

# Reinforcement Learning Training 2025

# Bellman's Equation

Foundation to solving MDP and RL problems.

## Recall (1)

- Markov decision process has transition probabilities

$$\Pr \left[ S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a \right]$$

which transitions the agent to state  $S_{t+1}$  and a reward of  $R_{t+1}$

- Cumulative reward at time  $t$

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$$

## Recall (2)

- A value function is an expected cumulative reward

$$v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$$

- An action-value function is an expected cumulative reward from taking action  $a$

$$q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$$

■ Note that  $v$  and  $q$  depend on the policy  $\pi$ .

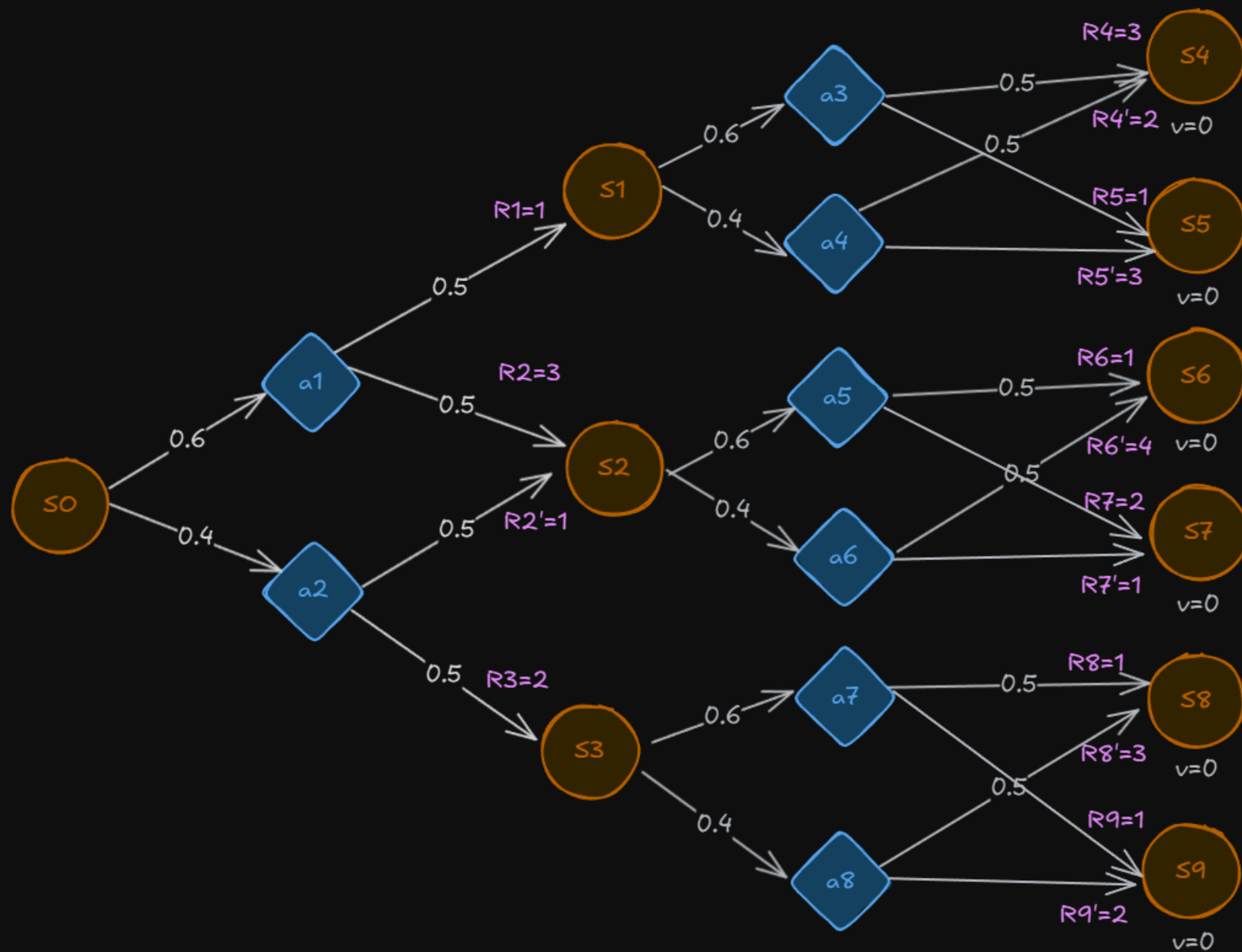
# Bellman Equation

- Allows relationships among  $v$  and  $q$ .

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

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \sum_{a'} \pi(a'|s') q(s', a') \right]$$

# Example



**Find**  $v(s_1)$

$$V(s_1) = 0.6 \left[ 0.5 \times (R_4 + \gamma \cancel{V(s_4)}) + 0.5 \times (R_5 + \gamma \cancel{V(s_5)}) \right] \\ + 0.4 \left[ 0.5 \times (R'_4 + \gamma \cancel{V(s_4)}) + 0.5 \times (R'_5 + \gamma \cancel{V(s_5)}) \right]$$

$$V(s_1) = 0.6 \left[ 0.5(3) + 0.5(1) \right] + 0.4 \left[ 0.5(2) + 0.5(3) \right] \\ = 2.2$$

**Find  $v(s_2)$  and  $v(s_3)$**

$$V(s_2) = 0.6 [ 0.5(1) + 0.5(2) ] + 0.4 [ 0.5(4) + 0.5(1) ] = 1.9$$
$$V(s_3) = 0.6 [ 0.5(1) + 0.5(1) ] + 0.4 [ 0.5(3) + 0.5(2) ] = 1.6$$

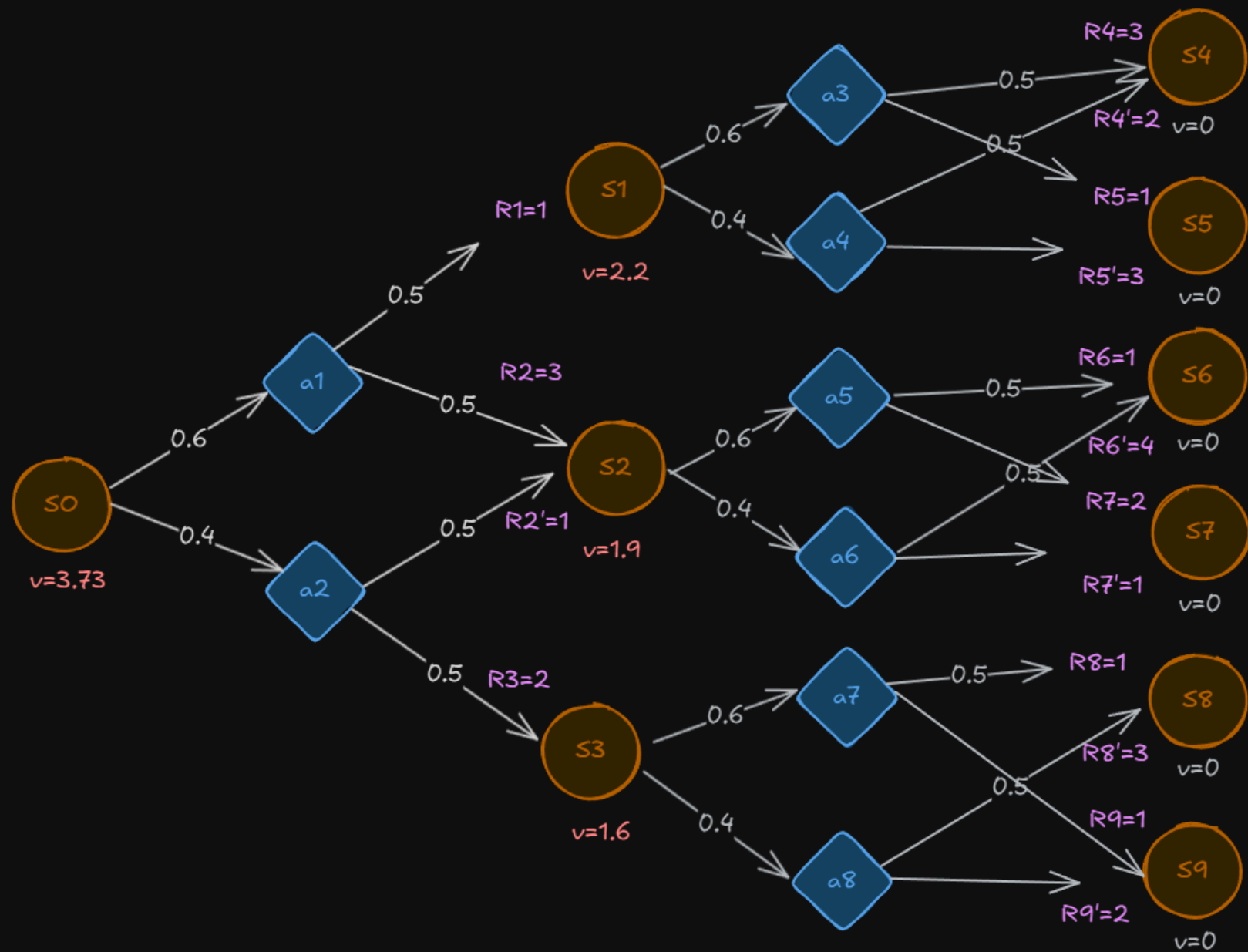


**Find**  $v(s_0)$

$$\begin{aligned} V(s_0) &= 0.6 [ 0.5(1 + 2.2) + 0.5(3 + 1.9) ] \\ &\quad + 0.4 [ 0.5(1 + 1.9) + 0.5(2 + 1.6) ] \\ &= 0.6(4.25) + 0.4(3.25) \\ &= 3.93 \end{aligned}$$

# Result

See values of  $v$



**Find**  $q(a_3)$

$$\begin{aligned} q(a_3) &= 0.5 \left[ R_4 + \gamma \sum_{a'} \cancel{(\dots)} \right] + 0.5 \left[ R_5 + \gamma \sum_{a'} \cancel{(\dots)} \right] \\ &= 0.5(3) + (0.5)(1) \\ &= 2 \end{aligned}$$

*Handwritten notes in orange:*  
An arrow points from the text "0 (No action)" to the first summation term.  
An arrow points from a "0" to the second summation term.

**Find**  $q(a_4) - q(a_8)$

$$q(a_4) = 0.5(2) + 0.5(3) = 2.5$$

$$q(a_5) = 0.5(1) + 0.5(2) = 1.5$$

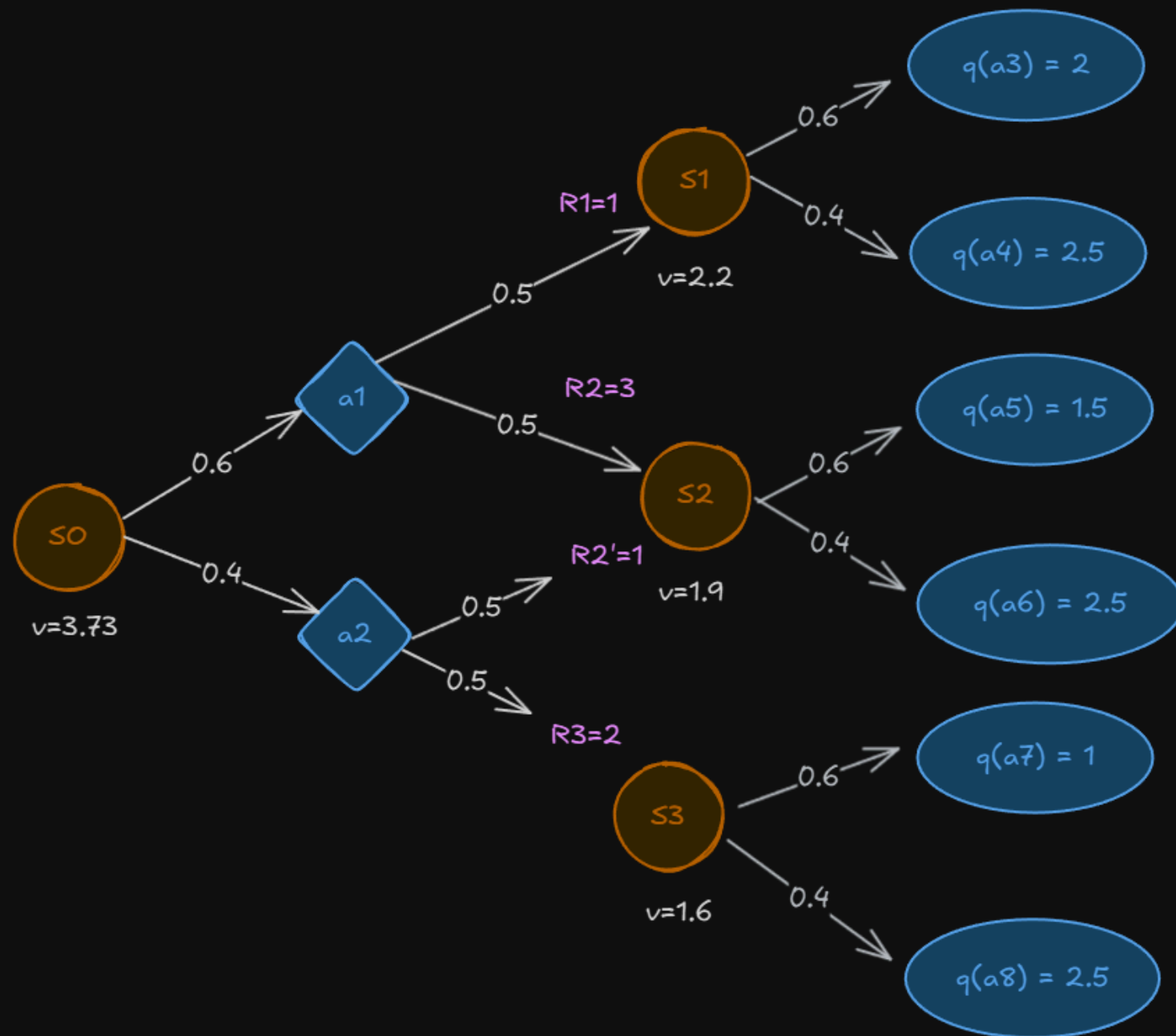
$$q(a_6) = 0.5(4) + 0.5(1) = 2.5$$

$$q(a_7) = 0.5(1) + 0.5(1) = 1$$

$$q(a_8) = 0.5(3) + 0.5(2) = 2.5$$

# Result

See  $q$

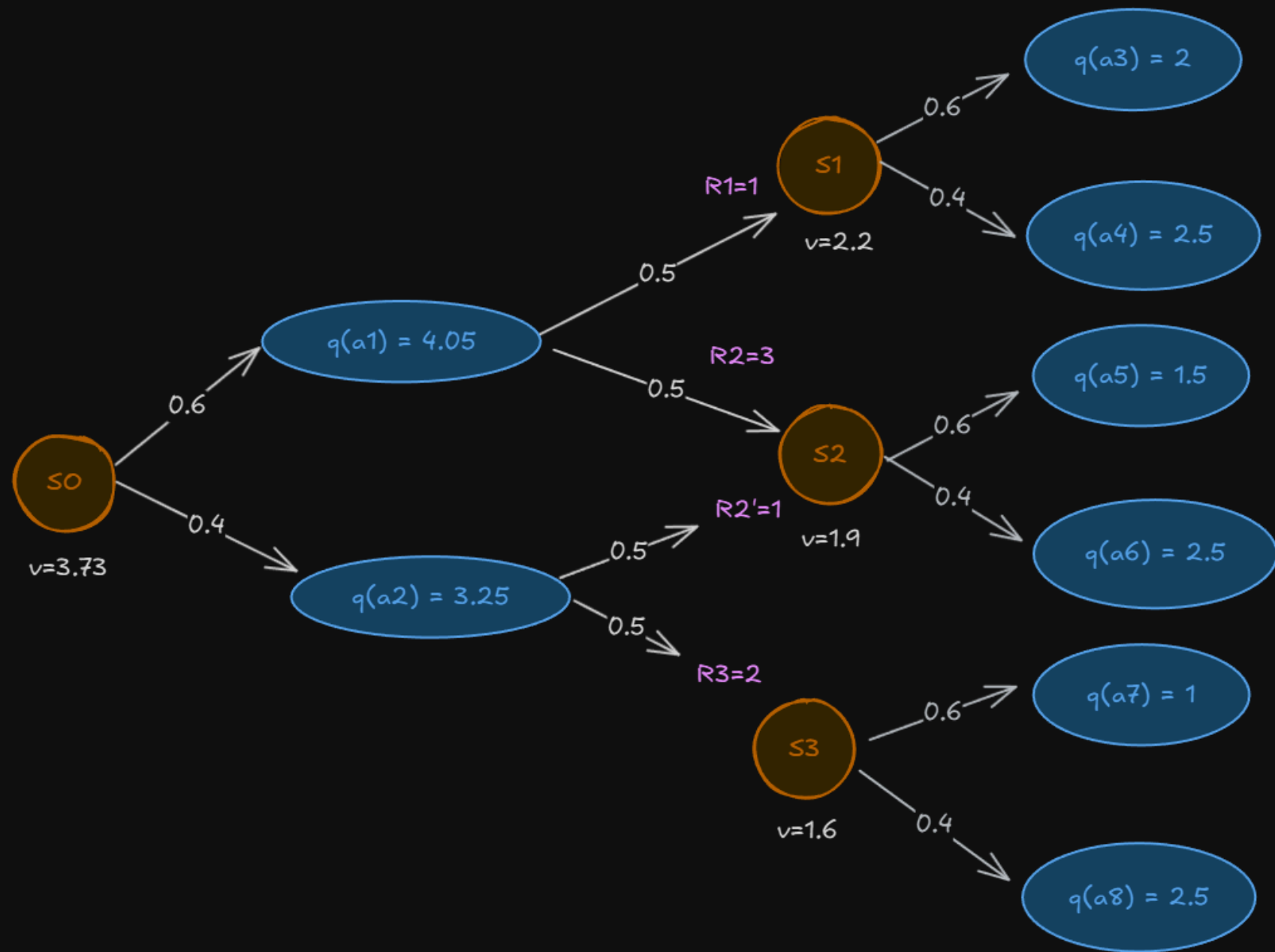


**Find**  $q(a_1), q(a_2)$

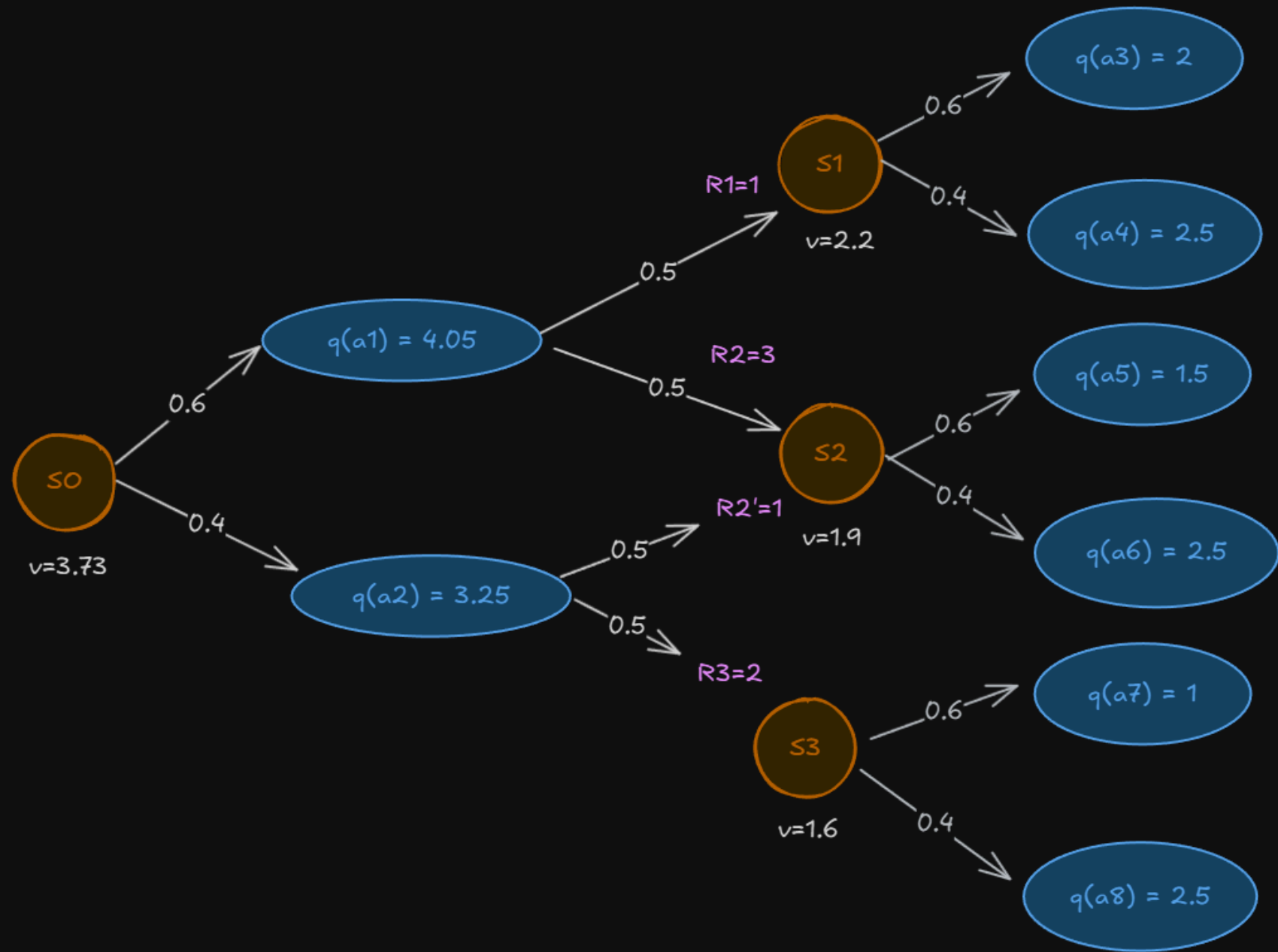
$$\begin{aligned} q(a_1) &= 0.5 \left[ 1 + \cancel{\gamma} (0.6(2) + 0.4(2.5)) \right] \\ &\quad + 0.5 \left[ 3 + \cancel{\gamma} (0.6(1.5) + 0.4(2.5)) \right] \\ &= 4.05 \end{aligned}$$

$$\begin{aligned} q(a_2) &= 0.5 \left[ 1 + \gamma (0.6(1.5) + 0.4(2.5)) \right] \\ &\quad + 0.5 \left[ 2 + \gamma (0.6(1) + 0.4(2.5)) \right] \\ &= 3.25 \end{aligned}$$

# Result

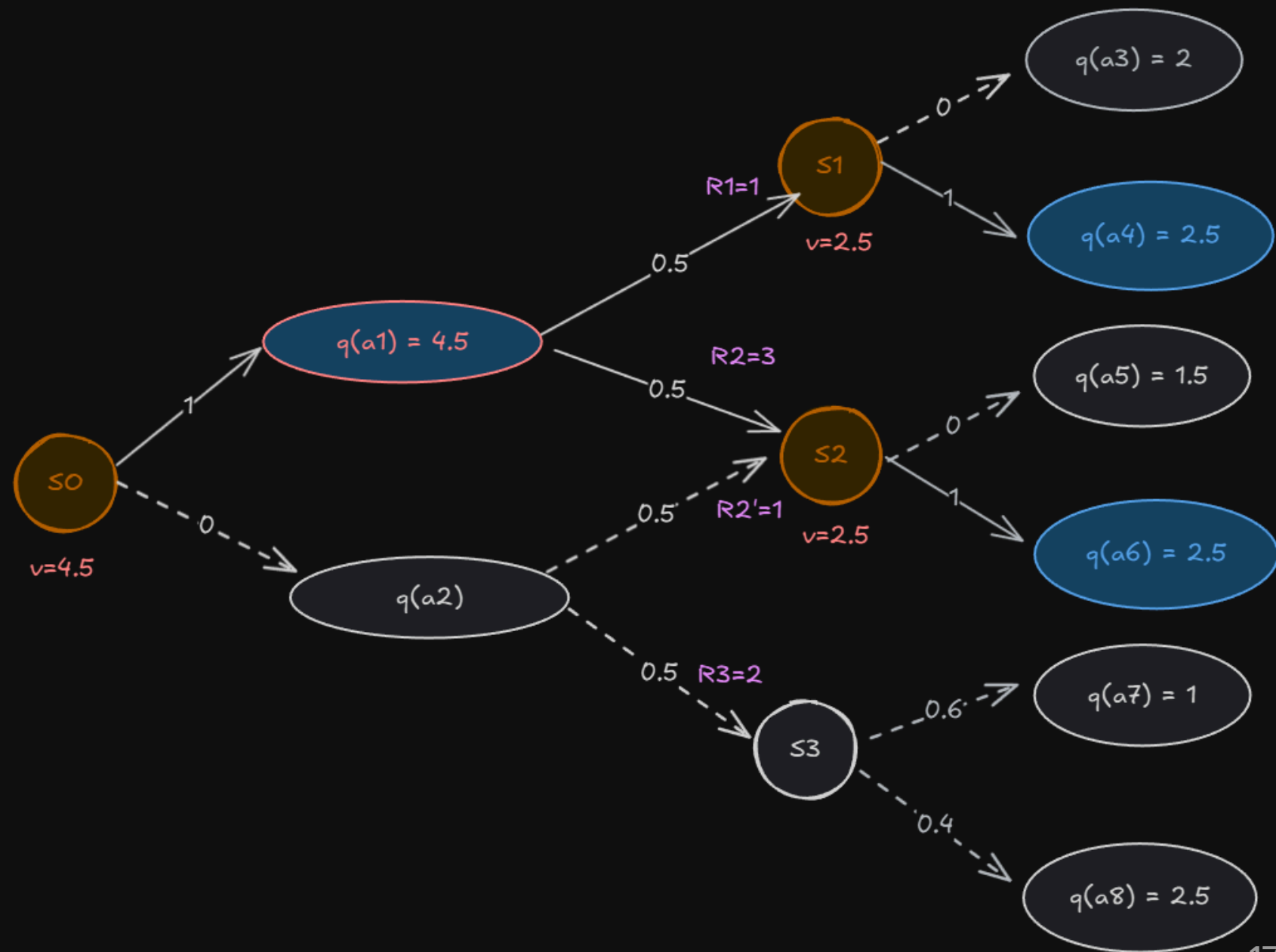


Can we do better?





# Optimal Policy



# Optimality Condition

- If agent is following the optimal policy  $\pi^*$  (something we want to find), then the value function will also be optimal.

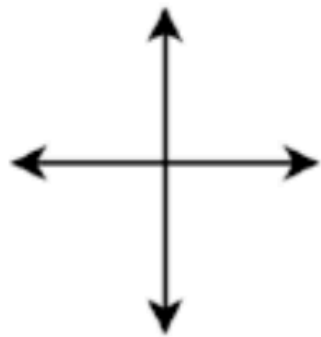
$$v^*(s) = \max_{\pi} v_{\pi}(s)$$

- It follows that

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

$$q^*(s, a) = \sum_{s', r} p(s', r | s, a) \cdot \left[ r + \gamma \max_{a'} q^*(s', a') \right]$$

# Gridworld Environment



Actions

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

16 states represented as "0" to "15"  
reward = -1 for all transitions

# Gridworld

- The top-left and bottom-right positions are terminal states, depicted as shaded cells.
- In each cell, the agent can move UP, RIGHT, DOWN, or LEFT.
- Actions deterministically move the agent in the specified direction, unless blocked by a wall.
- If the agent tries to move into a wall, it remains in its current position.
- The agent receives a reward of -1 for each step taken until it reaches a terminal state.

# Policy Evaluation / Prediction

From Bellman's equation.

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

We use iterative solution.

$$v_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$$

## ITERATIVE POLICY EVALUATION

Input  $\pi$ , the policy to be evaluated and convergence threshold  $\theta$

Initialize state values  $v(s) = 0$  or to any arbitrary values for all state  $s \in S$ .

However terminal state values should always be initialized to 0

Make a copy:  $v'(s) \leftarrow v(s)$  for all  $s$

Loop:

$$\Delta = 0$$

Loop for each  $s \in S$

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

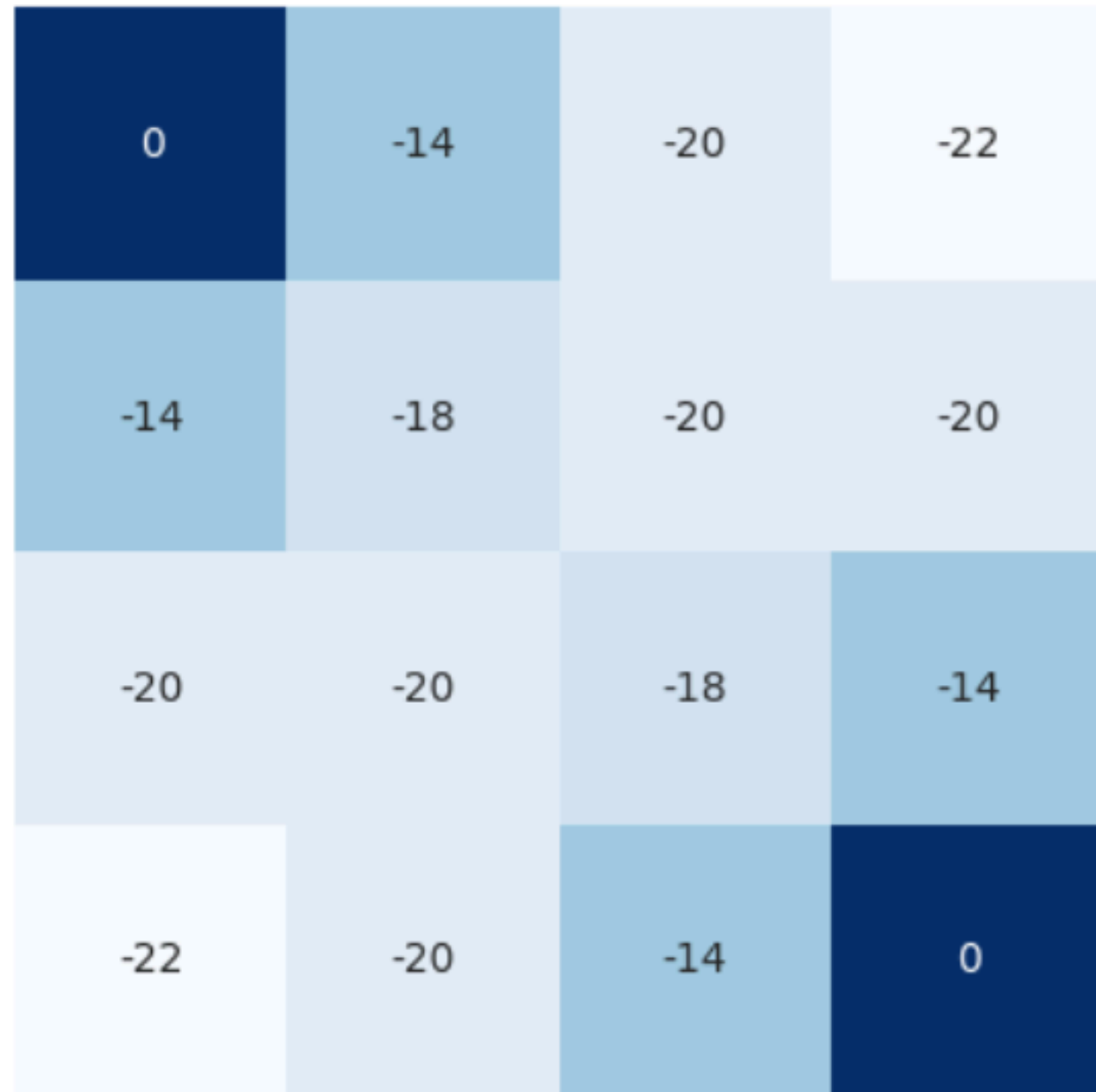
$$\Delta = \max(\Delta, |v(s) - v'(s)|)$$

$v(s) \rightarrow v'(s)$  for all  $s \in S$ , i.e., make a copy of  $v(s)$

Until  $\Delta < \theta$

# Result

- See a python code that uses a random policy.



# Validate Result

0	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	0

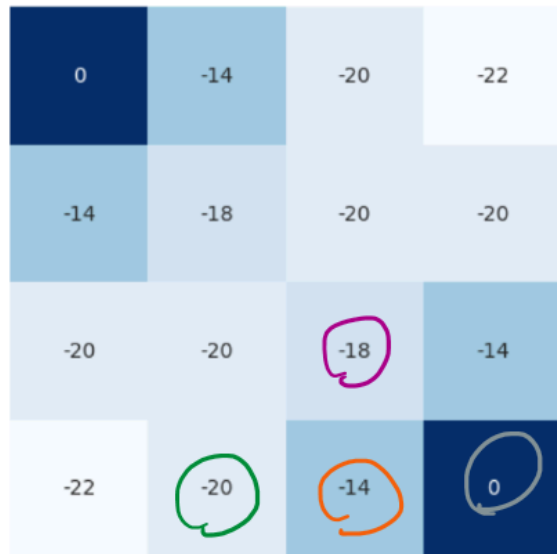
$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_{\pi}(s')]$$

$$14 = 0.5 [ 1(-1 + 1(-18)) ] + 0.5 [ 1(-1 + 1(-20)) ] + 0.5 [ 1(-1 + 1(-14)) ] + 0.5 [ 1(-1 + 1(0)) ]$$

$$14 = 0.5 (-19 -21 -15) = 14$$



# Policy Improvement



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

$$14 = 0.5 [ 1(-1 + 1(-18)) ] + 0.5 [ 1(-1 + 1(-20)) ] + 0.5 [ 1(-1 + 1(-14)) ] + 0.5 [ 1(-1 + 0) ]$$

$q(\text{up})$   
 $q(\text{left})$   
 $q(\text{down})$   
 $q(\text{right})$   
 Policy

# Policy Improvement

- You can improve the current policy by choose the action with maximum  $q$  value and define this to be the new policy.

$$\pi'(a|s) = \underset{a}{\operatorname{argmax}} q(s, a)$$

# Policy Improvement

- This is known as the *greedy* policy.
- Also Note the new state value  $v = -1$ .

$$\begin{aligned}
 14 &= 0.5 \left[ 1(-1 + 1(-18)) \right] + q(\text{up}) \\
 &\quad 0.5 \left[ 1(-1 + 1(-20)) \right] + q(\text{left}) \\
 &\quad 0.5 \left[ 1(-1 + 1(-14)) \right] + q(\text{down}) \\
 &\quad 0.5 \left[ 1(-1 + 0) \right] + q(\text{right}).
 \end{aligned}$$

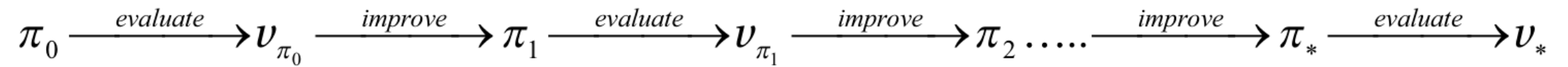
Current Policy =  $[0.5, 0.5, 0.5, 0.5]$

$$\pi'(a|s) = \underset{a}{\operatorname{argmax}} q(s, a)$$

$$-1 = \underline{1} \left[ 1(-1 + 0) \right] + q(\text{right}).$$

New Policy =  $[0, 1, 0, 0]$

# Policy Iteration



## POLICY ITERATION

Initialize state values  $v(s)$  and policy  $\pi$  arbitrarily

e.g.  $v(s) = 0, \forall s \in S$ , and  $\pi(a|s)$  as random

define convergence threshold  $\theta$

### Policy-Evaluation-Step

Loop:

$$\Delta = 0$$

Loop for each  $s \in S$

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

$$\Delta = \max(\Delta, |v(s) - v'(s)|)$$

$v(s) \leftarrow v'(s)$  for all  $s \in S$ , i.e. make a copy of  $v(s)$

Until  $\Delta < \theta$

### Policy-Improvement-Step

policy-changed  $\leftarrow$  false

Loop for each  $s \in S$

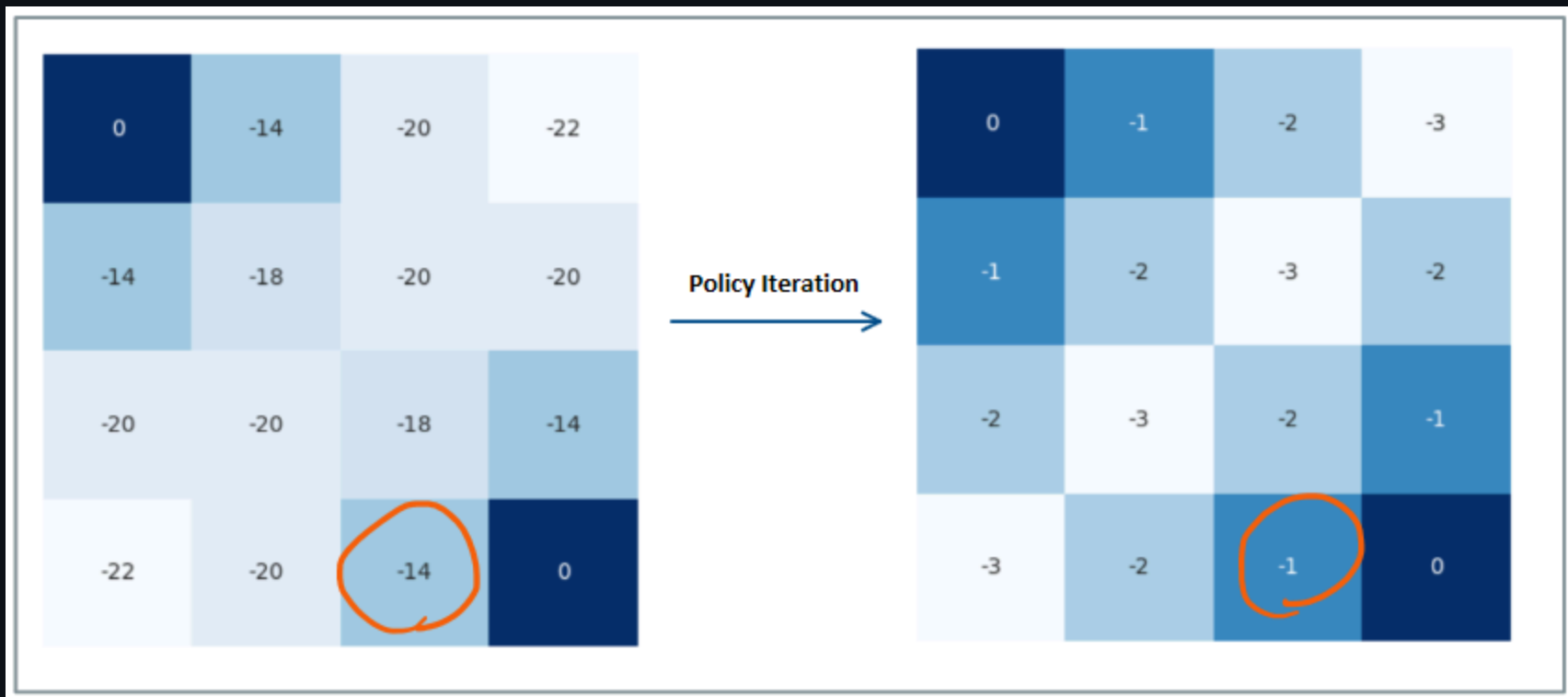
old-action  $\leftarrow \pi(s)$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma v(s')]$$

If  $\pi(s) \neq$  old-action then policy-changed = true

If policy-changed = true: then go back to Policy-Evaluation step followed by Policy Improvement

otherwise return  $v(s)$  as  $v_*$  and  $\pi(s)$  as  $\pi_*$



Optimal Policy

```
[[0.25 0.25 0.25 0.25]
```

```
[0.  0.  0.  1. ]
```

```
[0.  0.  0.  1. ]
```

```
[0.  0.  0.5 0.5 ]
```

```
[1.  0.  0.  0. ]
```

```
[0.5 0.  0.  0.5 ]
```

```
[0.  0.  0.5 0.5 ]
```

```
[0.  0.  1.  0. ]
```

```
[1.  0.  0.  0. ]
```

```
[0.5 0.5 0.  0. ]
```

```
[0.  1.  0.  0. ]
```

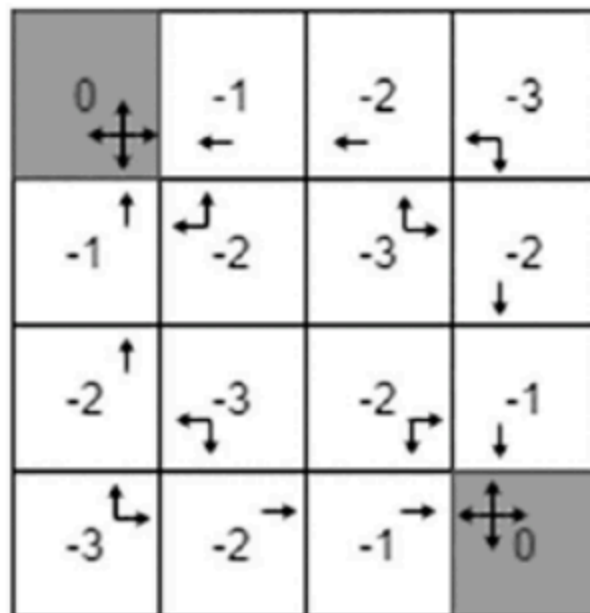
```
[0.  0.  1.  0. ]
```

```
[1.  0.  0.  0. ]
```

```
[0.  1.  0.  0. ]
```

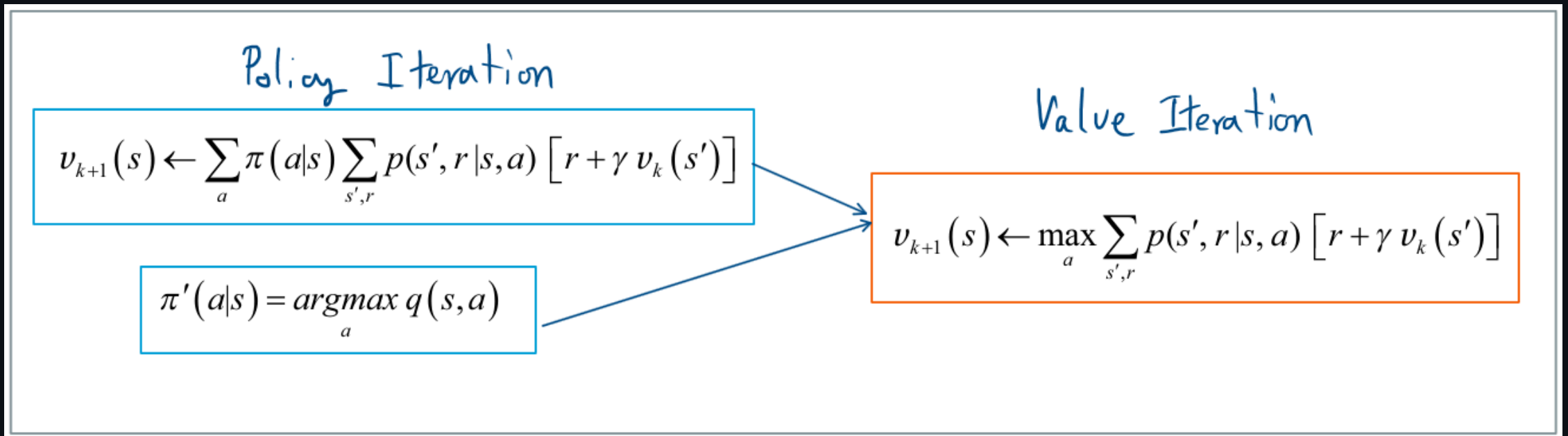
```
[0.  1.  0.  0. ]
```

```
[0.25 0.25 0.25 0.25]]
```



# Value Iteration

- Merge the policy evaluation and policy improvement together into single operation.





# Value Iteration

- The state values will keep improving.

$$v_0 \rightarrow v_1 \rightarrow \dots v_k \rightarrow v_{k+1} \rightarrow \dots \rightarrow v_*$$

- Once the values converges, we can find the optimum policy.

$$\pi_*(a|s) = \underset{a}{\operatorname{argmax}} \sum_{s',r} p(s',r|s,a) [r + \gamma v_*(s')]$$

## VALUE ITERATION

Initialize state values  $v(s)$  (Terminal states are always initialized to 0)

e.g.  $v(s) = 0, \forall s \in S$

define convergence threshold  $\theta$

Make a copy:  $v'(s) \leftarrow v(s)$  for all  $s$

Loop:

$$\Delta = 0$$

Loop for each  $s \in S$

$$v'(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v(s')]$$

$$\Delta = \max(\Delta, |v(s) - v'(s)|)$$

$v(s) \leftarrow v'(s)$  for all  $s \in S$ , i.e. make a copy of  $v(s)$

Until  $\Delta < \theta$

Output a deterministic policy, breaking ties deterministically.

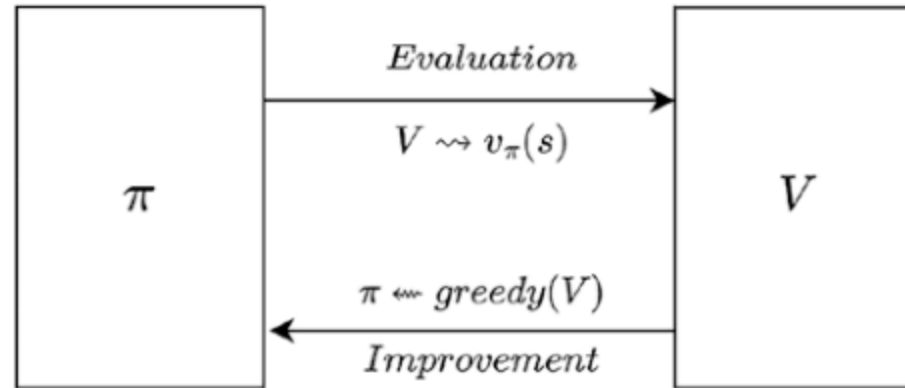
Initialize  $\pi(s)$ , an array of length  $|S|$

Loop for each  $s \in S$

$$\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v(s')]$$

# Generalized Policy Iteration

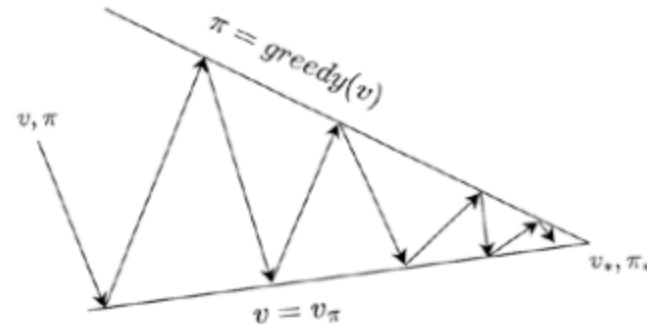
- Main idea



**Figure 3-13.** Iteration between the two steps. The first step is evaluation to get state values in sync with the policy being followed. The second step is that of policy improvement to do a greedy maximization for actions

# Generalized Policy Iteration

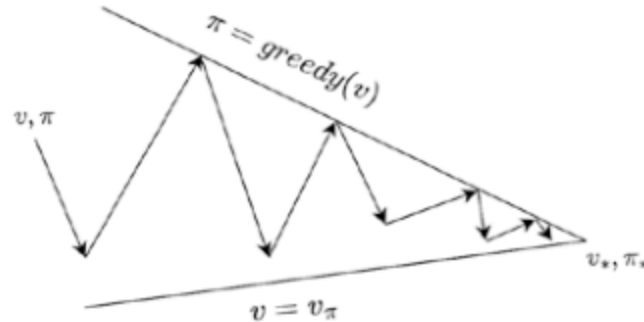
- Full convergence + change policy for all states.



**Figure 3-10.** Iteration between the two steps. The first step is evaluating to get the state values in sync with the policy being followed. The second step is improving the policy to do a greedy maximization for actions

# Generalized Policy Iteration

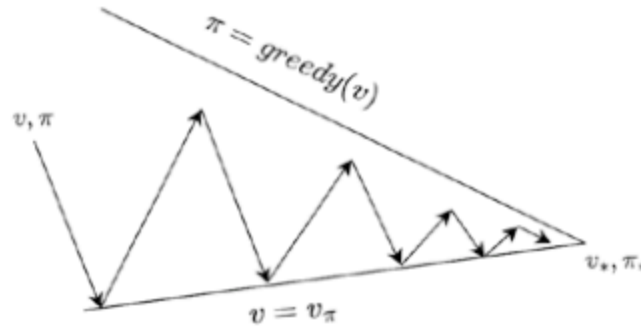
- **Partial** convergence + change policy for all states.



**Figure 3-11.** Iterating just enough to cover all the states in a single iteration leads to not touching the bottom line  $v = v_\pi$

# Generalized Policy Iteration

- Full convergence + change policy for **partial** states.



**Figure 3-12.** Iterating when the policy improvement step is not carried out over all the states, leading to not touching the top line  $\pi = \text{greedy}(v)$