

# EX1.

1) For  $V^\pi(s)$  to be  $-\infty$ , the agent should not reach any terminal state.

Let the policy  $\pi(s_t) = \begin{cases} \text{right}, & s_t = 3 \\ \text{left}, & s_t \neq 3 \end{cases}$

Hence,  $V^\pi(3) = -\infty$  since the agent will oscillate between grid 3 and 4.

2)  $Q^\pi(9, \text{right}) = -1 + V(10) = -1$ , since 10 is the terminal state,  $V(10)=0$ , one step reward is -1.

$$Q^\pi(12, \text{left}) = -1 + V(11) = -1 + (-13.39649575) \approx -14.396$$

3) ① Initialize  $V(s) = -1$  for  $s \in S - \{1, 10\}$ , and  $V(s) = 0$  for  $s \in \{1, 10\}$   
 $\pi(s)$  arbitrary for all  $s \in S - \{1, 10\}$

② Loop:

$$\Delta \leftarrow 0$$

Loop for each  $s \in S$ :

If  $s$  is terminal state ( $s = 1/10$ )

continue

$$V \leftarrow V_{\text{old}}$$

$$V(s) \leftarrow -1 + \sum_{a \in A} \pi(a|s) * V(s') \quad (\pi(a|s) = \frac{1}{4} \text{ for all } a \in A \text{ and } s' \text{ is the next step from } s \text{ taking } a)$$

$$\Delta \leftarrow \max(\Delta, |V - V_{\text{old}}|)$$

Until  $\Delta < \theta$  ( $\theta$  a small positive number determining the accuracy of estimation)

③ policy-stable  $\leftarrow$  true

For each state  $s \in S$  except terminal states:

If  $s$  is terminal state ( $s = 1/10$ )

continue

$$\text{old\_action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_{a \in A} (-1 + \sum_{s' \in S} p(s'|s, a) * V(s'))$$

If  $\text{old\_action} \neq \pi(s)$ , then  $\text{policy-stable} \leftarrow \text{false}$

If  $\text{policy-stable}$ , then stop and return  $V$  &  $\pi$  and  $\pi \leftarrow \pi_*$

else go to ②

4)  $V^\pi$ :

0	-1	-2	-3
-1	-1	-2	-3
-1	0	-1	-2
-2	-1	-2	-3

$\pi$ :

X	$\leftarrow$	$\leftarrow$	$\leftarrow$
$\uparrow$	$\downarrow$	$\downarrow$	$\leftarrow$
$\rightarrow$	X	$\leftarrow$	$\leftarrow$
$\uparrow$	$\uparrow$	$\uparrow$	$\uparrow$

"X": terminal state

The arrow points to the optimal move direction

# EX3.

Policy Iteration :

① Initialization

$$Q(s, a) = 0 \text{ for all } s \in S \text{ and } a \in A$$

$\pi(s) \in A(s)$  arbitrarily for all  $s \in S$

② Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each  $s \in S$ :

Loop for each  $a \in A$ :

$$\text{old\_q} \leftarrow Q(s, a)$$

$$Q(s, a) \leftarrow \sum_{s', r} p(s'|s, a) (r + \gamma \sum_{a'} \pi(a'|s') Q(s', a'))$$

$$\Delta \leftarrow \max(\Delta, |\text{old\_q} - Q(s, a)|)$$

Until  $\Delta < \theta$  ( $\theta$  a small positive number determining the accuracy of estimation)

### ③ Policy Improvement

$\text{policy-stable} \leftarrow \text{true}$

For each  $s \in S$ :

$$\text{old\_action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_{a \in A} Q(s, a)$$

If  $\text{old\_action} \neq \pi(s)$ , then  $\text{policy-stable} \leftarrow \text{false}$

If  $\text{policy-stable}$ , then stop and return  $Q \approx q_*$  and  $\pi \approx \pi_*$

Else go to ②

Value Iteration:

Initialization

$$Q(s, a) = 0 \text{ for all } s \in S \text{ and } a \in A$$

Loop:

$$\Delta \leftarrow 0$$

Loop for each  $s \in S$ :

Loop for each  $a \in A$ :

$$\text{old\_q} \leftarrow Q(s, a)$$

$$Q(s, a) \leftarrow \sum_{s', r} p(s'|s, a) (r + \gamma \max_{a' \in A} Q(s', a'))$$

$$\Delta \leftarrow \max(\Delta, |\text{old\_q} - Q(s, a)|)$$

Until  $\Delta < \theta$  ( $\theta$  a small positive number determining the accuracy of estimation)

For each  $s \in S$ :

$$\pi(s) \leftarrow \operatorname{argmax}_{a \in A} Q(s, a)$$

Return  $Q$  and  $\pi$