



### Problem 1.

Consider the following system with two states  $s_k \in \{s^1 = 0, s^2 = 1\}$ .

There are two possible actions:  $a^1$  and  $a^2$ . The transition probabilities can be expressed as:

$$p(s'|s, a^1) \begin{cases} 1 & s = 0, s' = 0 \\ 0 & s = 0, s' = 1 \\ 0 & s = 1, s' = 0 \\ 1 & s = 1, s' = 1 \end{cases} \quad p(s'|s, a^2) \begin{cases} 0 & s = 0, s' = 0 \\ 1 & s = 0, s' = 1 \\ 1 & s = 1, s' = 0 \\ 0 & s = 1, s' = 1 \end{cases}$$

Reward function is as follows:  $\begin{cases} \text{moving to state } s^2: +1 \\ \text{moving to state } s^1: 0 \\ \text{action } a^1 \text{ and } a^2: 0 \end{cases}$

Start with a random policy  $\pi^0(s^1) = a^1, \pi^0(s^2) = a^1, \gamma = 0.9, \theta = 0.85$ . Use Policy Iteration to compute  $\pi^1(s^1), \pi^1(s^2)$ . Use  $V_0(s^1) = V_0(s^2) = 0$ , for initialization of Policy Evaluation.

Solution:  $\pi^0(s) = \pi^0(s^2) = a^1$

$$V_{k+1}(s) = \sum_{s' \in \mathcal{S}, s^2 \in \mathcal{S}} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k(s')]$$

$$V_{k+1}(s^1) = \underbrace{P(s^1|s^1, \pi(s^1))}_{\substack{1 \\ \text{a}^1}} [R(s^1, \pi(s^1), s^1) + \gamma V_k(s^1)] \\ + \underbrace{P(s^2|s^1, \pi(s^1))}_{\substack{0 \\ \text{a}^1}} [R(s^1, \pi(s^1), s^2) + \gamma V_k(s^2)]$$

$$V_{k+1}(s^1) = R(s^1, a^1, s^1) + \gamma V_k(s^1)$$

$$\text{For } s^1 \rightarrow V_{k+1}(s^2) = \underbrace{P(s^1|s^2, \pi(s^2))}_{\substack{0 \\ \text{a}^1}} [R(s^2, \pi(s^2), s^1) + \gamma V_k(s^1)] \\ + \underbrace{P(s^2|s^2, \pi(s^2))}_{\substack{1 \\ \text{a}^1}} [R(s^2, \pi(s^2), s^2) + \gamma V_k(s^2)]$$

$$V_{k+1}(s^2) = R(s^2, a^1, s^2) + \gamma V_k(s^2)$$

$$V_0(s^1) = V_0(s^2) = 0$$

$$\begin{cases} V_1(s^1) = R(s^1, a^1, s^1) + \gamma V_0(s^1) = 0 + \gamma \cdot 0 = 0 \\ V_1(s^2) = R(s^2, a^1, s^2) + \gamma V_0(s^2) = 1 + \gamma \cdot 0 = 1 \end{cases}$$

$$\max_{s \in \mathcal{S}} |V_1(s) - V_0(s)| = \max \{ |V_1(s^1) - V_0(s^1)|, |V_1(s^2) - V_0(s^2)| \} \\ = \max \{ 0, 1 \} = 1 > \theta$$

$$V_2(s^1) = R(s^1, a^1, s^1) + \gamma V_1(s^1) = 0 + \gamma \cdot 0 = 0$$

$$V_2(s^2) = R(s^2, a^1, s^2) + \gamma V_1(s^2) = 1 + \gamma \cdot 1 = 1.9$$

$$\max \{ |V_2(s^1) - V_1(s^1)|, |V_2(s^2) - V_1(s^2)| \}$$

$$= \max \{ 0, 0.9 \} = 0.9 > \theta \leftarrow 0.85 \quad \times$$

$$V_3(s^1) = R(s^1, a^1, s^1) + \gamma V_2(s^1) = 0 + \gamma 0 = 0$$

$$V_3(s^2) = R(s^2, a^1, s^2) + \gamma V_2(s^2) = 1 + \gamma 0.9 = 2.7$$

$$\max \{ |V_3(s^1) - V_2(s^1)|, |V_3(s^2) - V_2(s^2)| \}$$

$$= \max \{ 0, 0.81 \} = 0.81 < \theta \leftarrow 0.85 \quad \checkmark$$

Policy Evaluation Step stops.

$$V^{\pi^0}(s^1) = 0, V^{\pi^0}(s^2) = 2.7$$

Policy Improvement:

using latest  $V(s)$  obtained from policy Evaluation, we have

$$\pi^1(s) = \operatorname{argmax}_{s' \in S} \sum_{s' \in S} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$$

$$\begin{aligned} \pi^1(s^1) = \operatorname{argmax} & \left\{ \underbrace{P(s^1|s^1, a^1)}_1 \left[ \underbrace{R(s^1, a^1, s^1)}_0 + \underbrace{\gamma V(s^1)}_0 \right] + \underbrace{P(s^2|s^1, a^1)}_0 \left[ \underbrace{R(s^1, a^1, s^2)}_1 + \underbrace{\gamma V(s^2)}_{2.7} \right] \right. \\ & \left. + \underbrace{P(s^1|s^1, a^2)}_0 \left[ \underbrace{R(s^1, a^2, s^1)}_0 + \underbrace{\gamma V(s^1)}_0 \right] + \underbrace{P(s^2|s^1, a^2)}_1 \left[ \underbrace{R(s^1, a^2, s^2)}_1 + \underbrace{\gamma V(s^2)}_{2.7} \right] \right\} \end{aligned}$$

$$\pi^1(s^1) = \operatorname{argmax} \{ 0, 3.43 \} = a^2$$

$$\begin{aligned} \pi^1(s^2) = \operatorname{argmax} & \left\{ \underbrace{P(s^1|s^2, a^1)}_0 \left[ \underbrace{R(s^2, a^1, s^1)}_0 + \underbrace{\gamma V(s^1)}_0 \right] + \underbrace{P(s^2|s^2, a^1)}_1 \left[ \underbrace{R(s^2, a^1, s^2)}_1 + \underbrace{\gamma V(s^2)}_{2.7} \right] \right. \\ & \left. + \underbrace{P(s^1|s^2, a^2)}_1 \left[ \underbrace{R(s^2, a^2, s^1)}_0 + \underbrace{\gamma V(s^1)}_0 \right] + \underbrace{P(s^2|s^2, a^2)}_0 \left[ \underbrace{R(s^2, a^2, s^2)}_1 + \underbrace{\gamma V(s^2)}_{2.7} \right] \right\} \end{aligned}$$

$$\pi^1(s^2) = \operatorname{argmax} \{ 3.43, 0 \} = a^1$$

**Problem 2.**

Consider the problem defined in Problem 1.

- a) Given  $\begin{bmatrix} V_0(s^1) \\ V_0(s^2) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  and  $\gamma = 0.9$ , perform Value Iteration method to compute  $V_1, V_2, V_3$ .
- b) Compute  $\pi(s = 0)$  and  $\pi(s = 1)$  associated with  $V_3$ .

Solution:

$$a) V_0(s^1) = V_0(s^2) = 0$$

$$V_{k+1}(s^1) = \max_{a \in \{a^1, a^2\}} \left[ \sum_{s^2 \in S} R(s^1, a, s^2) + \gamma P(s^2 | s^1, a) V_k(s^2) \right]$$

$$V_{k+1}(s^1) = \max_{a \in A} \left[ P(s^1 | s^1, a) (R(s^1, a, s^1) + \gamma V_k(s^1)) + P(s^2 | s^1, a) (R(s^1, a, s^2) + \gamma V_k(s^2)) \right]$$

$$V_{k+1}(s^1) = \max \left\{ \left[ P(s^1 | s^1, a^1) (R(s^1, a^1, s^1) + \gamma V_k(s^1)) + P(s^2 | s^1, a^1) (R(s^1, a^1, s^2) + \gamma V_k(s^2)) \right], \right. \\ \left. \left[ P(s^1 | s^1, a^2) (R(s^1, a^2, s^1) + \gamma V_k(s^1)) + P(s^2 | s^1, a^2) (R(s^1, a^2, s^2) + \gamma V_k(s^2)) \right] \right\}$$

$$V_{k+1}(s^1) = \max \{ \underbrace{R(s^1, a^1, s^1)}_0 + \gamma \underbrace{V_k(s^1)}_0, \underbrace{R(s^1, a^2, s^2)}_0 + \gamma \underbrace{V_k(s^2)}_0 \}$$

$$V_{k+1}(s^2) = \max \{ \underbrace{R(s^2, a^1, s^2)}_1 + \gamma \underbrace{V_k(s^2)}_0, \underbrace{R(s^2, a^2, s^1)}_0 + \gamma \underbrace{V_k(s^1)}_0 \}$$

$$V_1(s^1) = \max \{ \underbrace{R(s^1, a^1, s^1)}_0 + \gamma \underbrace{V_0(s^1)}_0, \underbrace{R(s^1, a^2, s^2)}_0 + \gamma \underbrace{V_0(s^2)}_0 \} = 1$$

$$V_1(s^2) = \max \{ \underbrace{R(s^2, a^1, s^2)}_1 + \gamma \underbrace{V_0(s^2)}_0, \underbrace{R(s^2, a^2, s^1)}_0 + \gamma \underbrace{V_0(s^1)}_0 \} = 1$$

$$V_2(s^1) = \max \{ \underbrace{R(s^1, a^1, s^1)}_0 + \gamma \underbrace{V_1(s^1)}_1, \underbrace{R(s^1, a^2, s^2)}_0 + \gamma \underbrace{V_1(s^2)}_1 \} = 1.9$$

$$V_2(s^2) = \max \{ \underbrace{R(s^2, a^1, s^2)}_1 + \gamma \underbrace{V_1(s^2)}_1, \underbrace{R(s^2, a^2, s^1)}_0 + \gamma \underbrace{V_1(s^1)}_1 \} = 1.9$$

$$V_3(s^1) = \max \left\{ \underbrace{R(s^1, a^1, s^1)}_0 + \underbrace{\gamma V_2(s^1)}_{1.9}, \underbrace{R(s^1, a^2, s^2)}_1 + \underbrace{\gamma V_2(s^2)}_{1.9} \right\} = 2.71$$

$$V_3(s^2) = \max \left\{ \underbrace{R(s^2, a^1, s^2)}_1 + \underbrace{\gamma V_2(s^2)}_{1.9}, \underbrace{R(s^2, a^2, s^1)}_0 + \underbrace{\gamma V_2(s^1)}_{1.9} \right\} = 2.71$$

b)

$$\pi(s) = \arg \max_{a \in \{a^1, a^2\}} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$$

$$\pi(s^1) = \arg \max_{a \in \{a^1, a^2\}} \left[ \sum_{s' \in S} R(s^1, a, s') + \gamma P(s'|s^1, a) V_R(s') \right]$$

$$\pi(s^1) = \arg \max \left\{ P(s^1|s^1, a^1) (R(s^1, a^1, s^1) + \gamma V_R(s^1)) + P(s^2|s^1, a^1) (R(s^1, a^1, s^2) + \gamma V_R(s^2)), \right. \\ \left. P(s^1|s^1, a^2) (R(s^1, a^2, s^1) + \gamma V_R(s^1)) + P(s^2|s^1, a^2) (R(s^1, a^2, s^2) + \gamma V_R(s^2)) \right\}$$

$$\pi(s^1) = \arg \max \left\{ \underbrace{R(s^1, a^1, s^1)}_0 + \underbrace{\gamma V(s^1)}_{2.71}, \underbrace{R(s^1, a^2, s^2)}_1 + \underbrace{\gamma V(s^2)}_{2.71} \right\} = a^2$$

Similarly :

$$\pi(s^2) = \arg \max \left\{ \underbrace{R(s^2, a^1, s^2)}_1 + \underbrace{\gamma V(s^2)}_{2.71}, \underbrace{R(s^2, a^2, s^1)}_0 + \underbrace{\gamma V(s^1)}_{2.71} \right\} = a^1$$

### Problem 3.

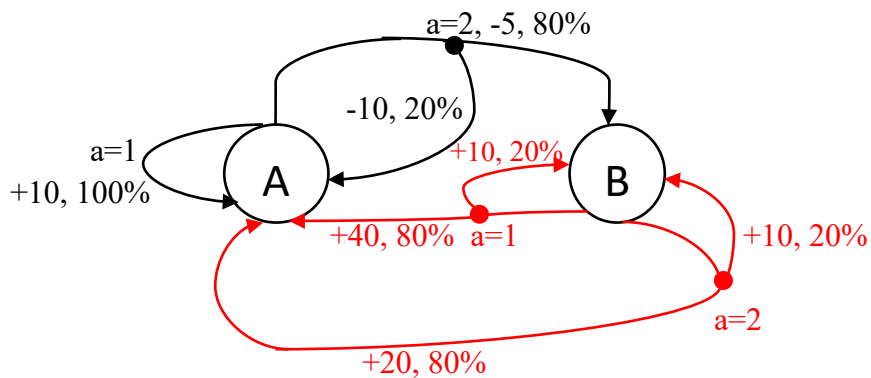
Consider the following MDP having two states: A, B. In each state, there are two possible actions: 1 and 2. The transition model and reward are shown in the diagram below.

Apply Policy Iteration to determine the optimal policy and state values of A and B.

Assume the initial policy is action 2 for both states,  $\gamma = 0.9$ .

For evaluation of policy, you need to solve two set of linear equations for the following form, instead of iterative steps of policy evaluation:

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



\*Here is an example of transition and reward from the diagram:

In state A, action 2 moves the agent to state B with probability 0.8 with the corresponding reward -5, and make the agent stay at state A with probability 0.2 and corresponding reward -10.

Solution:

$$V^\pi(s) = \sum P(s', r | s, a) [r + \gamma V^\pi(s')]$$

Policy Evaluation for  $\pi^0$

$$\begin{cases} V^\pi(A) = 0.8 [-5 + 0.9 V^\pi(B)] + 0.2 [-10 + 0.9 V^\pi(A)] \\ V^\pi(B) = 0.8 [+20 + 0.9 V^\pi(A)] + 0.2 [+10 + 0.9 V^\pi(B)] \end{cases}$$

$$\begin{aligned} 0.82 V^\pi(A) - 0.72 V^\pi(B) &= -6 \\ -0.72 V^\pi(A) + 0.82 V^\pi(B) &= 18 \end{aligned} \rightarrow \begin{cases} V^\pi(A) = 52.2 \\ V^\pi(B) = 67.8 \end{cases}$$

Policy Improvement for  $\pi^0$

$$\pi'(A) = \operatorname{argmax} \left\{ \overbrace{1 [-10 + 0.9 V^{\pi^0}(A)]}^{Q=1, 58.9}, \underbrace{0.8 [-5 + 0.9 V^{\pi^0}(B)] + 0.2 [-10 + 0.9 V^{\pi^0}(A)]}_{Q=2, 52.2} \right\} = 1$$

$$\pi'(B) = \operatorname{argmax} \left\{ \overbrace{0.8 [40 + 0.9 V^{\pi^0}(A)] + 0.2 [+10 + 0.9 V^{\pi^0}(B)]}^{Q=1, 83.7}, \underbrace{0.8 [20 + 0.9 V^{\pi^0}(A)] + 0.2 [10 + 0.9 V^{\pi^0}(B)]}_{Q=2, 67} \right\} = 1$$

$\pi' \neq \pi^0 \rightarrow \text{Continue}$



$$\pi'(A) = \pi'(B) = 1$$

Policy Evaluation for  $\pi'$

$$\begin{cases} V^{\pi'}(A) = 1 [10 + 0.9 V^{\pi'}(A)] \\ V^{\pi'}(B) = 0.8 [40 + 0.9 V^{\pi'}(A)] + 0.2 [10 + 0.9 V^{\pi'}(B)] \end{cases}$$

$$\begin{aligned} 0.1 V^{\pi'}(A) &= 10 \\ -0.72 V^{\pi'}(A) + 0.82 V^{\pi'}(B) &= 34 \end{aligned} \rightarrow \begin{cases} V^{\pi'}(A) = 100 \\ V^{\pi'}(B) = 129.2 \end{cases}$$

Policy Improvement for  $\pi'$

$$\pi^2(A) = \operatorname{argmax} \left\{ \overbrace{1 [10 + 0.9 V^{\pi'}(A)]}^{a=1, 100}, \underbrace{0.8 [-5 + 0.9 V^{\pi'}(B)] + 0.2 [-10 + 0.9 V^{\pi'}(A)]}_{a=2, 105} \right\} = 2$$

$$\pi^2(B) = \operatorname{argmax} \left\{ \overbrace{0.8 [40 + 0.9 V^{\pi'}(A)] + 0.2 [10 + 0.9 V^{\pi'}(B)]}^{a=1, 129.2}, \underbrace{0.8 [20 + 0.9 V^{\pi'}(A)] + 0.2 [10 + 0.9 V^{\pi'}(B)]}_{a=2, 113} \right\} = 1$$

$$\pi^2 \neq \pi' \rightarrow \text{continue}$$

$$V^* = V_5$$

$$\pi^* = \underset{a \in A}{\operatorname{argmax}} \overset{\text{Row-Wise}}{R^a} + \gamma M(a) V^*$$

$$\pi^* = \operatorname{argmax} \left\{ \begin{bmatrix} 0.1 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 0.2 & 0.8 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2.2738 \\ 3.0431 \end{bmatrix}, \begin{bmatrix} -1 \\ 0.8 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 2.2738 \\ 3.0431 \end{bmatrix} \right\}$$

$$= \operatorname{argmax} \left\{ \begin{bmatrix} 2.7003 \\ 3.2388 \end{bmatrix}, \begin{bmatrix} 1.0464 \\ 3.4696 \end{bmatrix} \right\} = \begin{bmatrix} a^1 \\ a^2 \end{bmatrix}$$

#### Problem 4.

Consider the following maze with 14 states and a goal. The agent can take one of the following four actions at any given state  $A = \{UP, Down, Right, Left\}$ . The state transitions are deterministic; for example  $P(S'=10 | S=12, a=U) = 1$ .

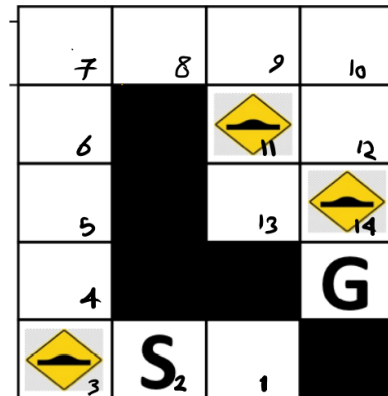
The reward is as follows:

$$\begin{cases} -1 & \text{taking any action} \\ +20 & \text{moving to goal} \\ -10 & \text{moving to bump} \end{cases}$$

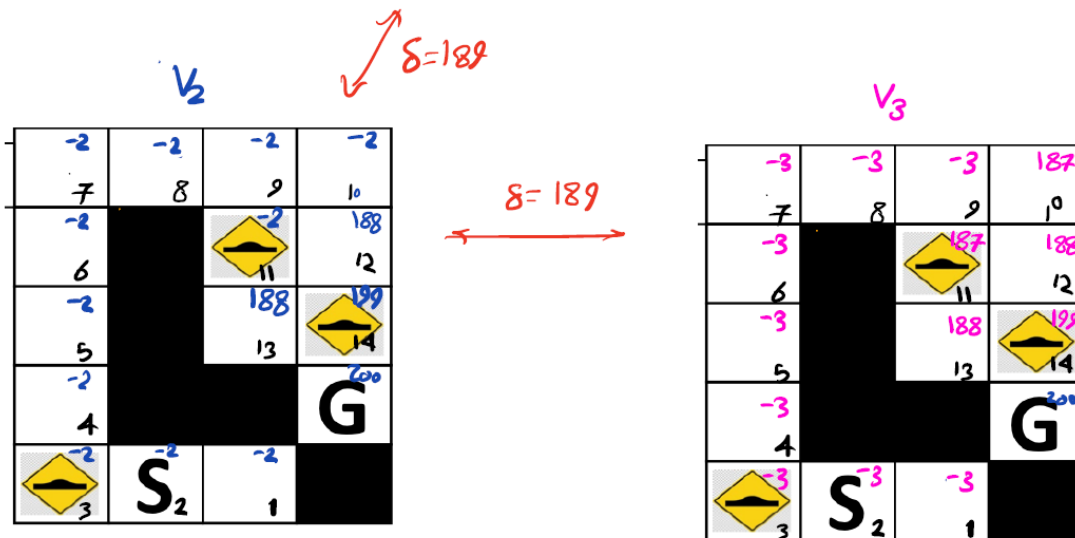
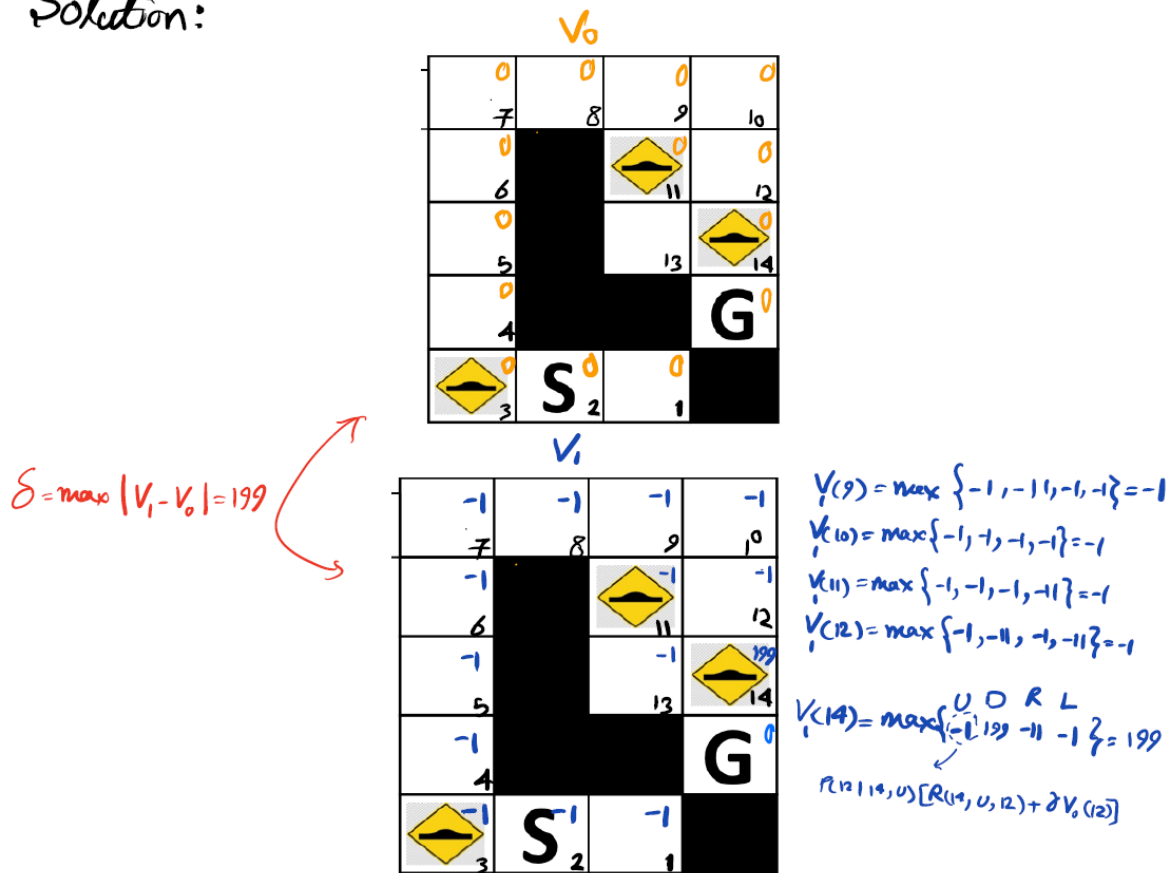
a) Using  $\gamma=1$  and  $\theta=0.5$ , perform **Vector-form** Value Iteration method with  $V_0(s)=0$  to compute  $V^*$ .

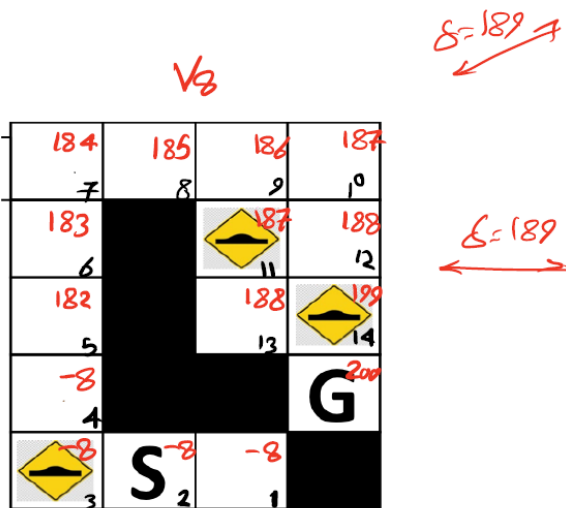
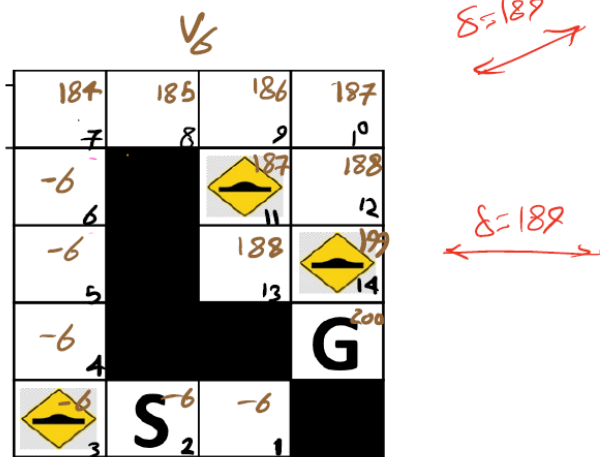
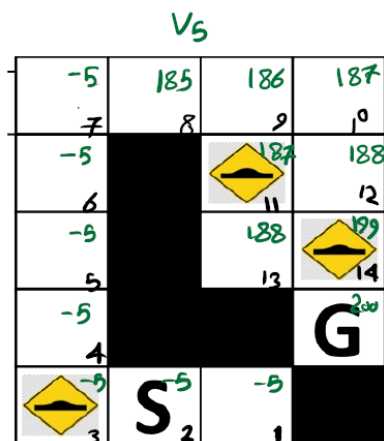
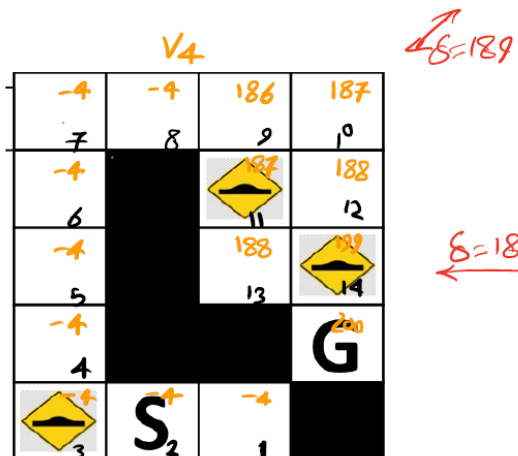
b) Compute the optimal Policy.

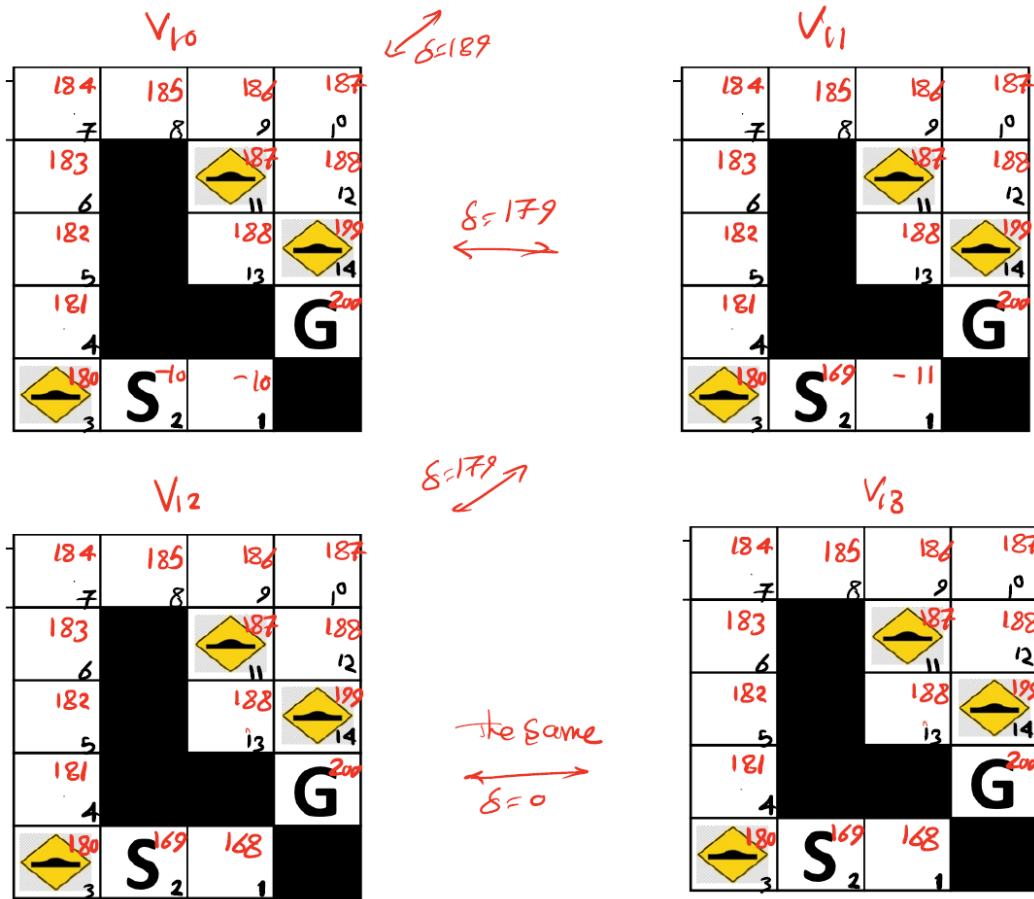
\* Show all intermediate state values in maze, without details of calculation.



Solution:







$$V_{13} = V^* \rightarrow \pi^*(s) = \operatorname{argmax}_{s'} \sum_{s''} P(s''|s, a) [R(s, a, s'') + \gamma V^*(s'')]$$



**Problem 5.**

For an MDP defined by state space  $S$ , action space  $A$ , reward  $R(s, a, s')$  and transition probability  $p(s'|s, a)$ , write the following:

a) For a given policy  $\pi$ , write

- $V^\pi(s)$  based on  $V^\pi$
- $V^\pi(s)$  based on  $Q^\pi$
- $Q^\pi(s, a)$  based on  $V^\pi$
- $Q^\pi(s, a)$  based on  $Q^\pi$

b) For the optimal policy  $\pi^*$ , write

- $V^*(s)$  based on  $V^*$
- $V^*(s)$  based on  $Q^*$
- $Q^*(s, a)$  based on  $V^*$
- $Q^*(s, a)$  based on  $Q^*$

- An example of response:

$$v^\pi(s) = \sum_{s'} p(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

Solution:

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

$$V^{\pi}(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma Q^{\pi}(s', \pi(s))]$$

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

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

$$b) V^*(s) = \max_{a \in A} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^*(s')]$$

$$V^*(s) = \max_{a \in A} \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a' \in A} Q^*(s', a')]$$

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^*(s')]$$

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a' \in A} Q^*(s', a')]$$