

Ex-02

i) Formulating an MDP:

a) The State space  $S$  in 4 rooms.  
Consists of 4 rooms themselves.

i.e. State:  $(t, y)$   $t, y \in [0, 10]$ .  
 $(t, y) \neq (5, 1), (j, 5), (K, 4)$

$i = 0, 2, 3, 4, 5, 6, 7, 8, 9, 10$ .

$j = 0, 2, 3, 4$ .

$K = 6, 7, 9, 10$ .

$\equiv x \in [0, 10], y \in [0, 10] - [(0, 5), (2, 5), (3, 5), (4, 5), (5, 5), (6, 4), (7, 4), (9, 4), (10, 4), (5, 6), (5, 2), (5, 3), (5, 4), (5, 6), (5, 7), (5, 9), (5, 10)]$

Whereas, action Space  $A$  consists of 4 possible actions move.

Action = { North, South, East, west }.

b) In 4 room domain has 17 walls.

So state  $s_n = 121 - 17$

$= 104$ . when not

Action is 'East' or 'West' blocked by walls

(S, a) Pair there are 3 Non zero  $P(S_{i+1} | S_i, a)$  value.

No. of non zero when non bloled by walls  $= 104 \times 4 \times 3$   
 $= 1248$ .

when blocked by walls  $= 104 \times 4 \times 1 = 416$ .

Appox. the no. of non zero rows can be 832.

Q1 C)

## PSEUDOCODE

Initialize an empty table for transition probabilities

For each cell in the grid (excluding the boundary):

  If the cell is not a wall:

    For each action in [UP, DOWN, LEFT, RIGHT]:

      Calculate the resulting cell after the action

      If the resulting cell is within the grid and not a wall:

        If the action does not take the agent out of bounds:

          Add an entry to the table with:

- The current cell as the current state  $s$
- The resulting cell as the next state  $s'$
- The action  $a$
- The transition probability (0.8 for the intended direction, 0.1 for the perpendicular directions)
- The reward  $r$  (which is 0 unless transitioning into the goal state, which gives a reward of 1)

  If the current cell is the goal state (10, 10):

    Add an entry to the table with:

- The current cell as the current state  $s$
- The start state (0, 0) as the next state  $s'$  for all actions
- The action  $a$
- A transition probability of 1
- The reward  $r$  of 0

②

2) 1 point RL objective.

- a) Treated pole balancing as an episodic task but also used discounting, with all rewards zero except -1 upon failure.

Defining failure time as  $T$

Sum of all failure discounted rewards,  $t < T$  will be 0

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots + \gamma^{T-t-1} R_{T-1}$$

$\therefore G_t = 0$  for  $t < T$  as  $R_{t+1} = R_{t+2} = 0$ .

As the agent receives a reward (-1) is a terminal time step  $G_t$  is simply immediate reward

$$G_t = -1$$

$\therefore$  we can formulate

$$G_t = -\gamma^{T-t-1}$$

While continuing formulation. (Continuous task)

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$= \gamma^k (-1)$$

$$G_t = -\gamma^k \quad \text{Same as expression above}$$

$$\therefore G_t = \gamma^k (-1) + \gamma^{2k} (-1) + \gamma^{3k} (-1) \dots$$

③

$$\therefore G_t = -\gamma^k [1 + \gamma^k + \gamma^{2k} + \dots]$$

$$\therefore \boxed{G_t = -\gamma^k \cdot \frac{1}{1-\gamma^k}}$$

Thus the episodic case agent focuses on single outcome at episode end, while in continuing case,

- b) For episodic task while discounting  $\gamma=1$ ,  
 where reward is 0 at every time step  
 $T$  is a terminal state.  
 and  $R_T = 1$  for the final time step.

The expected reward

$$G_t = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots + \gamma^{T-1} R_T$$

$$R_1 = R_2 = R_3 = 0$$

$$\therefore G_t = 0 + 1(0) + 1^2(0) + \dots + 1 \cdot 1$$

$$\therefore \boxed{G_t = 1}$$

$\therefore G_t$  always be 1, so it will never have change,

$\Rightarrow$  Thus, we need to know to provide intermediate reward to let agent know and learn to achieve better.

3) 1 point : Discounted return.

a)

$$\gamma = 0.5, T = 5$$

$$R_1 = -1 \quad R_4 = 3$$

$$R_2 = 2 \quad R_5 = 2$$

$$R_3 = 6$$

} (given)

To find:

$$G_0, G_1, G_2, G_3, G_4, G_5 = ?$$

Soln:

$$G_t = R_{t+1} + \gamma G_{t+1}$$

— (eqn 3.9)

$$\therefore t = 5$$

$$G_5 = R_6 + \gamma G_6 = 0$$

$$\therefore \boxed{G_5 = 0}$$

$$t = 4$$

$$G_4 = R_5 + \gamma G_5$$

$$= 2 + 0.5(0)$$

$$\therefore \boxed{G_4 = 2}$$

$$t = 3$$

$$G_3 = R_4 + \gamma G_4$$

$$= 3 + (0.5) 2$$

$$\boxed{G_3 = 4}$$

$$t = 2$$

$$G_2 = R_3 + \gamma G_3$$

$$= 6 + (0.5) 4$$

$$\boxed{G_2 = 8}$$

$$t=1$$

$$\begin{aligned} G_1 &= R_2 + \gamma G_2 \\ &= 2 + (0.5) 8 \\ \boxed{G_1 &= 6} \end{aligned}$$

$$t=0$$

$$\begin{aligned} G_0 &= R_1 + \gamma G_1 \\ &= -1 + (0.5) 6 \\ \boxed{G_0 &= 2} \end{aligned}$$

$$G_0 = 2, G_1 = 6, G_2 = -8, G_3 = 4, G_4 = 2, G_5 = 0 \text{ etc.}$$

$$b) \quad \gamma = 0.9$$

$$R_1 = 2$$

$$R_2 = R_3 = R_4 \dots R_n = 7$$

$$\text{To find: } G_0, G_1 = ?$$

$$\text{Soln: } G_t = \gamma G_{t+1} + R_{t+1}$$

$$\therefore t=1$$

$$\begin{aligned} G_1 &= R_2 + \gamma G_2 \\ G_2 &= R_3 + \gamma G_3 \\ &= 7 + \gamma(7) + \gamma^2(7) \\ &= 7(1 + \gamma + \gamma^2 + \dots) \\ &= 7 \left( \frac{1}{1-\gamma} \right) = 7 \left( \frac{1}{1-0.9} \right) \end{aligned}$$

$$\boxed{G_2 = 70}$$

$$G_1 = 7 + (0.9)(70)$$

$$\boxed{G_1 = 70}$$

$$t=0$$

$$\begin{aligned} G_0 &= R_1 + \gamma G_1 \\ &= 2 + (0.9)(70) \end{aligned}$$

$$\boxed{G_0 = 65}$$

$$\therefore \boxed{\begin{matrix} G_0 = 65 \\ G_1 = 70 \end{matrix}}$$



4) Case I: Up motion

$$\begin{aligned}
 C_{up} &= 50 + \sum_{n=2}^{101} \gamma^{n-1} R_n \\
 &= 50 + (-\gamma) \frac{(1 - \gamma^{100})}{1 - \gamma} = \frac{50 - \gamma(1 - \gamma^{100})}{1 - \gamma}
 \end{aligned}$$

Case II: Down motion

$$\begin{aligned}
 C_{down} &= -50 + \sum_{n=2}^{101} \gamma^{n-1} R_n \\
 &= -50 + \frac{\gamma(1 - \gamma^{100})}{1 - \gamma}
 \end{aligned}$$

So  $C_{up} - C_{down}$ ,

$$\begin{aligned}
 &= \frac{50 - \gamma(1 - \gamma^{100})}{1 - \gamma} - \left( -50 + \frac{\gamma(1 - \gamma^{100})}{1 - \gamma} \right) \\
 &= 100 - 2 \frac{\gamma(1 - \gamma^{100})}{1 - \gamma}
 \end{aligned}$$

Using Wolfram alpha,

$$-1.047 < \gamma < 0.9843$$

And,  $C_{down} - C_{up}$

$$\begin{aligned}
 &= -50 + \frac{\gamma(1 - \gamma^{100})}{1 - \gamma} - \frac{50 - \gamma(1 - \gamma^{100})}{1 - \gamma} \\
 &= -100 + 2 \frac{\gamma(1 - \gamma^{100})}{1 - \gamma}
 \end{aligned}$$

for discount factor bet<sup>n</sup>.

$0.9843 < \gamma < 1$  is better to take down motion.

5) 1 point: Modifying the reward function.

$$a) G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \dots \text{ (eqn 3.8)}$$

Adding constant to all rewards.

$$G'_t = (R_{t+1} + c) + \gamma (R_{t+2} + c) + \gamma^2 (R_{t+3} + c) \dots$$

Factoring out c to see the effect:

$$G'_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots \\ c(1 + \gamma + \gamma^2 \dots)$$

The sum of Series

$$G_t = \frac{1}{1-\gamma}$$

if  $\gamma < 1$

So, adding a constant across all states.

$$\boxed{V_c = \frac{c}{1-\gamma}}$$

This added value  $V_c$  is constant across all states and does not depend on the policy, so it does not affect the relative values.

This means policy that maximizes the original return  $G_t$  will also maximize return  $G'_t$ .

$$V_{\pi c}(s) = V_{\pi}(s) + V_c \\ = V_{\pi}(s) + \frac{c}{1-\gamma} //$$



⑧  
b) For an episodic task.

$$G_t = (R_1 + c) + \gamma(R_2 + c) + \gamma^2(R_3 + c) \dots + \gamma^{n-1}(R_n + c)$$

$$G_t = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots + \gamma^{n-1} R_n + c(1 + \gamma + \gamma^2 + \dots + \gamma^{n-1})$$

The sum of series

$$\frac{1 - \gamma^n}{1 - \gamma}$$

$$\begin{aligned} V_{\pi_c}(s) &= E_{\pi} [G_t | s_t = s] \\ &= E_{\pi} [G_t | s_t = s] + E_{\pi} [c_T | s_t = s] \\ &= V_{\pi}(s) + \underbrace{c E_{\pi} [T | s_t = s]} \end{aligned}$$

$V_c$  now depends on  $n$ , length of the episode

$\therefore$  The shifts ~~extra~~ in the return will not be constant.

$\therefore$  This will have an effect.

6) 1 point Bellman eq<sup>n</sup>:

a)

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) [\gamma + \gamma V_{\pi}(s')] \quad (3.14)$$

for all  $s \in S$

Since the policy is equiprobable action has a probability of 0.25.

$$\begin{array}{ccccc} & & 2.3 & \leftarrow V_n & \\ v_w \rightarrow & 0.1 & \textcircled{0.7} & 0.4 & \leftarrow V_e \\ & & -0.4 & \leftarrow V_s & \end{array}$$

$\therefore \gamma$  is less than 1

$$\therefore \gamma = 0.933$$

$\therefore$  Equation for centre state  $V_c$ , reward is 0.

$$V_c = 0.25 (\gamma + \gamma V_n) + 0.25 (\gamma + \gamma V_s) + 0.25 (\gamma + \gamma V_e) + 0.25 (\gamma + \gamma V_w)$$

$$V_c = 0.25 (0 + (0.93)(2.3)) + 0.25 (0 + (0.93)(-0.4)) + 0.25 (0 + (0.93)(0.4)) + 0.25 (0 + (0.93)(0.1))$$

$$\therefore \boxed{V_c = 0.7}$$

This matches the value of centre state in the figure.

$$b) \quad q_*(s, a) = \sum_{s'} p(s', r | s, a) [r + \gamma \max_a q_*(s', a)]$$

$$\gamma = 0.9$$

$$\text{centre value} = 17.8$$

$$\begin{array}{ccc} & 19.8 & \\ 19.8 & \textcircled{17.8} & 16 \\ & 16 & \end{array}$$

$$= \frac{1}{2} \times 1(0 + 0.9 \times 19.8)$$

$$+ \frac{1}{2} \times 1(0 + 0.9 \times 19.8)$$

$$V_c \approx 17.82$$

$$\boxed{V_c = 17.8}$$

7) Guessing and Verifying value fun<sup>n</sup>:

a) For 3-State MDP is,

$$V = 0.5$$

To verify the value fun<sup>n</sup>

$$\begin{aligned} V_{\pi}(S) &= \sum \pi(a|S) \sum P(S', r | S, a) [r + \gamma V_{\pi}(S')] \\ &= \frac{1}{2} \times 1 \times (1 + 1 \times 0) + \frac{1}{2} \times 1 \times (0 + 1 \times 0) \\ &= \frac{1}{2} + \frac{1}{2} \times 0 = \frac{1}{2} \end{aligned}$$

$$\therefore \boxed{V_{\pi}(S) = 0.5} \quad \therefore$$

$\therefore \text{LHS} = \text{RHS}$

b) Value fun<sup>n</sup> for 7 state MDP

$$\therefore V_{\pi}(S) = \frac{1}{2} \times 1 \times (0 + 1 \times 0) + \frac{1}{2} \times 1 \times (0 + 1 \times V_{\pi}(S)(B)) = \frac{1}{2} V_{\pi}(S)(B)$$

$$\therefore V_{\pi}(S)(B) = \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(A)) + \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(C))$$

$$\therefore V_{\pi}(S)(C) = \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(B)) + \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(D))$$

$$\therefore V_{\pi}(S)(D) = \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(C)) + \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(E))$$

$$\therefore V_{\pi}(S)(E) = \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(D)) + \frac{1}{2} \times 1 \times (0 + V_{\pi}(S)(F))$$

$$\therefore \boxed{V(A) = \frac{1}{6}, V(B) = 2 \cdot V(A) = \frac{1}{3}, V(C) = 3 \cdot V(A) = \frac{1}{2}, V(D) = 4 \cdot V(A) = \frac{2}{3}, V(E) = \frac{5}{6}}$$

c) For this value function for arbitrary is,

$$V_{\pi}(S) = \frac{i}{n-1} \quad \therefore \text{where } [i = 1, 2, 3, \dots, (n-2)]$$

8) Solving for Value Function.

a) Equation for the 2 States in the recycling robot arbitrary policy  $\pi(a|s)$

$$V_{\pi} = \sum_s \pi(a|s) \sum_{s', r} P(S'=r|s, a) (r + \gamma V(s'))$$

So,

$$\begin{aligned} V_{low} = & \pi(\text{search} | low) [(1-\beta)(3 + \gamma V_{high}) + \beta(\text{search} + \gamma V_{low})] \\ & + \pi(\text{wait} | low) [1 \times (\gamma \text{wait} + \gamma V_{low})] \\ & + \pi(\text{recharge} | low) [0 + 1 \times \gamma V_{high}] \end{aligned}$$

$$\begin{aligned} V_{high} = & \pi(\text{search} | high) [\alpha(\gamma \text{search} + V_{high}) + (1-\alpha)V_{search} + \gamma V_{low}] \\ & + \pi(\text{wait} | high) [\gamma \text{wait} + V_{high}] \end{aligned}$$

b)  $\pi(\text{search} | high) = 1$ ,  $\pi(\text{wait} | low) = 0.5$   
 $\pi(\text{recharge} | low) = 0.5$ ,  $\alpha = 0.7$ ,  $\beta = 0.8$ ,  $\gamma = 0.9$   
 $V_{search} = 10$ ,  $V_{wait} = 3$

$$V_{low} = 0.5 \times 1 \times (3 + 0.9 V_{low}) + 0.5 \times 1 \times (0.9 + V_{high})$$

~~$$V_{high} = 0.55 V_{low} = 1.5 + 0.45 V_{high} \quad \text{--- (i)}$$~~

$$V_{high} = 0.8 \times (10 + 0.9 \times V_{high}) + 0.2 \times (10 + 0.9 V_{low})$$

$$\therefore 0.55 V_{low} = 1.5 + 0.45 V_{high} \quad \text{--- (i)}$$

$$\therefore 0.28 V_{high} = 10 + 0.18 V_{low} \quad \text{--- (ii)}$$

Solving (i) & (ii) we get,

$\begin{aligned} V_{high} &= 79.03 \\ V_{low} &= 67.40 \end{aligned}$
---



c) So,

$$\pi(\text{wait low}) = 0$$

$$\pi(\text{recharge low}) = 1 - 0$$

$$\begin{aligned} V_{\text{low}} &= 0(3 + 0.9 V_{\text{low}}) + (1-0) 0.9 V_{\text{high}} \\ &= 30 + 0.9 \theta V_{\text{low}} + 0.9 - 0.9 \theta V_{\text{high}} \end{aligned}$$

$$(1 - 0.9 \theta) V_{\text{low}} = 30 + 0.9(1 - \theta) V_{\text{high}}$$

from (b)

$$0.28 V_{\text{high}} = 10 + 0.18 V_{\text{low}}$$

Solving we get

$$V_{\text{low}} = \frac{9 - 8.16 \theta}{0.118 - 0.9 \theta} = \frac{90.67 - 69.86 \theta}{0.118 - 0.9 \theta}$$

$$\begin{aligned} \therefore V_{\text{low}} &= 76.28 \\ V_{\text{high}} &= 84.75 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{when } \theta = 0$$