
CS 5180

Reinforcement Learning

Exercise 2

Sep 27, 2024

Anuj Patel 002874710



Ques 1) Part a: State Space (S)

The state space S consists of all possible positions the agent can occupy in the environment, that are represented as coordinates $p(x,y)$. here $x \in [0,10]$ and $y \in [0,10]$ and $y \in [0,10]$. The agent cannot take occupied cells.

Action Space (A):

The action space A consists of four possible actions: Up, Down, Left, Right.

Ques 1(b)

(1) calculate the number of states:

(2) we will consider here four room examples

• Total grid = $11 \times 11 \Rightarrow 121$ cells

• number of free cells $\Rightarrow 104$ cells (7 walls)

(2) Total action $\Rightarrow 4$.

(3) Total number of State-action pairs

$$= S \times A$$

$$\Rightarrow 104 \times 4$$

$$= 416$$
 pairs

(4) we will consider a deterministic movement. the agent move to the intended adjacent cells.

→ Possible outcome per action

- Intended move
- blocked by wall.

So there are 2 possible next state.

(c5) Reward function

- The agent will receive a reward of +1 when reaches a goal and 0 otherwise.
- Given s and a , the reward r is calculated by resulting s'

(c6) Total non zero entries =

$$= 15 \times 14 \times \text{possible } s' \text{ per } (s, a)$$

$$= 104 \times 4 \times 2$$

$$= 832 \text{ rows}$$

- So approximately 832 rows in conditional probability table $P(s'|s, a)$ have non-zero probability

≡ code.txt

```
1 start grid from the given array
2 define action space A with movements: Up, Down, Left, Right
3 Define state space S as all positions (i, j) where grid[i][j] == 0
4 Set goal_state = (10, 10)
5 Set start_state = (0, 0)
6
7 For each state s in S:
8     For each action a in A:
9         If s == goal_state:
10             s_prime = start_state
11             r = 0
12             probability = 1
13             Output (s, a, s_prime, r, probability)
14             Continue to next action
15             Compute intended position s_intended based on action a:
16             s_intended = (s_i + delta_i, s_j + delta_j)
17             If s_intended is within grid bounds and grid[s_intended] == 0:
18                 s_prime = s_intended
19             Else:
20                 s_prime = s // Collision with wall or boundary
21             If s_prime == goal_state:
22                 r = +1
23             Else:
24                 r = 0
25             probability = 1
26             Output [(s, a, s_prime, r, probability)]
```

Ques:- Part (a)

$\gamma = 0.5$, $T = 0.5$ and $R_1 = -1$, $R_2 = 2$, $R_3 = 6$, $R_4 = 3$, $R_5 = 2$
discounted return: $G_t = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots + \gamma^{T-t-1} R_T$

- Backward calculation

$G_5 = 0$ (no rewards are received after time 5)

$$G_4 = R_5 + \gamma G_5 \Rightarrow 2 + 0.5(0) = 2$$

$$G_3 = R_4 + \gamma G_4 \Rightarrow 3 + 0.5(2) = 4$$

$$G_2 = R_3 + \gamma G_3 \Rightarrow 6 + 0.5(4) = 8$$

$$G_1 = R_2 + \gamma G_2 \Rightarrow 2 + 0.5(8) = 6$$

$$G_0 = R_1 + \gamma G_1 \Rightarrow -1 + 0.5(6) = 2$$

Overall

$$G_0 = 2, G_1 = 6, G_2 = 8, G_3 = 4, G_4 = 2, G_5 = 0$$

Part (b)

$\gamma = 0.9$, $R_1 = 2$ followed by an infinite sequence of 7

For infinite sequence where rewards are constant after first reward we can use this formula:

$$\begin{aligned}G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + G_\infty \\&= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots)\end{aligned}$$

If all reward from R_2 onward are 7.

~~Endless loop~~

$$G_1 = 7 + \gamma \cdot 7 + \gamma^2 \cdot 7 + \dots + G_\infty$$

$$G_1 = 7(1 + \gamma + \gamma^2 + \dots)$$

$$G_1 = 7 \cdot \frac{1}{1-\gamma}$$

$$G_I = 7 \times \frac{1}{1-0.9} = 70$$

$$G_O = R_1 + 2G_I = 2 + 0.9(70) = 65$$

$$\boxed{G_I = 70} \text{ and } \boxed{G_O = 65}$$

due cases

Failure case

Reward $R_t = 0$ at all times except $R_{T-1} = -1$
Upon failure at time T

Episodic task

here ~~an episode ends upon failure. The return at each time step t is~~

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

all the rewards are zero except at failure

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

$$\text{stop} = -\gamma^{T-t-1}$$



continuing task

here the process goes on indefinitely without terminal state. failure is treated as a transition back to state x_0

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

non zero Reward $R_T = -1$, then process continue after T

$$G_t = -\gamma^{T-t-1} + \sum_{k=T-t}^{\infty} \gamma^k x_0$$

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

\rightarrow difference:

- both episodic and continuing tasks give the same result at each time step t

other episode:

- for the episodic task, the episode ends at failure and no further rewards are considered
- for continuing task, the agent continues operate after failure. It potentially occurs more rewards or penalties

\rightarrow

Point (b)

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

What is wrong?

- The reward structure does not differentiate between quick and slow solution. with $\gamma=1$ and only a reward at the end, every successful path through the maze gets the same total reward of t_1 , regardless of steps taken.
- the agent has no incentive to find the shortest path.
- all successful trajectories appear equally good to the agent.

how to improve:

- 1) Introduce discounting factor
- 2) add small negative reward for each step
- 3) Also add time-dependent reward.

Ques 4: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ → eq 3.8

Now we will add c to all rewards.

The new rewards are $R'_t + c = R_{t+1} + c$.

→ New G'_t using new rewards $R'_t + c$ is

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

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

$$G'_t = crt + \sum_{k=0}^{\infty} \gamma^k c$$

$$G'_t = G_t + \frac{c}{1-\gamma}$$

This shows that adding a constant c to all rewards shift the return by constant $\frac{c}{1-\gamma}$,

but the relative difference between return states are same.

→ Adding a constant to all rewards does not affect the value of any states

→ The value function for any state is under the modified reward system shifts by

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

value remains the same.

Q-4
cb)

Episodic value function

$$V_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t R_{t+1} | S_0 = s \right]$$

Add a constant c to each reward.

$$V'_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t (R_{t+1} + c) | S_0 = s \right]$$

We can write as

$$V'_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t R_{t+1} | S_0 = s \right] + \underbrace{\sum_{t=0}^T \gamma^t c}_{\Rightarrow V_{\pi}(s)}$$

$$\sum_{t=0}^T \gamma^t c = c \times 1 - \frac{\gamma^{T+1}}{1-\gamma}$$

Thus modified function values.

$$V'_{\pi}(s) = V_{\pi}(s) + c \times \frac{1 - \gamma^{T+1}}{1 - \gamma}$$

- In episodic task adding a constant c to all rewards does affect the total value because the sum of $\sum_{t=0}^T \gamma^t$ is finite. The agent's decision will not change because the relative difference between state values but rewards will change.

Example:

If we add $c = -1$ to every reward the agent will follow the same policy to reach the goal, but the total rewards will decrease by $\frac{1-\gamma^{T+1}}{1-\gamma}$. Also, if we add

$c = +1$, the total reward increase by same amount, even though the agent follow the same optimal path.

(a)

Q4e5 State values - center = 0.7, Up = 2.3, Right = 0.4
Down = -0.4, Left = 0.7

Probability of each action = $\frac{1}{4}$

$$\gamma = 0.9$$

The bellman eq for equiprobable random policy at center state $v_1(s)$ (eq 3.14.)

$$v_1(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_1(s')]$$

$$= \frac{1}{4} [r \times 2.3] + \frac{1}{4} [r \times 0.4] + \frac{1}{4} [r \times (-0.4)] + \frac{1}{4} [r \times 0.7]$$

~~$$= \frac{1}{4} [0.9 \times 2.3 + 0.9 \times 0.4 + 0.9 \times (-0.4) + 0.9 \times 0.7]$$~~

$$= \frac{1}{4} [0.9 \times 2.3 + 0.9 \times 0.4 + 0.9 \times (-0.4) + 0.9 \times 0.7]$$

$$= 0.675 \approx 0.7$$

This aligns with the center value given

Q-5 Part b (optimal policy)

State value = center = 17.8, Up = 18, Right = 16.4
 Down = 16, Left = 18.2

$$\gamma = 0.9$$

^{optimally}
 Bellman equation for v^* at center state.

$$v^*(s) = \max_a \left[\sum_{s', r} p(s', r | s, a) [\gamma + r + v^*(s')] \right] \dots 3.19$$

Assuming transitions are deterministic under the optimal policy and there are no explicit rewards

$$v^*(s) = \max [r \cdot v^*(s') \text{ for each action } a]$$

$$v^*(s) = \max [0.9 \times 18, 0.9 \times 16.4, 0.9 \times 16, 0.9 \times 18.2]$$

$$v^*(s) = 16.38$$

The highest value we got using calculated optimal value is 16.38 (left). This does not align with the stated value 17.8.

here I divides equation into two
to simplify

Ques:-

(a) Bellman equation for state values:-

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

$P(s'|s, a)$ → the probability of transitioning to state s' from state s after taking action a

$R(s, a, s')$: Reward received after transitioning from s to s' having action a

→ States and action for high

high: states - battery is high
low: battery is low

Actions:-

In high (A_{high}): {Search, wait}

In low (A_{low}): {Search, wait, recharge}

→ Bellman for high state

$$V_{\text{high}} = \sum_a \pi(a|\text{high}) \sum_{s'} P(s'|\text{high}, a) [R(\text{high}, a, s') + \gamma V_{\pi}(s')]$$

expand equation:

→ Action: Search

probability of taking action $\pi(\text{search}|\text{high})$

→ Transitions and rewards

Stay in high:-

probability α

Reward r_{search}

Next state value = $V^{\pi}(\text{high})$

Transition to low

Prob = $1 - \alpha$

Reward = r_{search}

Next state value = $V^{\pi}(\text{low})$

Eq:-

$$\pi(\text{search}|\text{high}) [\alpha(r_{\text{search}} + \gamma V^{\pi}(\text{high})) + (1-\alpha)(r_{\text{search}} + \gamma V^{\pi}(\text{low}))]$$

→ Action: wait

Probability of taking action $\pi(\text{wait}|\text{high})$

→ Transitions and rewards

Stays in high

Prob = $1 - \alpha$

Reward = r_{wait}

Next state value = $V^{\pi}(\text{high})$

p9

$$\pi(\text{wait}|\text{high}) [\alpha(r_{\text{wait}} + \gamma V^{\pi}(\text{high}))]$$

$$\Rightarrow V_{\text{final}}^{\pi} = \Pi(\text{Search|high}) [\alpha(\text{Search} + \gamma V^{\pi}(\text{high})) + (1-\alpha)(\text{Search} + \gamma V^{\pi}(\text{low})) + \Pi(\text{Wait|high}) (\text{Wait} + \gamma V^{\pi}(\text{high}))]$$

\rightarrow For state low

$$V^{\pi}(\text{low}) = \sum_a \Pi(a|\text{low}) \sum_{s'} P(s'| \text{low}, a) [R(s') + \gamma V^{\pi}(s')]$$

except in eq:

1) Action: search

- Probability of taking action: $\Pi(\text{Search|low})$

Transition and rewards.

• Stay in low:

probability: β

Reward: γSearch

Next state: $V^{\pi}(\text{low})$.

• Battery depletes

Prob: $1 - \beta$

Reward: -3 penalty

Next state: $V^{\pi}(\text{high})$.

$$\text{eq. } \Pi(\text{Search|low}) [\beta(\text{Search} + \gamma V^{\pi}(\text{low})) + (1-\beta)(-3 + \gamma V^{\pi}(\text{high}))]$$

2) Action: wait

Prob :- $\pi(\text{coait}(10\omega))$
 Transition & rewards
 - τ_{wait} & $r\sqrt{v_{10\omega}}$

- Stay in low:-

$$\text{Prob} = 1$$

$$\text{Reward} = \tau_{\text{wait}}$$

$$\text{Next state value: } \sqrt{\pi(\text{coait}(10\omega))} \cdot \tau_{\text{wait}} + r\sqrt{v_{10\omega}}$$

$$\pi(\text{coait}(10\omega))[\tau_{\text{wait}} + r\sqrt{v_{10\omega}}]$$

3) Action- Recharge.

Prob :- $\pi(\text{recharge}(10\omega))$ must be 1 - $\pi(\text{coait}(10\omega))$

Transition and reward.

- Transition to high

$$\text{Prob} := 1$$

$$\text{Reward} := 0$$

$$\text{Next state value: } \sqrt{\pi(\text{chigh})}$$

$$\pi(\text{recharge}(10\omega))[\tau + r\sqrt{\pi(\text{chigh})}]$$

final eq:-

$$\sqrt{v_{10\omega}} = \pi(\text{seach}(10\omega))[\beta(\tau_{\text{search}} + r\sqrt{v_{10\omega}}) +$$

$$\text{final eq} := (\alpha + C_1 - \beta)C_3 + r\sqrt{\pi(\text{chigh})} + \pi(\text{coait}(10\omega))[\tau_{\text{wait}} + r\sqrt{v_{10\omega}}] + \pi(\text{recharge}(10\omega))[\tau + r\sqrt{\pi(\text{chigh})}]$$

Ques 6

→ Part (b)

Given Parameters

$$\gamma = 0.9$$

• Transition probabilities

$$\alpha = 0.8$$

$$\beta = 0.6$$

• Reward

$$r_{\text{search}} = 10$$

$$r_{\text{wait}} = 3$$

$$\text{Rescue reward} = -3$$

• Policy

$$\pi(\text{Search} | \text{high}) = 1$$

$$\pi(\text{Wait} | \text{low}) = 0.5$$

$$\pi(\text{Recharge} | \text{low}) = 0.5$$

$$\pi(\text{Search} | \text{low}) = 0$$

→ Bellman equation for high.

→ will put into the above question find eq.

$$V_{\text{high}} = \alpha [r_{\text{search}} + \gamma V_{\text{high}}] + (1-\alpha) [r_{\text{wait}} + \gamma V_{\text{low}}]$$

$$0.28V_{\text{high}} - 0.18V_{\text{low}} = 10 \rightarrow c_1$$

→ Bellman for low

$$V_{\text{low}} = \pi(\text{wait} | \text{low}) [r_{\text{wait}} + \gamma V_{\text{low}}] + \pi(\text{recharge} | \text{low}) [0 + \gamma V_{\text{high}}]$$

$$0.55V_{\text{low}} - 0.45V_{\text{high}} = 1.5 \rightarrow c_2$$

→ multiply eq c_1 by 55 and eq c_2 by 28

$$15 \cdot 4V(\text{high}) - 9 \cdot 9V(\text{low}) = 55.0$$

$$15 \cdot 4V(\text{low}) - 12 \cdot 6V(\text{high}) = -44.2$$

After solving will get

$$V(\text{high}) \approx 79.04$$

$$V(\text{low}) \approx 67.40$$

→ Check that it satisfies the bellman equation

$$V(\text{high}) = 10 + 0.72 \times V(\text{high}) + 0.18V(\text{low}) \\ = 79.04.$$

$$V(\text{low}) = 1.5 + 0.45 \times V(\text{low}) + 0.45 \times V(\text{high}) \\ = 67.38$$

→ This shows that expected cumulative rewards for the robot starting in each state is under given policy

→ When battery is high, cumulative reward is 79.04

→ When battery is low, cumulative reward is 67.40

Extra credit

questions

→ bellman equation with θ

$$\begin{aligned} V_{\text{high}} &= \gamma_{\text{search}} + \gamma [QV_{\text{high}} + c_1 - \theta] V_{\text{low}} \\ &= 10 + 0.9(0.8 \cdot V_{\text{high}}) + 0.2 V_{\text{low}} \\ &= 10 + 0.72 V_{\text{high}} + 0.18 V_{\text{low}} \\ \theta &= 0.28 V_{\text{high}} - 0.18 V_{\text{low}} = 1.0 \rightarrow c_1 \end{aligned}$$

$$V_{\text{low}} = \theta [\gamma_{\text{wait}} + \gamma V_{\text{high}}] + c_1 - \theta [V_{\text{high}}]$$

after simplify

$$\Rightarrow V_{\text{low}}(c_1 - \theta) = 30 + 0.9(c_1 - \theta) V_{\text{high}} \rightarrow c_2$$

add V_{high} in above and then solve
we get.

$$0.28 V_{\text{high}} - 0.18 \left(\frac{30 + 0.9(c_1 - \theta) V_{\text{high}}}{1 - 0.9\theta} \right) = 10$$

→ To maximize V_{low} we consider $\theta = 0$ and $\theta = 1$

→ I) At $\theta = 0$

$$\pi(\text{credit} | \text{low}) = 0$$

$$\begin{aligned} V_{\text{low}} &= 0 + \frac{0.9 \times 1 \times V_{\text{high}}}{1 - 0} \\ &= 0.9 V_{\text{high}} \end{aligned}$$

SUBSTITUTE V_{low} into c_1

$$0.28V_{\text{high}} - 0.18 \times 0.9V_{\text{high}} = 10$$

$$V_{\text{high}} = 84.75$$

$$V_{\text{low}} = 76.28$$

(2) At $\theta = 1$, we have $V_{\text{high}} = 15 - 4P_{\text{high}}$

$$\Pi(\text{charge}) = 1 - e^{-0.2(15 - 4P_{\text{high}})}$$

Put in eq 2, we get $1 - e^{-0.2(15 - 4P_{\text{high}})} = 0$

$$V_{\text{low}} = 3 = 30$$

→ Substitute V_{low} in eq (1) $\rightarrow V_{\text{high}} = 15 - 4P_{\text{high}}$

$$V_{\text{high}} = 0.28V_{\text{high}} + 15 - 4P_{\text{high}}$$
$$V_{\text{high}} \approx 55$$

- Overall, $V_{\text{low}} < V_{\text{high}}$ so optimal $\theta = 0$
- V_{low} is higher when $\theta = 0$
- optimal policy
 $\theta = 0$; choose recharge in the low state
- value function at $\theta = 0$: $V_{\text{high}} = 84.75$

$$V_{\text{high}} = 84.75$$

$$V_{\text{low}} = 76.28$$

To maximize V_{low} , I set $\theta = 0$. So the robot will always recharge in low state.