

Theory
Que 1

Initial values given :-

$$V(S_0) = 0 ; V(S_1) = 0 ; V(S_2) = 0 ; V(S_3) = 0$$

(a) We know that

$$V^{k+1}(s) = \max_a \left(R(s,a) + \gamma \sum P(s'|s,a) V^k(s') \right)$$

R - Reward
P - Prob
a - actions
γ - discount

Iteration 1

$$V^1(S_0) = \max_a \left(1 + 0.9 \left(0.5 * 0 + 0.5 * 0 \right) \right) = 1$$

$$V^1(S_1) = 2$$

$$V^1(S_2) = 3$$

$$V^1(S_3) = 10$$

In first iteration, values come out to be same as rewards. Initial values are for all states.

Iteration 2

$$\begin{aligned} V^2(S_0) &= \max \left(1 + 0.9 \left(0.5 * V^1(S_1) + 0.5 * V^1(S_2) \right) \right) \\ &= 1 + 0.9 (0.5 * 2 + 0.5 * 3) \\ &= 1 + 0.9 (1 + 1.5) \\ &= 1 + 0.9 (2.5) \\ &= 3.25 \end{aligned}$$

$$\begin{aligned} V^2(S_1) &= \max \left(2 + 0.9 \left(0.5 * V^1(S_1) + 0.5 * V^1(S_2) \right) \right) \\ &= \max \left(2 + 0.9 (0.5 * 2 + 0.5 * 10) \right) \\ &= \max \left(2 + 0.9 (3) \right) \\ &= \max \left(2 + 0.9 (6), 2 + 0.9 (3) \right) \\ &= 7.4 \end{aligned}$$

$$v^2(s_2) = \max \left(3 + 0.9 \left(1 * v'(s_0) \right) \right) \\ = 3 + 0.9 * 2 = 3.9$$

$$v^2(s_3) = \max \left(10 + 0.9 \left(1 * v'(s_3) \right) \right) \\ = 10 + 0.9(1 * 10) \\ = 10 + 9 = 19$$

Iteration 3

$$v^2(s_0) = \max \left(1 + 0.9 \left(0.5 * 7.4 + 0.5 * 3.9 \right) \right) \\ = 1 + 0.9(3.7 + 1.95) \\ = 6.085$$

$$v^2(s_1) = \max \left(2 + 0.9 \left(0.5 * 7.4 + 0.5 * 19 \right), 2 + 0.9(3.9) \right) \\ = \max(2 + 0.9(13.2), 2 + 0.9(3.9)) \\ = \max(13.88, 5.51) \\ = 13.88$$

$$v^2(s_2) = \max \left(3 + 0.9 \left(1 * 13.25 \right) \right) \\ = 3 + 0.9(13.25) = ~~5.825~~ 5.925$$

$$v^2(s_3) = \max \left(10 + 0.9 \left(1 * 19 \right) \right) \\ = 10 + 0.9(19) = 27.1$$

(b) It can be noticed that the optimal policy for state S_1 is :- the movement from S_1 to S_3 . This is because S_3 has a larger reward.

$$\begin{aligned} V'(S_3) &= R(S_3) + \gamma * V'(S_3) \\ &= 10 + 0.9 * V'(S_3) \\ &= 100 \end{aligned}$$



Now, optimal value for state S_1 ,

$$\begin{aligned} V'(S_1) &= 2 + 0.9 (0.5 * V'(S_1) + 0.5 * V'(S_3)) \\ &= ~~85.454545~~ 85.454545... \end{aligned}$$

Thus, it is clear that best policy for state S_1 is moving from S_1 to S_3 .

(c) (i) False

Since ~~the value~~ the value will get updated in every iteration, it won't converge.

$$V(S_t)^{t+1} = \max_a (R(S_t, a) + \gamma \sum_{S'} P(S' | S_t, a) V(S')^t)$$

(ii) False

MDP does not converge. The value keeps getting updated in every iteration.

(iii) True

If $\gamma = 0$, then $V(s_t)^t = \max_a (R(s, a))$
So, value will be same as the reward for that state.

(iv) True,

Ayclic MDP \Rightarrow having no cycles.

After every iteration in an ayclic MDP, at least one state gets to optimal value.

Then after N iterations, ayclic MDP converges.

(v) ~~True~~ False

Since it is given that there are no absorbing goal states, so, MDP won't converge.