-1 + 0.2(-1.152)) + 0.8(-1 + 0.2(-0.8736)) \\ = -1.18586$$

$$A = \max(-1.17488, -1.18586) \\ = -1.17488$$

$$\rightarrow B: L \rightarrow 0.8(-1 + 0.2(-1.152)) + 0.2(-1 + 0.2(-0.8736))$$

$$= -1.21926$$

$$U \rightarrow 0.2(-1 + 0.2(-0.8736)) + 0.8(-4 + 0.2(16))$$

$$= -0.874944$$

$$B = \max(-1.21926, -0.874944)$$

$$B = -0.874944$$

$$\rightarrow C: R \rightarrow 0.75(-1 + 0.2(-0.805)) + 0.25(-3 + 0.2(16))$$

$$= -0.82075$$

$$D \rightarrow 0.8(-1 + 0.2(-1.152)) + 0.2(-1 + 0.2(-0.805))$$

$$= -1.21652$$

$$C = \max(-0.82075, -1.21652)$$

$$C = -0.82075$$

-0.82	16
-1.1748	-0.875

IV 4th iteration

$$\rightarrow A: U \rightarrow 0.2(-1 + 0.2(-1.17488)) + 0.8(-1 + 0.2(-0.82075))$$

$$= -1.17832$$

$$R \rightarrow 0.2(-1 + 0.2(-1.17488)) + 0.8(-1 + 0.2(-0.874944))$$

$$= -1.18699$$

$$A = \max(-1.17832, -1.18699)$$

$$= -1.17832$$

$$\rightarrow B: L \rightarrow 0.8(-1 + 0.2(-1.17488)) + 0.2(-1 + 0.2(-0.874944))$$

$$= -1.22298$$

$$U \rightarrow 0.2(-1 + 0.2(-0.874944)) + 0.8(-4 + 0.2(16))$$

$$= -0.874998$$

$$B = \max(-1.22298, -0.874998)$$

$$B = -0.874998$$

$$C: R \rightarrow 0.75(-1 + 0.2(-0.82075)) + 0.25(-3 + 0.2(16))$$

$$= -0.823112$$

$$D \rightarrow 0.8(-1) + 0.2(-1.17488) + 0.2(-1 + 0.2(-0.82075))$$

$$= -1.22081$$

$$D = \max(-0.823112, -1.22081)$$

$$D = -0.823112$$

-0.823112	16
-1.17832	-0.874998

Summary

0	16	→	-0.84	16	↔	-0.805	16
0	0		-1	-0.7		-1.152	-0.8736

-0.823	16
-1.178	-0.875

-0.82	16
-1.1748	-0.875

Q4. The optimal path after value iteration is

$A \rightarrow C \rightarrow T$

and this is different than what I answered in Q1 ($A \rightarrow B \rightarrow T$)

Q4. The optimal path found after the VI algo is

$A \rightarrow C \rightarrow T$, and this is different from my initial answer of $A \rightarrow B \rightarrow T$.

It was because of the values of step cost and maybe because of the value of discounting factor.

Q.5 The part in the path from A to T is based on the transitions from $B \rightarrow T$ and $C \rightarrow T$. As given the probability of $B \rightarrow T$ is high (0.8) but the step cost is also high (-4) than the prob. of $C \rightarrow T$, (0.25) and the step cost (-3) is lower. Thus we can say that when the reward of final state is below a threshold, the preferred state will be the one that minimizes the step cost (magnitude) so it will try to conserve its total value. But when the final reward is high, then we can risk the high step-cost of $B \rightarrow T$, because of greater probability there is now much of larger chance of increasing the value. Whereas $C \rightarrow T$ will have very low probability and the benefit of high Reward will be wasted.

For Reward = 16, optimal path = $A \rightarrow C \rightarrow T$

$\exists x$, s.t. $x > 16$

and Reward $> x$

optimal path will be $A \rightarrow B \rightarrow T$