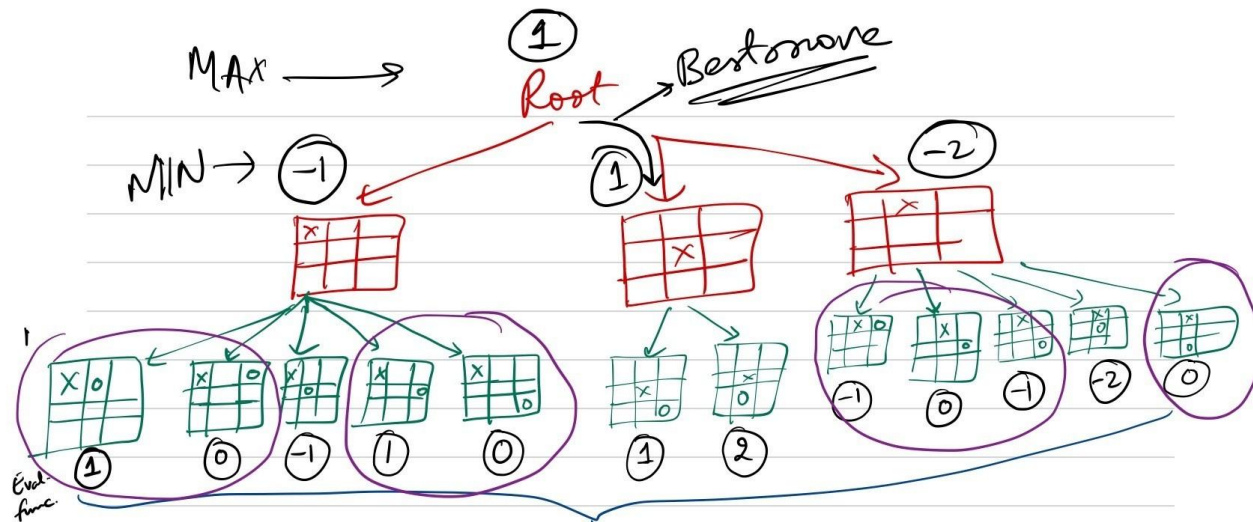


Q1



1. Drawing the above tree
2. Evaluation at depth 2
3. MIN - MAX values
4. Pruning the circled nodes: Optimal ordering is the center-X branch to be the first branch & then the other branches with the first child being less than the minimum of the center-X branch

## Q2. a.

### ANSWER

$$x > 8$$

### EXPLANATION

Depending on the value of  $x$ , there can be three different values for taking Action 1:

If  $x < -8$ , then Action 1 results in  $-8$ ;

If  $-8 \leq x \leq 11$  then Action 1 results in  $x$ ;

If  $x > 11$  then Action 1 results in  $11$ .

Action 2 always results in a utility of  $8$

Hence Action 1 is optimal for Player 1 if  $x > 8$

## b.

### ANSWER

$$x > 9$$

### EXPLANATION

Action 2 gives Player 1 a utility of  $10$ , so the average of  $x$  and  $11$  must be greater than  $10$ .  $(x + 11)/2 > 10 \rightarrow x > 9$

## c.

**ANSWER**

No

**EXPLANATION**

The minimax value can never be strictly greater than the expectimax value for the same tree because in minimax Player 2 always chooses the worst possible move for Player 1, while in expectimax, those same nodes average that value with other higher values. Thus, the utility at a node under expectimax is always at least as high as the utility of the same node under minimax.

**Q3.**

3.1 Transition Matrix:

<u>for action a:</u>				<u>for action b</u>			
	I	II	III		I	II	III
I	0.2	0.8	0	I	0.9	0	0.1
II	0.8	0.2	0	II	0	0.9	0.1
III	0	0	0	III	0	0	0

3.2.  $|A^S| = 2^2 = 4$

Optimal policy: if in 1 do b  
if in 2 do a }  $\rightarrow$  Intuition: Maximize Reward

3.3  $V_0^I = -1$  |  $V_1^I = -1 + \max((0.8)(-2) + (0.2)(-1), (0.9)(-1) + (0.1)(0)) = -1.9$   
 $V_0^{II} = -2$  |  $V_2^{II} = -2 + \max((0.8)(-1) + (0.2)(-2), (0.9)(-2) + (0.1)(0)) = -3.2$

### 3.4 → Initial Policy Eval [b, b]

$$V^I = -1 + 0.9V^I + 0.1V^{III} \Rightarrow V^I = -10$$

$$V^{II} = -2 + 0.9V^{II} + 0.1V^{III} \Rightarrow V^{II} = -20$$

$$V^{III} = 0$$

#### → Evaluation

Can be done with  $\text{argmax}(\sum T \times V)$  or  $R + \sum T \cdot V$

•  $R + \sum T \cdot V$

$$V^I = -19 \text{ (for action a)} < -10 \text{ (for action b)} \text{ [Don't change policy]}$$

$$V^{II} = -14 \text{ (for action a)} > -20 \text{ (for action b)} \text{ [Change policy to b]}$$

•  $\text{argmax}(\sum T \times V)$

$$\begin{aligned} \pi_{V^I} &= \text{argmax} \{ (0.9)^b(-10) + (0.1)^a(0), (0.8)^a(-20) + (0.2)^b(-10) \} \\ &= \text{argmax} \{ -9, -18 \} = b \end{aligned}$$

$$\begin{aligned} \pi_{V^{II}} &= \text{argmax} \{ 0.9^b(-20) + (0.1)^a(0), (0.8)^a(-10) + (0.2)^b(-20) \} \\ &= \text{argmax} \{ -18, -12 \} = a \end{aligned}$$

### → Policy Eval [b, a]

$$V^I = -1 + 0.9V^I + 0.1V^{III} = -10$$

$$V^{II} = -2 + 0.8V^I + 0.2V^{II} \quad V^{II} = -10/0.8 = -12.5$$

$$V^{III} = 0$$

#### → Evaluation

•  $R + \sum T \cdot V$

$$V^I = -1 + (0.8)(-12.5) + (0.2)(-10) = -13 \text{ (for action a)} < -10 \text{ (for action b)} \text{ [No Change]}$$

$$V^{II} = -2 + (0.9)(-12.5) + (0.1)(0) = -13.25 \text{ (for action b)} < -12.5 \text{ [for action a]}$$

•  $\text{argmax}(\sum T \times V)$

$$\begin{aligned} \pi_{V^I} &= \text{argmax} \{ (0.8)^a(-12.5) + (0.2)^b(-10), (0.9)^b(-10) + (0.1)^a(0) \} \\ &= \text{argmax} \{ -12, -9 \} = b \end{aligned}$$

$$\begin{aligned} \pi_{V^{II}} &= \text{argmax} \{ (0.8)^a(-10) + (0.2)^b(-12.5), (0.9)^b(-12.5) + (0.1)^a(0) \} = \{ \\ &= \text{argmax} \{ -10.5, -11.25 \} = a \end{aligned}$$

No Change in Policy Terminate

3.5. <-10, -12.5, 0>

### 3.6. Policy evaluation equations become unsolvable

$$V(1) = -1 + 0.8*V(2) + 0.2*V(1)$$

$$V(2) = -2 + 0.8*V(1) + 0.2*V(2)$$

Discounting helps in making these equations solvable and choice of the discount factor determines the policy. If a small discount factor is chosen then the future plays a negligible role as the agent gets greedy and will probably choose action b in state 2.

#### Q4. a.

- $T(A, \text{south}, C) = 1$   
The action south is taken twice from state A, and both times results in state C.  
 $2/2 = 1$
- $T(B, \text{east}, C) = 1$   
The action east is taken twice from state B, and both times results in state C.  
 $2/2 = 1$
- $T(C, \text{south}, E) = 0.75$   
The action south is taken four times from state C, and results in state E three times.  $3/4 = 0.75$
- $T(C, \text{south}, D) = 0.25$   
The action south is taken four times from state C, and results in state D one time.  
 $1/4 = 0.25$

#### b.

$$V(A) = 8, V(B) = -2, V(C) = 4, V(D) = -10, V(E) = 10$$

#### EXPLANATION

The estimated value of  $\hat{V}^{\pi}(s)$  is equal to the average value achieved starting from that state.

$\hat{V}^{\pi}(A)$ : Episodes 1 and 4 start from state A and both result in a utility of 8.  $\frac{8+8}{2} = 8$

$\hat{V}^{\pi}(B)$ : Episodes 2 and 3 start from state B. Episode 2 results in -12, while episode 3 results in 8.  $\frac{8-12}{2} = -2$

$\hat{V}^{\pi}(C)$ : State C is visited in every episode. The remaining rewards from C in episodes 1, 3, and 4 total 9, while the remaining rewards in episode 2 total -11.  $\frac{9+9+9-11}{4} = 4$

$\hat{V}^{\pi}(D)$ : State D is only visited in episode 2 and has a remaining utility of -10.

$\hat{V}^{\pi}(E)$ : State E is visited in episodes 1, 3, and 4 and has a remaining utility of 10 in each state.  $\frac{10+10+10}{3} = 10$

#### Q5.

##### 5.1

M=0.6

The clockwise action was taken 5 times from A, and went to B 3 times.

N=0

The transition (A,clockwise,B) happened 3 times and had reward 0 every time.

O=0.4

The clockwise action was taken 5 times from A, and went to C 2 times.

P=-10

Both of the occurrences of (A,clockwise,C) had reward -10.

## 5.2

$$Q(A, \text{clockwise}) = -4.984$$

### EXPLANATION

$$V_k(B) = \max(Q_k(B, \text{clockwise}), Q_k(B, \text{counterclockwise})) = -3.76$$

$$V_k(C) = \max(Q_k(C, \text{clockwise}), Q_k(C, \text{counterclockwise})) = 0.72$$

$$\begin{aligned} Q(A, \text{clockwise}) &= T(A, \text{clockwise}, B) \times (R(A, \text{clockwise}, B) + \gamma V_k(B)) + \\ &T(A, \text{clockwise}, C) \times (R(A, \text{clockwise}, C) + \gamma V_k(C)) \\ &= .6 \times (0 + .5 \times -3.76) + .4 \times (-10 + .5 \times .72) = -4.984 \end{aligned}$$

$$Q(A, \text{counterclockwise}) = -5.336$$

### EXPLANATION

$$\begin{aligned} Q(A, \text{counterclockwise}) &= T(A, \text{counterclockwise}, B) \times \\ &(R(A, \text{counterclockwise}, B) + \gamma V_k(B)) + \\ &T(A, \text{counterclockwise}, C) \times (R(A, \text{counterclockwise}, C) + \gamma V_k(C)) \\ &= .4 \times (0 + .5 \times -3.76) + .6 \times (-8 + .5 \times .72) = -5.336 \end{aligned}$$

$$Q(B, \text{clockwise}) = -4.024$$

### EXPLANATION

$$V_k(A) = \max(Q_k(A, \text{clockwise}), Q_k(A, \text{counterclockwise})) = -4.24$$

$$\begin{aligned} Q(B, \text{clockwise}) &= T(B, \text{clockwise}, A) \times (R(B, \text{clockwise}, A) + \gamma V_k(A)) + \\ &T(B, \text{clockwise}, C) \times (R(B, \text{clockwise}, C) + \gamma V_k(C)) \\ &= .8 \times (-3 + .5 \times -4.24) + .2 \times (0 + .5 \times .72) = -4.024 \end{aligned}$$

$$Q(B, \text{counterclockwise}) = -9.624$$



**EXPLANATION**

$$\begin{aligned}
Q(B, \text{counterclockwise}) &= T(B, \text{counterclockwise}, A) \times \\
&\quad (R(B, \text{counterclockwise}, A) + \gamma V_k(A)) + \\
&\quad T(B, \text{counterclockwise}, C) \times (R(B, \text{counterclockwise}, C) + \gamma V_k(C)) \\
&= .8 \times (-10 + .5 \times -4.24) + .2 \times (0 + .5 \times .72) = -9.624
\end{aligned}$$

$$Q(C, \text{clockwise}) = 0.376$$

**EXPLANATION**

$$\begin{aligned}
Q(C, \text{clockwise}) &= T(C, \text{clockwise}, A) \times (R(C, \text{clockwise}, A) + \gamma V_k(A)) + \\
&\quad T(C, \text{clockwise}, B) \times (R(C, \text{clockwise}, B) + \gamma V_k(B)) \\
&= .6 \times (0 + .5 \times -4.24) + .4 \times (6 + .5 \times -3.76) = .376
\end{aligned}$$

$$Q(C, \text{counterclockwise}) = -8.328$$

**EXPLANATION**

$$\begin{aligned}
Q(C, \text{counterclockwise}) &= T(C, \text{counterclockwise}, A) \times \\
&\quad (R(C, \text{counterclockwise}, A) + \gamma V_k(A)) + \\
&\quad T(C, \text{counterclockwise}, B) \times (R(C, \text{counterclockwise}, B) + \gamma V_k(B)) \\
&= .2 \times (0 + .5 \times -4.24) + .8 \times (-8 + .5 \times -3.76) = -8.328
\end{aligned}$$

## Q6.

### EXPLANATION

The only value that gets updated is  $\hat{V}^\pi(B)$ , because the only transition observed starts in state B.

$$\hat{V}^\pi(A) = 1$$

$$\hat{V}^\pi(B) = .5 * 2 + .5 * (-2 + 8) = 4$$

$$\hat{V}^\pi(C) = 8$$

$$\hat{V}^\pi(D) = 10$$

$$\hat{V}^\pi(E) = 10$$

## Q7.

### EXPLANATION

For each  $s, a, s', r$  transition sample, you update the Q value function as follows:

$$Q(s, a) = (1 - \alpha) * Q(s, a) + \alpha * (R(s, a, s') + \gamma * \max_{a'} Q(s', a'))$$

First we update  $Q(A, \text{counterclockwise}) = .5 \times 3.153 + .5 \times (8 + .5 \times 2.73) = 6.259$

Then we update  $Q(C, \text{counterclockwise}) = .5 \times 2.133 + .5 \times (0 + .5 \times 6.259) \approx 2.631$

Note that we use the updated value of  $Q(A, \text{counterclockwise})$  in the second update.

Because there are only two samples, the other four values stay the same.