

1: MDP

The following MDP represents a simple transportation problem.

- States: $s_1, s_2, s_3, s_4, s_5, s_6, s_7$; s_1 is the start state
- Actions: In states s_1, s_2 and s_3 , two actions are possible: WALK and TELEPORT
States s_4, s_5, s_6 and s_7 are terminal/absorbing states.
- Reward: $R(s_1) = R(s_2) = R(s_3) = -0.1$
 $R(s_4) = -4$
 $R(s_5) = 2$
 $R(s_6) = 5$
 $R(s_7) = -2$

The transition probabilities are given by the following tables:

ACTION: WALK		Ending State						
		s_1	s_2	s_3	s_4	s_5	s_6	s_7
Starting State	s_1	0.5	0.5	0	0	0	0	0
	s_2	0	0	0	0.5	0.5	0	0
	s_3	0	0	0.5	0	0	0.5	0

ACTION: TELEPORT		Ending State						
		s_1	s_2	s_3	s_4	s_5	s_6	s_7
Starting State	s_1	0	0.5	0.5	0	0	0	0
	s_2	0	0	0.25	0.75	0	0	0
	s_3	0	0	0	0	0	0.5	0.5

Assuming $V_0(s_1) = 0, V_0(s_2) = 0, V_0(s_3) = 0$ and for the terminal states:

$$V_0(s_4) = R(s_4) = -4, V_0(s_5) = R(s_5) = 2, V_0(s_6) = R(s_6) = 5, V_0(s_7) = R(s_7) = -2$$

(a) After one step of value iteration, what is the utility of $V_1(s_1), V_1(s_2)$ and $V_1(s_3)$? Please show all of your work.

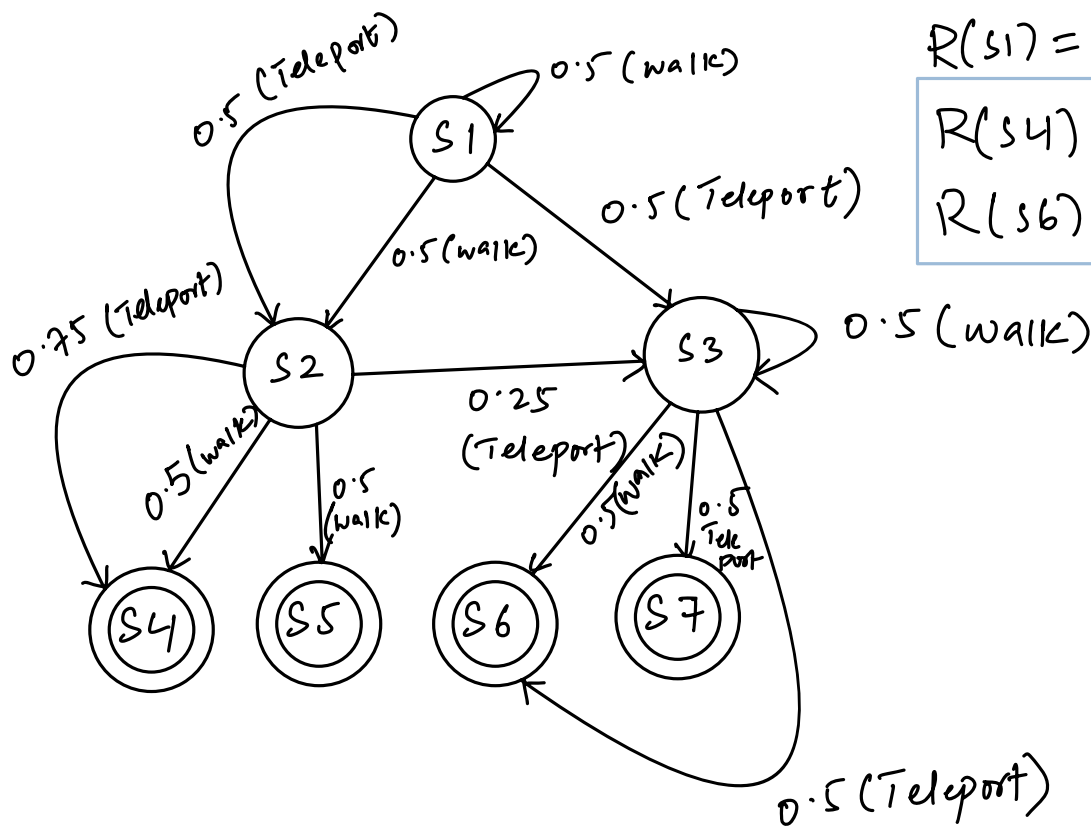
(b) After two steps of value iteration. What is the utility of $V_2(s_1), V_2(s_2)$ and $V_2(s_3)$? Please show all of your work.

(c)(Programming) Please write a program to perform value iteration on this problem. Recall that the value iteration algorithm iterates and updates the utilities using:

$$U_{t+1}(S_i) = \max_a \left(R(S_i) + \gamma \sum_{S_j} P(S_j|S_i, a) \times U_t(S_j) \right)$$

- What is the utility of the states s_1, s_2 , and s_3 after two iterations with $\gamma = 1$? Verify that the program output matches the result in question b.
- Assuming $\gamma = 0.95$, what is the utility of the states s_1, s_2 , and s_3 when the algorithm converges? What is the optimal policy for this MDP?

Utility: sum of (discounted) reward
 V : Maximize expected utility



$$R(S1) = R(S2) = R(S3) = -0.1$$

$$R(S4) = -4 \quad R(S5) = 2$$

$$R(S6) = 5 \quad R(S7) = -2$$

Terminal.

$$V_0(S1) = V_0(S2) = V_0(S3) = 0 \quad \# \text{ initial expected utility (max sum of rewards)}$$

$$V_0(S4) = R(S4) = -4, \quad V_0(S5) = R(S5) = 2, \quad V_0(S6) =$$

$$R(S6) = 5, \quad V_0(S7) = R(S7) = -2$$

a) # iteration 1

$$V_1(S1) = \max_a \left(R(S1) + \gamma \sum_{Sj} P(Sj | S1, a) * V_0(Sj) \right)$$

$$= \max_a \left(\begin{bmatrix} -0.1 + (0.5 \cdot 0) + (0.5 \cdot 0) \\ -0.1 + (0.5 \cdot 0) + (0.5 \cdot 0) \end{bmatrix} \right)$$

walk
Teleport

$$V_1(s_1) = -0.1$$

$$V_1(s_2) = \max_a \left(R(s_2) + \gamma \sum_{s_j} P(s | s_2, a) * V_0(s_j) \right)$$

$$= \max_a \left(\begin{bmatrix} -0.1 + (0.5 * -4) + (0.5 * 2) \\ -0.1 + (0.75 * -4) + (0.25 * 0) \end{bmatrix} \right)$$

$$V_1(s_2) = -1.1$$

$$V_1(s_3) = \max_a \left(R(s_3) + \gamma \sum_{s_j} P(s | s_3, a) * V_0(s_j) \right)$$

$$= \max_a \left(\begin{bmatrix} -0.1 + (0.5 * 5) + (0.5 * 0) \\ -0.1 + (0.5 * 5) + (0.5 * -2) \end{bmatrix} \right)$$

$$V_1(s_3) = 2.4$$

b) # iteration # 2

$$V_2(s_1) = \max_a \left(R(s_1) + \gamma \sum_{s_j} P(s_j | s_1, a) * V_1(s_j) \right)$$
$$= \max_a \left(\begin{bmatrix} -0.1 + (0.5 * -0.1) + (0.5 * -1.1) & \text{Walk} \\ -0.1 + (0.5 * -1.1) + (0.5 * 2.4) & \text{Teleport} \end{bmatrix} \right)$$

$$V_2(s_1) = 0.55$$

$$V_2(s_2) = \max_a \left(R(s_2) + \gamma \sum_{s_j} P(s_j | s_2, a) * V_1(s_j) \right)$$
$$= \max_a \left(\begin{bmatrix} -0.1 + (0.5 * -4) + (0.5 * 2) & \text{Walk} \\ -0.1 + (0.75 * -4) + (0.25 * 2.4) & \text{Teleport} \end{bmatrix} \right)$$

$$V_2(s_2) = -1.1$$

$$V_2(s_3) = \max_a \left(R(s_3) + \gamma \sum_{s_j} P(s_j | s_3, a) * V_1(s_j) \right)$$
$$= \max_a \left(\begin{bmatrix} -0.1 + (0.5 * 5) + (0.5 * 2.4) & \text{Walk} \\ -0.1 + (0.5 * 5) + (0.5 * -2) & \text{Teleport} \end{bmatrix} \right)$$

$$V_2(s_3) = 3.6$$

(c)

(i) same as part (b)

(ii) s_1 - utility : 1.45
 s_2 - utility : -1.05
 s_3 - utility : 4.33

Policy: Teleport

Policy: Walk

Policy: Walk