

HOMEWORK-2

NAME : MADHAN THANGAVEL

ID : B00814926 (mthonga1)

4) MDP

a) states : { Top, Bottom }

Action : { DRIVE, STOP }

Transition Func : $T(s, a, s')$ \Rightarrow probability that action - a taken from state - s leads to state - s' .ie : $P(s' | s, a)$

Reward Func : +3 (Top), +4 (bottom)

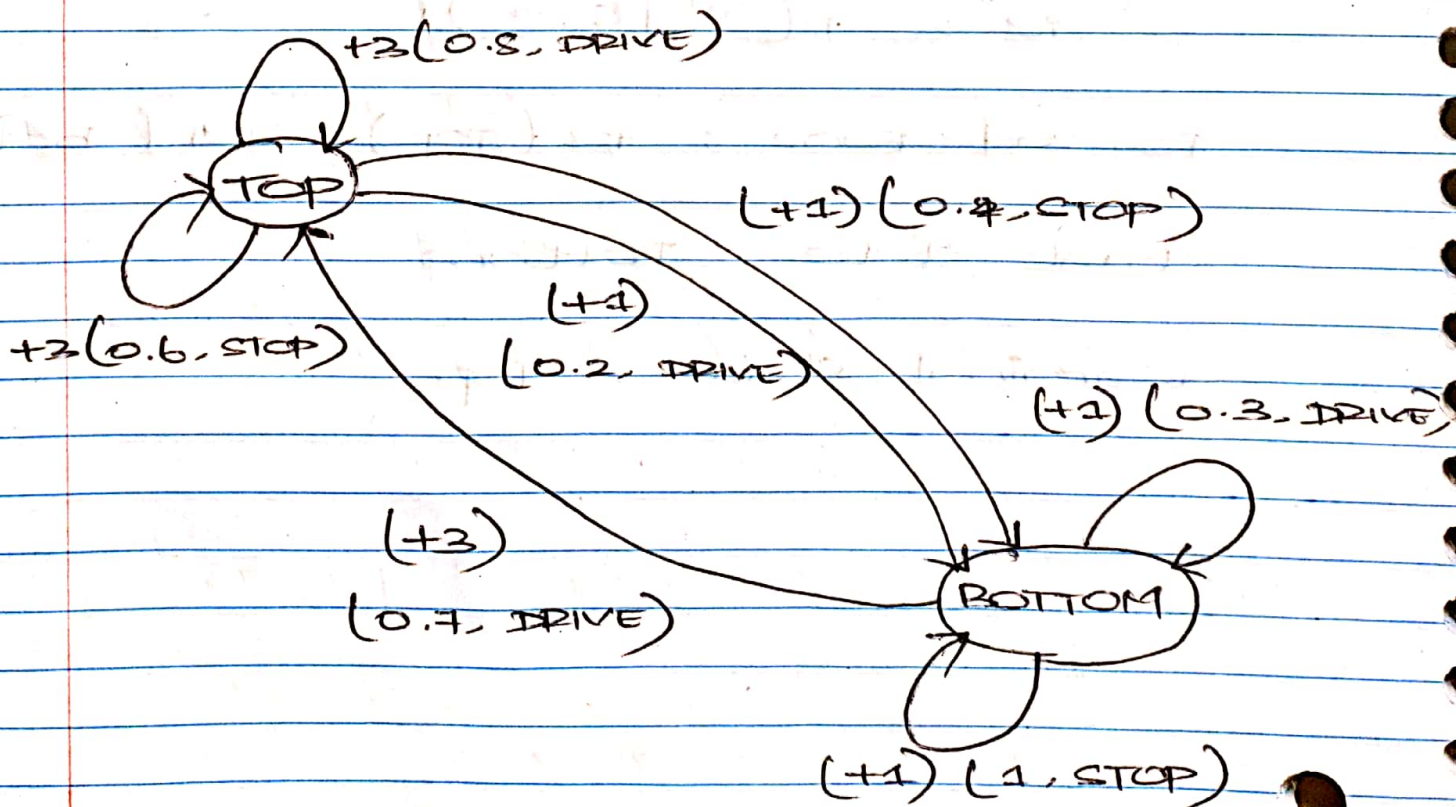
start state : Bottom

Terminal state : Top

(GIVEN) TRANSITION VALUES :

s	a	s'	$T(s,a,s')$	$R(s,a,s')$	COST
TOP	DRIVE	TOP	0.8	+3	-1
TOP	DRIVE	BOTTOM	0.2	+1	-1
TOP	STOP	TOP	0.6	+3	
TOP	STOP	BOTTOM	0.4	+1	
BOTTOM	DRIVE	TOP	0.7	+3	-1
BOTTOM	DRIVE	BOTTOM	0.3	+1	-1
BOTTOM	STOP	BOTTOM	1.0	+1	

→ TRANSITION DIAGRAM :



→ VALUE ITERATION :

We know that,

$$V_{k+1}(s) \leftarrow \max_a \sum T(s, a, s')$$

$$[R(s, a, s') + \gamma V_k(s')]$$

$\gamma \Rightarrow$ Discount Factor.

so let us compute $v_1(s)$ and $v_2(s)$ for both the states.

$v_1(s)$

$$\star \text{ TOP: } v_1(s) = \max [0.8(3-1) + 0.2(1-1), \\ 0.6(3) + 0.4(4)]$$

$$= \max [4.6, 2.2]$$

$$v_1(s) = 4.6$$

\star BOTTOM :

$$v_2(s) = \max [0.7(3-1) + 0.3(1-1), \\ 4(4)]$$

$$= \max [1.4, 4]$$

$$v_2(s) = 4$$

$V_2(c)$

* TOP :

$$\begin{aligned} V_2(c) &= \text{Max} \left[0.8 \left((0.8(3-1) + 0.8(2.2)) \right. \right. \\ &\quad \left. \left. + 0.2(0.8(4)) \right), \right. \\ &\quad \left. 0.6(3 + 0.8(2.2)) + 0.4(1 + 0.8) \right] \\ &= \text{Max} \left[(0.8(3.76) + 0.16), (0.6(4.76) \right. \\ &\quad \left. + 0.32) \right] \\ &= \text{Max} [3.168, 3.176] \\ &= 3.176 // \end{aligned}$$

* BOTTOM :

$$\begin{aligned} V_2(c) &= \text{Max} \left[0.7 \left((3-1) + 0.8(1.4) \right) \right. \\ &\quad \left. + 0.3(1-1 + 0.8(1.4)) \right), \\ &\quad \left. 1(1 + 0.8(1)) \right] \\ &= \text{Max} \left[(1.4 + 1.12) + 0.3(1.12), 1.8 \right] \\ &= \text{Max} [5.88, 1.8] \\ &= 5.88 // \end{aligned}$$

We know that V_0 is zero for both the states.

Therefore

V	TOP	BOTTOM
V_0	0	0
V_1	2.2	2.4
V_2	3.176	5.88

d) We know that,

$$V^{\pi}(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^{\pi}(s')]$$

$\gamma \Rightarrow$ Discount Factor.

$$\text{As } V_0^{\pi}(s) = 0$$

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

$$V_0^{\pi}(s) = 0$$

$$V_1^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') R[s, \pi(s), s'] + \gamma V_0^{\pi}(s')$$

(top)

$$= 0.8 [(3-1) + 0.2(0) + 0.6(2) + 0.4]$$

$$= 1.6 + 1.8 + 0.4$$

$$= 3.8$$

(bottom)

$$= 0.7 [(3-1) + 0.3(0) + 1(1)]$$

$$= 1.4 + 1$$

$$= 2.4$$

$$V_2^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') R[s, \pi(s), s'] + \gamma V_1^{\pi}(s')$$

(top)

$$= [0.8(0.8 \times 2) + (0.8)(3.8)]$$

$$= 1.28 + 3.04$$

$$= 4.32$$

Bottom.

$$= [0.7(2) + 0.8(1.4) + 0.3(0.8) + 4.4(2.8)]$$

$$= [2.52 + 3.04]$$

$$= 5.56$$

V	TOP	BOTTOM
$V_0^n(s)$	0	0
$V_1^n(s)$	3.8	2.4
$V_2^n(s)$	4.32	5.56

e) POLICY EXTRACTION

We know that

$$Q^*(s, a) = \sum_{s'} \tau(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$\pi^*(s) = \arg_{a \max} Q^*(s, a)$$

Therefore,

$$\pi^*(s) = \arg_{a \max} E [T(s, a, s') [R(s, a, s') + \gamma V^*(s')]]$$

Now convert $V^*(s)$ into $Q(s, a)$

$$Q[\text{TOP}, \text{DRIVE}] = [0.8(2 + 0.8^* 3.176) + 0.2(0.8(5.88))]$$

$$\Rightarrow 5.09$$

$$\begin{aligned} Q[\text{TOP}, \text{STOP}] &= [0.6(2 + 0.8)(3.176) + 0.4(1 + 0.8(5.88))] \\ &= [3.32 + 5.104] \\ &= 8.42. \end{aligned}$$

$$Q[\text{BOTTOM}, \text{DRIVE}] =$$

$$[0.7(2 + 0.8(5.88)) + 0.3(0.8(5.88))]$$

$$= [4.69 + 1.41]$$

$$\Rightarrow 6.10$$

$$Q[\text{BOTTOM}, \text{STOP}] = [1(1 + 0.8(5.88))]$$

$$\Rightarrow 5.70$$

$$\pi^*(s) = \arg \max_a \sum_{s'} p(s'|s, a) u(s')$$

where

$$u(s^*) = p(s) + \gamma \max_a \sum_{s'} p(s'|s, a) \times u(s')$$

so,

$$\pi^*(\text{TOP}) = \arg \max [Q[\text{TOP}, \text{DRIVE}], Q[\text{TOP}, \text{STOP}]]$$

$$= \arg \max [5.09, 8.42]$$

$$= 8.42 (\text{stop})$$

$$\pi^*(\text{BOTTOM}) = \arg \max [Q[\text{BOTTOM}, \text{DRIVE}], Q[\text{BOTTOM}, \text{STOP}]]$$

$$= \arg \max [6.10, 5.70]$$

$$= 6.10 (\text{drive})$$

Hence the policy

$\pi^*(s) = \{ \text{stop, pause} \}$ in this iteration.

2) MAP colouring

a) code reads

neighbours (y, x) :- occ (x, y)

block 1 //

→ 1 { colour (x, c) : colours (c) }

1 :- countries (x).

Why because, to ensure all the country has exactly one colour with aggregates.

colour (x, red) | colour (x, green)

| colour (x, blue) | colour (x, white)

:- countries (x).

b) 'not' represents two countries cannot have same color.

If x is the country and c is the color, then it is unknown that country is having color. Which then it is having a color.

c) When neighbour rules are removed,

It throws an error / scanning crashing

"atom does not occur in my rule head neighbours (x_1, x_2).

When the action is not performed, then the clingo response by stating atom not present.