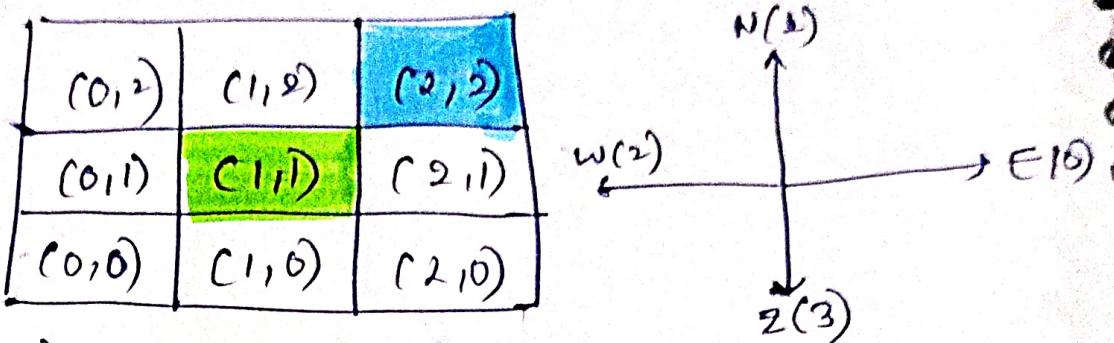


Ques 2: (b) Considering an example of Value Iteration for 3 iterations.

Let take an example of 3×3 grid



Goal state = $(2,2)$

Obstacle state = $(1,1)$

East :- $(x+1, y)$

West :- $(x-1, y)$

North :- $(x, y+1)$

South :- $(x, y-1)$

$$\gamma = 0.9$$

$$\theta \in \{0, 1, 2, 3\}$$

Rewards :-

$$\text{Step} = -1$$

$$\text{Collision} = -100$$

$$\text{goal} = +50$$

Bellman Optimality Equation :-

$$V(s) = \max_a \sum_a \pi(a|s) \left(\sum_{s'} p(s'|s, a) r + \gamma \sum_{s'} p(s'|s, a) V(s') \right),$$

$\forall s \in S$

the value function is initialized $v_0(s) = 0$ for all states, this will help us an initial estimate and allow value iteration to iteratively refine the value function using Bellman Optimality Equation.

Iteration 0:

$$v_0(s) = 0 \quad \forall s$$

Iteration 1:

since $v_0 = 0$,

$$v_1(s) = \max_a \sum_a \pi(a|s) \left(\sum_{s'} p(s'|s, a) r + \gamma \sum_{s'} p(s'|s, a) v_0(s') \right)$$

substituting $v_0(s') = 0$,

$$v_1(s) = \max_a \sum_a \pi(a|s) \left(\sum_{s'} p(s'|s, a) \gamma \right)$$

$$= \max_a \sum_{s'} p(s'|s, a) \pi \quad \begin{cases} \text{best policy is} \\ \text{deterministic \&} \\ \text{greedy,} \end{cases}$$

CASE I: near goal state

Falling Example:

$$\begin{aligned} 1) & (1, 2, 0) \\ & (x, y, 0) \end{aligned}$$

$$\begin{aligned} & Y_1 = \underbrace{0.5 \times 50}_{\text{goal}} + \underbrace{0.1 \times (-100)}_{\text{boundary}} + \underbrace{0.1 \times (-100)}_{\text{obstacle}} \\ & = 40 - 10 - 10 = 20 \end{aligned} \quad \boxed{\text{Forward Action}}$$

$$2) \quad (0, 0, 0)$$

$v_0(s)$

Reward = -1
 Deterministic (period 2)] Turn left +

$$\sum_{a_1} r \cdot p(a_1 | s_1, a) = 1(-1) = -1$$

Reward = -1
 Deterministic (period 2)] Turn right
 $\sum_{a_1} r \cdot p(a_1 | s_1, a) = 2(-1) = -2$

$$V_1(1, 2; 0) = \max(-1, -2) = -1$$

for iteration α :

$$V_2(s) = \max_a \left(\sum_{s'} r \cdot p(r | s, a) + \gamma \sum_{s'} p(s' | s, a) V_1(s') \right)$$

$$V_1(1, 2, 0) = -1$$

→ for forward action:-

$$\sum_{a_1} p(a_1 | s, a) = 0.8 \times 50 + 0.1 \times (-10) + 0.1 \times (-100) \\ = 20$$

future value term:

$$\gamma \sum_{s'} p(s' | s, a) V_1(s') = 0 \quad \begin{array}{l} \text{[all outcomes} \\ \text{are terminal]} \end{array}$$

$$Q(s, \text{forward}) = 20$$

\Rightarrow Turn left \rightarrow
turn left \rightarrow next state = $(1, 2, 1)$

From iteration 2:

$$V_1(2, 2, 1) = -1$$

Immediate reward:-

$$\sum_a r P(a|s, a) = -2$$

Future value:

$$\gamma \sum_{s'} p(s'|s, a) V_1(s') = 0.9 \times (-1) = 0.9$$

$$Q(s, \text{turn left}) = -1 - 0.9 = -1.9$$

\Rightarrow similarly,

Turn Right \rightarrow

Turn right \rightarrow next-state = $(2, 2, 3)$

Immediate reward:-

$$\sum_a r P(a|s, a) = -1$$

Future value:

$$\gamma \sum_{s'} p(s'|s, a) V_1(s') = 0.9 \times (-1) = 0.9$$

$$Q(s, \text{turn right}) = -1.9$$

$$V_2(1, 2, 0) = \max(20, -1.9, -1.9) = 20$$

Situation 3:-

$$V_2(1,2,0) = 20$$

$$V_3(s) = \max_a \left(\sum_s r p(s|s,a) + \gamma \sum_{s'} p(s'|s,a) V_2(s') \right)$$

Cx:-

$$s = (2,2,0) \quad \text{action forward}$$

immediate reward

$$\sum_a r p(a|s,a) = 0.8 \times 50 + \frac{0.1 \times -100}{x-100} + 0.1$$

all outcomes are terminal,
no future value

$$Q(\text{forward}) = 20$$

Turn Left

Turn left \rightarrow new state $= (1,2,1)$

$$V_2(1,2,1) = -1.9$$

immediate reward

$$\sum_a r p(a|s,a) = -1$$

future rewards:-

$$\gamma \sum_{s'} p(s'|s,a) V_2(s') = 0.9 \times (-1.9) \\ = -1.71$$

$$Q(\text{turn left}) = -1.71$$

Turn Right

Turn right \rightarrow new state $(1, 2, 3)$

immediate reward:-

$$\sum a p(a|s, a) = -1$$

future reward

$$\gamma \sum_{s'} p(s'|s, a) v_2(s') = 0.9 \times (-1.9) = -1.71$$

$$\text{Total} = -1 - 1.71 = -2.71$$

$$Q(\text{Turn Right}) = -2.71$$

$$V_3(1, 2, 0) = \max(-20, -2.71, -2.71) = -20$$

Case II - distant corner state
for example $s = (0, 0, 0)$

$V_0(s) = 0 \rightarrow$ initialization & iteration 0

Iteration 1

$V_0 = 0$, future term disappears.

Forward Action

$$V_1(s) = \max_a \sum_r \gamma p(r|s, a)$$

$$= 0.8 \times 1 + 0.1 \times -1 + 0.1 \times -10$$

$$= -0.8 - 0.1 - 10 = -10.9$$

$$Q(\text{forward Action}) = -10.9$$

Turn left - { Probability = $\frac{1}{2}$ } Deterministic
 $\sum_a r_p(a|s, a) = 1(-1) = -1$

$$\boxed{Q(s, \text{Turn Left}) = -1}$$

Turn right

(similarly)

$$\sum_a r_p(a|s, a) = 1(-1) = -1$$

$$\boxed{Q(s, \text{Turn Right}) = -1}$$

Total maximum

$$V_2(0, 0, 0) = \max(-11.71, -1.9, -1.9) \\ = -1.9$$

Iteration 3

$$V_2(0, 0, 0) = -1.9$$

$$V_2(1, 0, 0) = -1.9$$

$$V_2(0, 1, 0) = -1.9$$

Immediate rewards :-

$$\sum_a r_p(a|s, a) = 0.8 \times (-1) + 0.1 \times (-1) + \\ 0.1 \times (-10) \\ = -10.9$$

Future rewards :-

$$\gamma \sum_{s'} p(s'|s, a) V_2(s') = \\ 0.9 \times (0.8 \times (-1.9) + 0.1 \times (-1.9)) \\ = -1.539$$

$$Q(\text{Forward}) = -10.9 - 1.539 \\ = -12.439$$

Turn Right :-

Immediate reward:

$$\sum_r r p(r|s,a) = -1$$

future rewards:

$$\gamma \sum_{s'} p(s'|s,a) v_2(s') = 0.9(-1.9) = -1.71$$

$$Q(\text{Turn Right}) = -1 - 1.71 = -2.71$$

Turn Left :-

Immediate reward:

$$\sum_a r p(a|s,a) = -1$$

Future rewards:-

$$\gamma \sum_{s'} p(s'|s,a) v_2(s') = 0.9(-1.9) \\ = -1.71$$

$$Q(\text{Turn Left}) = -1.71 - 1 \\ = -2.71$$

Total maximum:

$$V_3(0,0,0) = \max(-12.439, -2.71, -2.71)$$

$V_3(0,0,0) = -2.71$

State	v_1	v_2	v_3	
(1, 1, 0)	20	20	20	
(0, 0, 0)	-1	-1.9	-2.7	

from this, we can observe, that states who are near goal state are stable at value 20, but corner distant states become increasingly negative due to cumulative step cost and collision risk. the near goal state having positive stable value because terminal reward dominates.

Optimal Policy

again using same example,
using the 3x3 grid:

$$\text{Goal} = (d, 2)$$

$$\text{Obstacle} = (1, 1)$$

$$\gamma = 0.9$$

$$\text{step} = 1$$

After performing value iterations until convergence, we obtain optimal value $V^*(s)$

optimal value function satisfies Bellman optimality equation:

$$V^*(s) = \max_a \left(\sum_r p(r|s,a)r + \gamma \sum_{s'} p(s'|s,a)V^*(s') \right)$$

$$\tilde{Q}^*(s,a) = \sum_r p(r|s,a)r + \gamma \sum_{s'} p(s'|s,a)V^*(s')$$

$$\boxed{V^*(s) = \max_a \tilde{Q}^*(s,a)}$$

Optimal policy:

$$\boxed{\pi^*(s) = \operatorname{argmax}_a \tilde{Q}^*(s,a)}$$

this action maximizes the expected discounted return.

Ex In 3x3 grid:

after value iterations converge, we obtained:

$$V^*(1,2,0) = 20$$

$$V^*(0,0,0) = -2.71$$

Case 1: Start S = (2, 2, 0)
(Position (2, 2), facing East)

Optimal policy :-

$$\pi^*(s) = \arg\max_a Q^*(s, a)$$

where,

$$Q^*(s, a) = \sum_r p(s'|s, a)r + \gamma \sum_{s'} p(s'|s, a)V^*(s')$$

Action: Forward

Possible transitions:

- 0.8 \rightarrow (2, 2) Goal \rightarrow reward = 50
- 0.1 \rightarrow Boundary \rightarrow reward = -100
- 0.1 \rightarrow Obstacle \rightarrow reward = -100

Since goal is terminal, future value is 0

$$Q^*(2, 2, 0, \text{forward}) = 0.8(50) + 0.1(-100) \\ + 0.1(-100) \\ = 20$$

Action: Turn Left

turning gives immediate reward
= -1

New State Value:

$$V^*(1, 2, 1) = -2.71$$

$$Q^*(1, 2, 0, \text{TurnLeft}) = -1 + 0.9(-2.71) \\ = -4 - 2.439 = -3.439$$

Action: Turn Right

similarity:

$$\theta^*(1, 2, 0, \text{Turn Right}) = -1 + 0.9(-3.71) \\ \approx -3.439$$

Compare:

$$\theta^*(\text{forward}) = 20$$

$$\theta^*(\text{Turn Left}) \approx -3.439$$

$$\theta^*(\text{Turn Right}) \approx -3.439$$

$$\boxed{\pi^*(1, 2, 0) = \text{forward}}$$

case 2

$$s = (0, 0, 0)$$

Position, $(0, 0)$, facing East

$$V^*(0, 0, 0) = -2.71$$

Action: Forward

$$\theta^*(0, 0, 0, \text{forward}) =$$

~~0.8(1 - 0.9(-2.71))~~

$$0.8(-1 + 0.9(-2.71)) + 0.1[-1 + 0.9(-2.71)]$$

$$+ 0.1(-100)$$

$$= -2.7512 - 0.3439 - 10$$

$$= -13.0951$$

Action: Turn Left

$$Q^*(0,0,0, \text{Turn Left}) \Rightarrow -1 + 0.9 V^*(0,0,1)$$
$$\approx -1 + 0.9(-2.71)$$
$$\approx -3.439$$

Action: Turn Right

$$Q^*(0,0,0, \text{Turn Right}) \Rightarrow -1 + 0.9 V^*(0,0,3)$$
$$\approx -1 + 0.9 \times (-2.71)$$
$$\approx -3.439$$

Note

$$Q^*(\text{Turn Left}) = Q^*(\text{Turn Right})$$
$$\approx -3.439$$

$$Q^*(\text{forward}) = -13.095$$

$$-3.439 > -13.095$$

$$\boxed{\pi^*(0,0,0) = \text{Turn Left or Turn Right}}$$

for manual computation of first three iterations, a reduced 3×3 grid was used to clearly demonstrate the Bellman process, the same implementation generalizes to 10×10 grid environment

c) Policy iteration consists of steps :-

(1) Policy evaluation $\Rightarrow V_{\pi^k} = r_{kk} + \gamma P_{kk} V_{\pi^k}$

(2) Policy improvement $\Rightarrow \pi_{k+1} = \arg \max_{\pi} (r_{kk} + \gamma P_{kk} V_{\pi^k})$

\Rightarrow Initializing an arbitrary policy (Assumption):

$\pi_0(s) = \text{Forward } V_s$

In every state, the agent choose Forward.

Now, policy evaluation:

For fixed policy π , value function:

$$V^\pi(s) = \sum_a p(a|s, \pi(s)) r + \gamma \sum_{s'} p(s'|s, \pi(s)) V^\pi(s')$$

equation is solved iteratively

Ex

$s = (1, 2, 0)$ in grid size (3×3)

under policy π_0 , action = forward

transitions :-

0.8 \rightarrow Goal \rightarrow reward = 50

0.1 \rightarrow Boundary \rightarrow reward = -100

0.1 \rightarrow Obstacle \rightarrow reward = -100

size, goal & terminal

$$V^{\pi_0}(1, 2, 0) = 0.8(50) + 0.1(-100) + 0.1(-100)$$

$$\approx 20$$

$$\boxed{V^{\pi_0}(1, 2, 0) \approx 20}$$

state $c = (0, 0, 0)$

Transitions

$$0.8 \rightarrow (1, 0) \rightarrow \text{reward} = -1$$

$$0.1 \rightarrow (1, 0) \rightarrow \text{reward} = -1$$

$$0.1 \rightarrow \text{collision} \rightarrow \text{reward} = -100$$

$$V^{\pi_0}(0, 0, 0) = 0.8[-1 + 0.9 V^{\pi_0}(1, 0, 0)] + 0.1[-1 + 0.9 V^{\pi_0}(0, 1, 1)] \\ + 0.1(-100)$$

$$= 0.8(-1) + 0.1(-1) + 0.1(-100)$$

$$= -0.8 - 0.1 - 10$$

$$= -10.9$$

$$V^{\pi_0}(0, 0; 0) = -10.9$$

Policy Improvement

$$Q^\pi(s, a) = \sum_r p(r|s, a)r + \gamma \sum_{s'} p(s'|s, a)V^\pi(s')$$

we compare all actions and choose the best.

we already computed,

$$Q^{\pi_0}(\text{forward}) = -10.9$$

Turn-left:

Turning gives, New state $\approx (0, 0, 1)$

Reward ≈ -1

probability ≈ 1

$$Q^{\pi_0}(\text{TurnLeft}) = -1 + 0.9 V^{\pi_0}(0, 0, 1)$$

$$\text{Initially, } V^{\pi_0}(0, 0, 1) \approx 0$$

$$= -1 + 0.9(0) = -1$$

Similarly,

Turn Right:

$$Q^{\pi_0}(\text{TurnRight}) = -1 + 0.9(0) = -1$$

Now,

$$-1 > 10.9$$

policy improves,

$$\pi_1(0, 0, 0) = \text{Turn Left (or Turn Right)}$$

second policy evaluation
assuming, $\pi_1(0, 0, 0) \rightarrow \text{turn left}$

$$V^{\pi_1}(0, 0, 0) = -1 + 0.9 V^{\pi}(0, 0, 1)$$

If $(0, 0, 1)$ also turn similarly:

Eventually values propagate:

$$V^{\pi_1}(0, 0, 0) = -1 + 0.9 (-1 + 0.9 (-1))$$
$$= -1 - 0.9 - 0.81$$

$$\text{ft } 8: \quad -1(-1 + 0.9 + 0.9^2)$$

Final improvement check:
now encompassing forward using updated
values;

$$Q(\text{forward}) = 0.8[-1 + 0.9(-2.71)] +$$
$$0.1[-1 + 0.9(-2.71)] + 0.1(-100)$$

$$\text{backward computed} = -13.095$$

$$\text{turning gives } -1 + 0.9(-2.71) = -3.439$$

$\delta_{\text{in}0}$:

$$-3.439 > -13.095$$

Policy remain same \rightarrow stable

final optimal policy:

$\pi^*(1, 2, 0) = \text{forward}$

$\pi^*(0, 0, 0) = \text{Turn Right or Turn Left}$