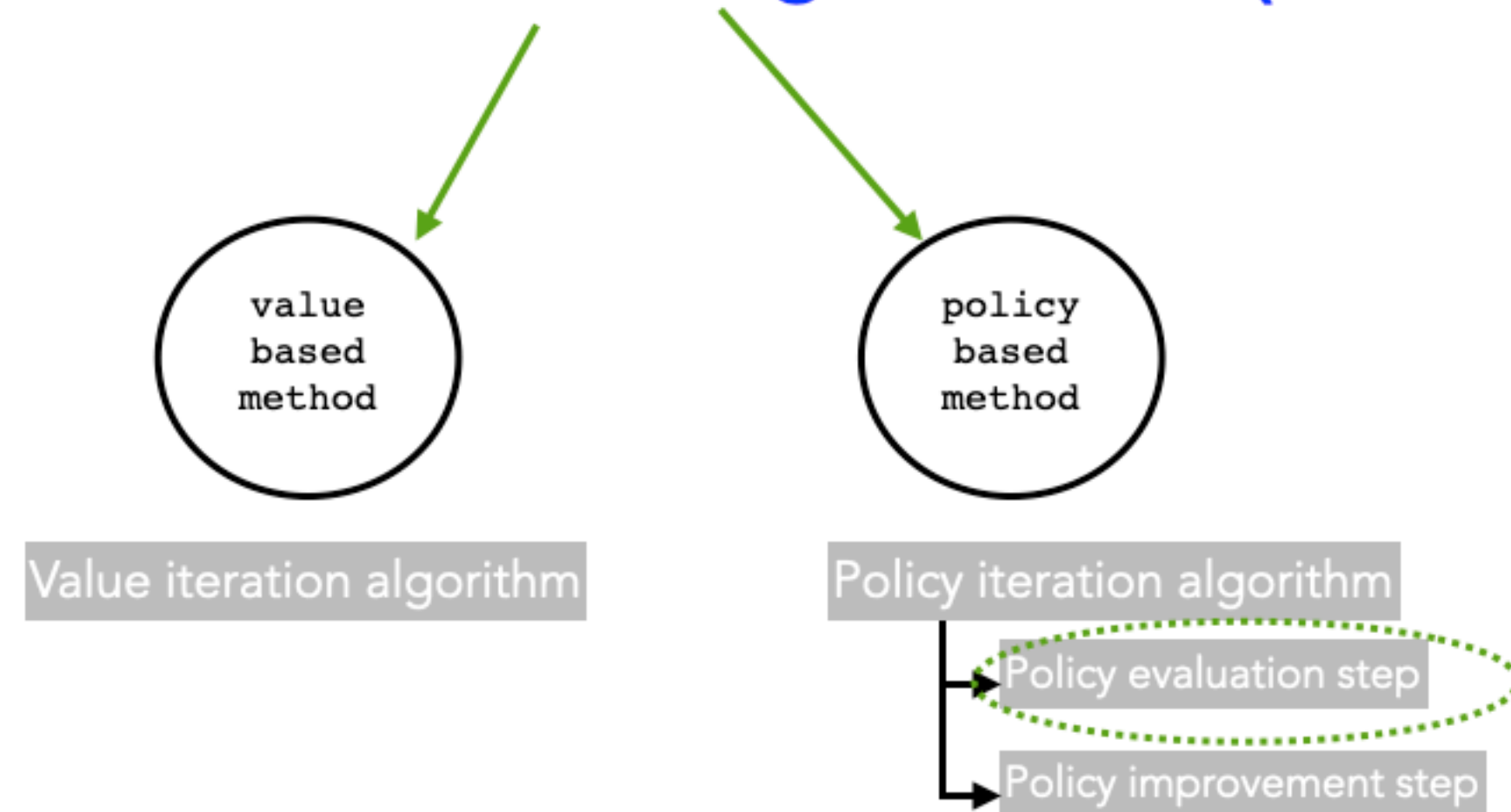


CS143: Artificial Intelligence

Reinforcement Learning Solvers (so far)



Policy Iteration Algorithm



What is the MDP for Frozen Lake?

- The transition function encodes the movements in the Lake
 - In Frozen Lake, model or dynamics is known
 - $P(s' | a, s)$ is read as agent transitions from state s to s' by taking action a
- How many states are there?
 - How can you show the probability transitions $P(s' | a, s)$ from each state to the next?

- Let's simplify and consider a simpler problem where there only a couple of states
 - Let's consider a simple dice game with 2 states only



MDP for a Simple Dice Game with 2 States

MDP for a Simple Dice Game

- For each round $t = 1, 2, 3, \dots$
 - You choose **stay** or **quit**
 - If **quit**, you get \$10 and we end the game
 - If **stay**, you get \$4 and then I roll a 6-sided dice
 - If the dice results in 1 or 2, we end the game
 - Otherwise, if the dice results in 3, 4, 5, or 6, we continue to the next round

WHAT IS THE BEST STRATEGY FOR THIS GAME?

stay

quit

MDP for a Simple Dice Game

- You choose **stay** or **quit**
 - If **quit**, you get \$10 and we end the game
 - If **stay**, you get \$4 and then I roll a 6-sided dice
 - If the dice results in 1 or 2, we end the game
 - Otherwise, if the dice results in 3, 4, 5, or 6, we continue to the next round

WHAT IS THE BEST STRATEGY FOR THIS GAME?

stay

Expected rewards:

12.0

quit

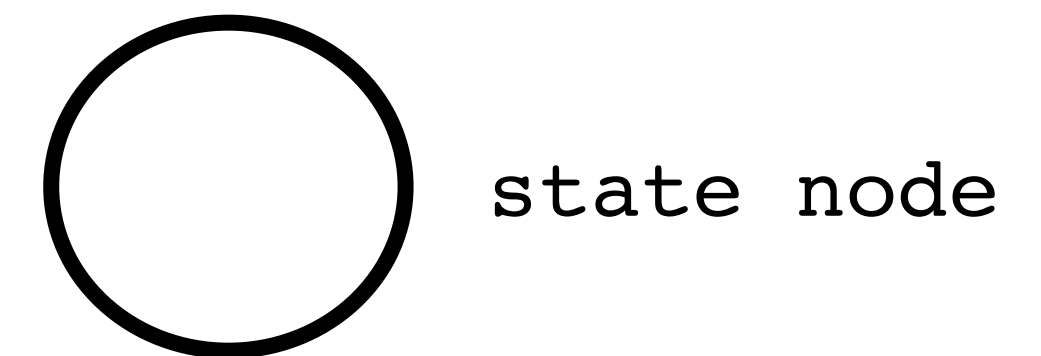
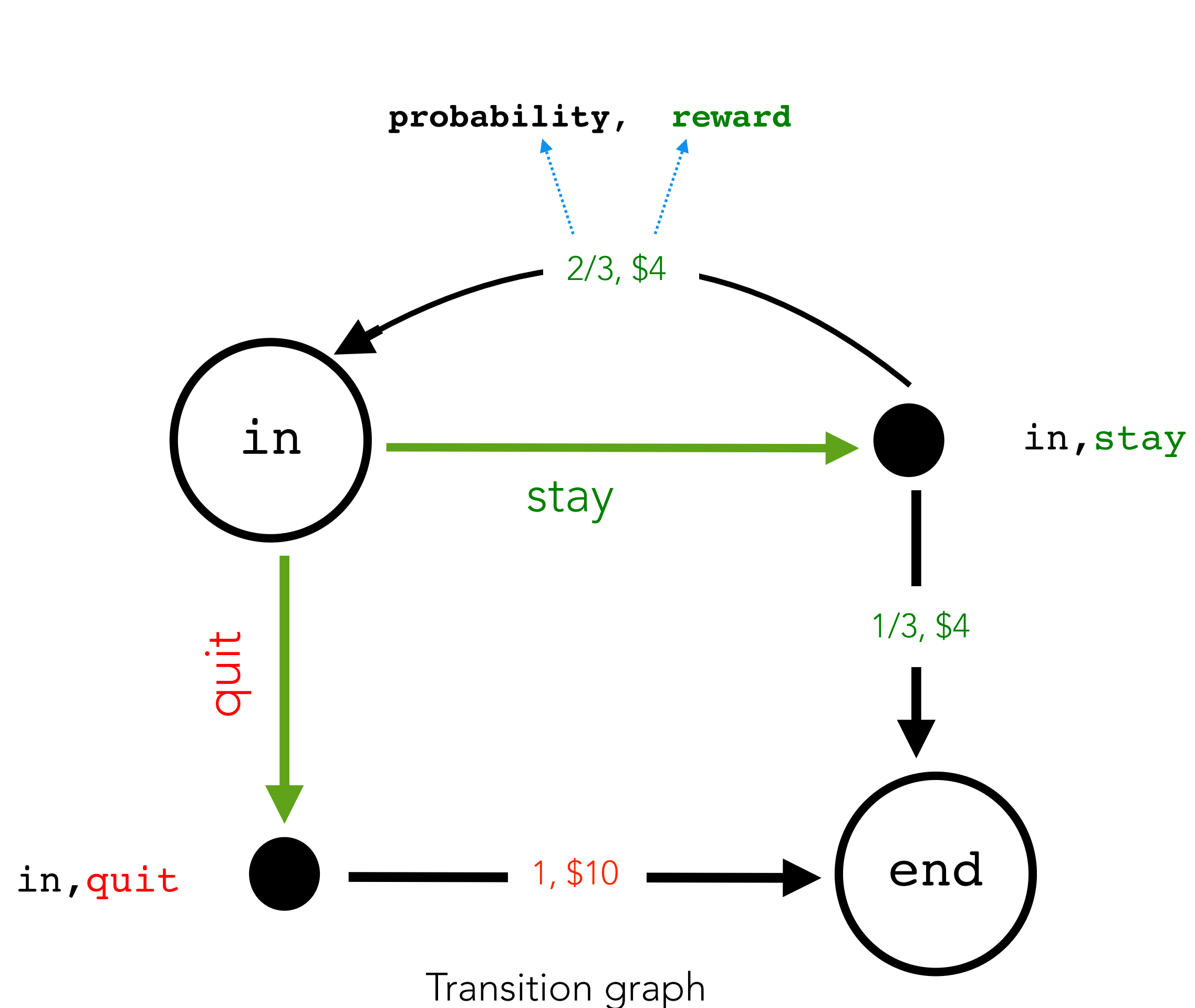
Expected rewards:

10.0

If you **quit**, then you'll get a reward of \$10. Therefore, the **stay** strategy is preferred, even though sometimes you'll get less than \$10.

MDP for a Simple Dice Game

- Let us formalize the dice game as a Markov decision process (MDP)



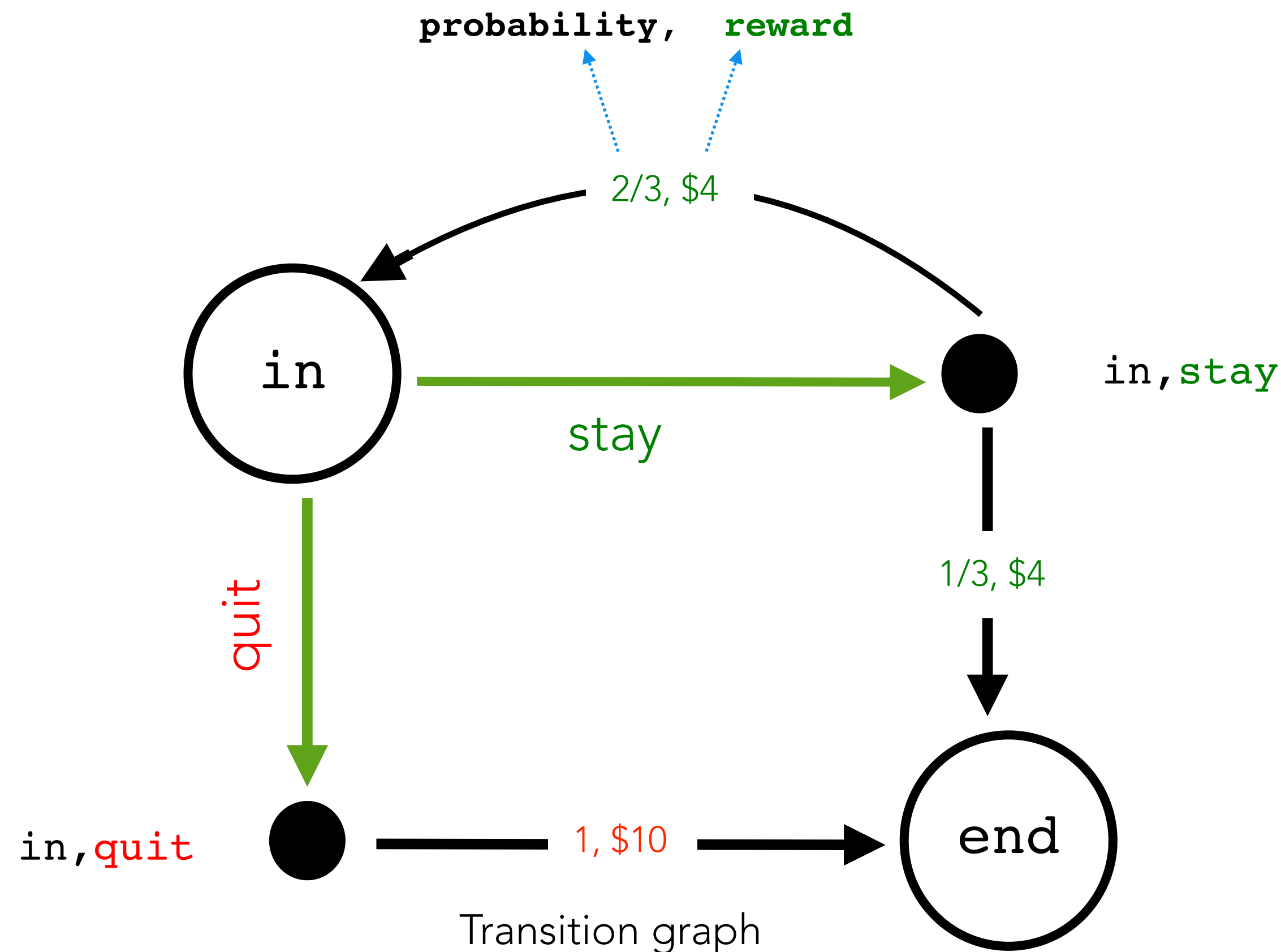
- Edges coming out of a state node are the possible actions from that state node leading to a action node



- Edges coming out of a action node are the possible **random outcomes** of that action, which end up back in state nodes
- Sum of **probabilities** coming out of a action node is 1.0

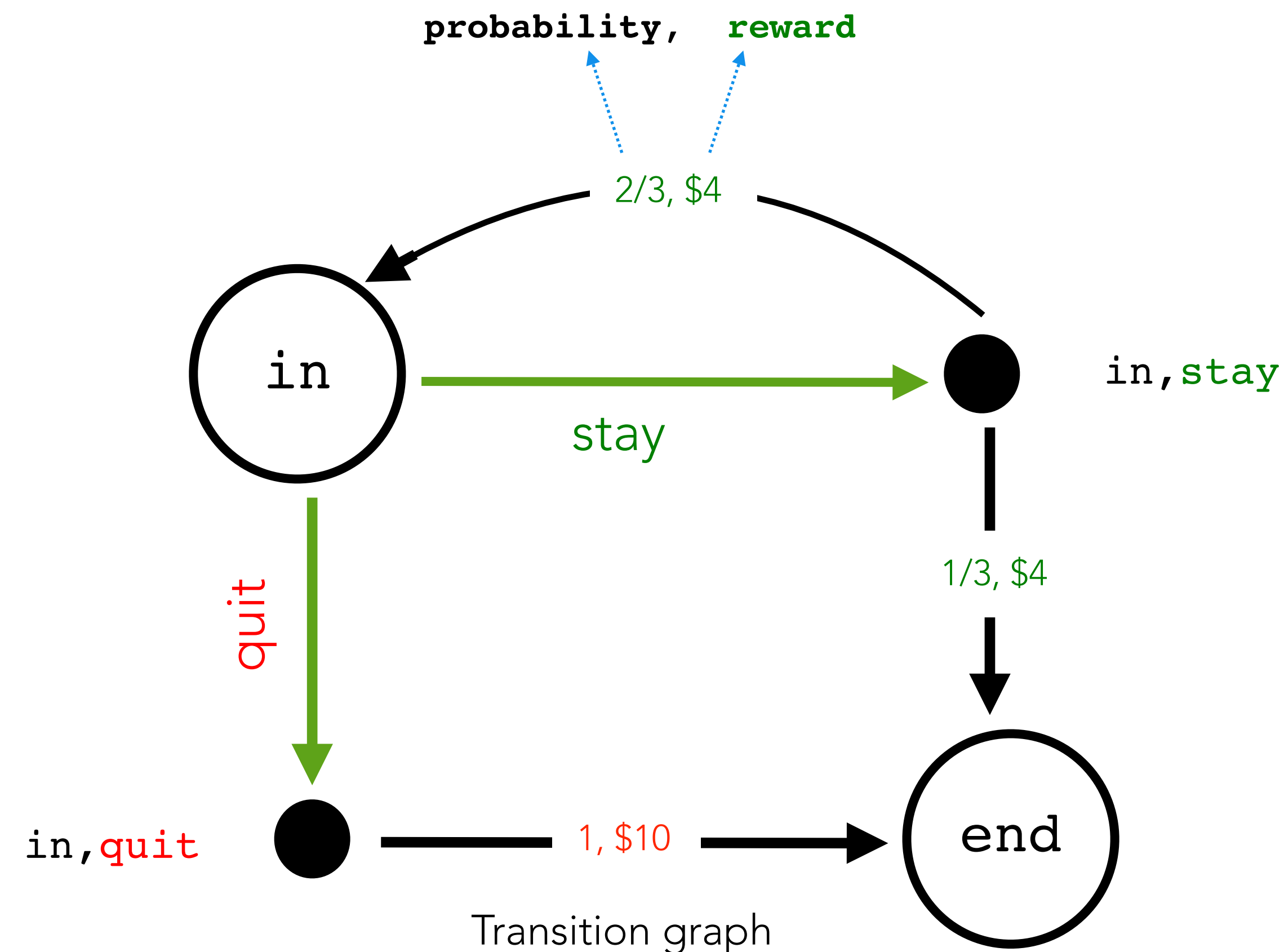
Optimal state-value or value function

- **value function** $V_{opt}(s)$: estimate of *how good* it is to be in a given state s and acting optimally
- $V_{opt}(s)$ is the expected reward received when starting in state s and following optimal policy π_{opt} thereafter



Optimal action-value or q-value function

- **q-value function** $Q_{opt}(s, a)$: estimate of the expected return obtained by taking action a in state s , and thereafter following the optimal policy π_{opt} optimally thereafter



Relationship: value and q-value functions

- Looks like from the previous to visualizations, the relationship between the **value function** and the **Q-value function** can be expressed recursively. For example, from the previous tree graph, it can be seen that the value function can be computed by taking the maximum over multiple Q-values corresponding to different actions. Mathematically, this can be written as follows:

$$V_{opt}(s) = \max \left\{ \begin{array}{l} Q_{opt}(s, quit) \\ Q_{opt}(s, stay) \end{array} \right\}$$

THERE IS ONLY ONE NON-TERMINAL STATE: *in* TERMINAL STATE: *end*

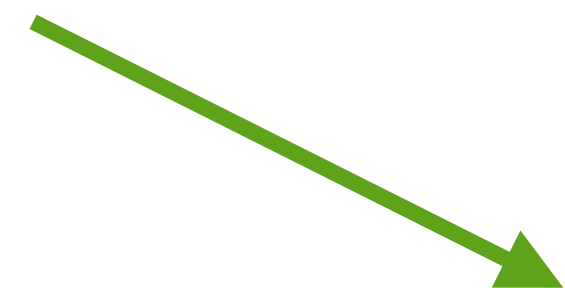
$V_{opt}(in) = ?$

$V_{opt}(end) = 0$

$$V_{opt}(in) = \max \left\{ \begin{array}{l} P(in | s, quit)[R(s, quit, in) + \gamma V_{opt}(in)] + P(end | s, quit)[R(s, quit, end) + \gamma V_{opt}(end)] \\ P(in | s, stay)[R(s, stay, in) + \gamma V_{opt}(in)] + P(end | s, stay)[R(s, stay, end) + \gamma V_{opt}(end)] \end{array} \right\}$$

$$= \max \left\{ \begin{array}{l} P(end | s, quit)[R(s, quit, end) + \gamma V_{opt}(end)] \\ P(in | s, stay)[R(s, stay, in) + \gamma V_{opt}(in)] + P(end | s, stay)[R(s, stay, end) + \gamma V_{opt}(end)] \end{array} \right\}$$

$Q_{opt}(in, quit)$



$$\frac{2}{3}(4 + \gamma V_{opt}(in)) + \frac{1}{3}(4 + \gamma V_{opt}(end))$$

\downarrow $p(s' | s, a)$

\downarrow $r(s, a, s')$

\downarrow Discount

\downarrow $p(s' | s, a)$

\downarrow $r(s, a, s')$

\downarrow Discount

$Q_{opt}(in, stay)$

Relationship: value and q-value functions

- In general, we can write the value function can be computed by taking the maximum over multiple Q-values corresponding to different actions. Mathematically, this can be written as follows:

$$\begin{aligned} V_{opt}(s) &= \max_a Q_{opt}(s, a) \\ &= \max_a \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma V_{opt}(s')] \end{aligned}$$

Question: Then how do you solve an MDP?

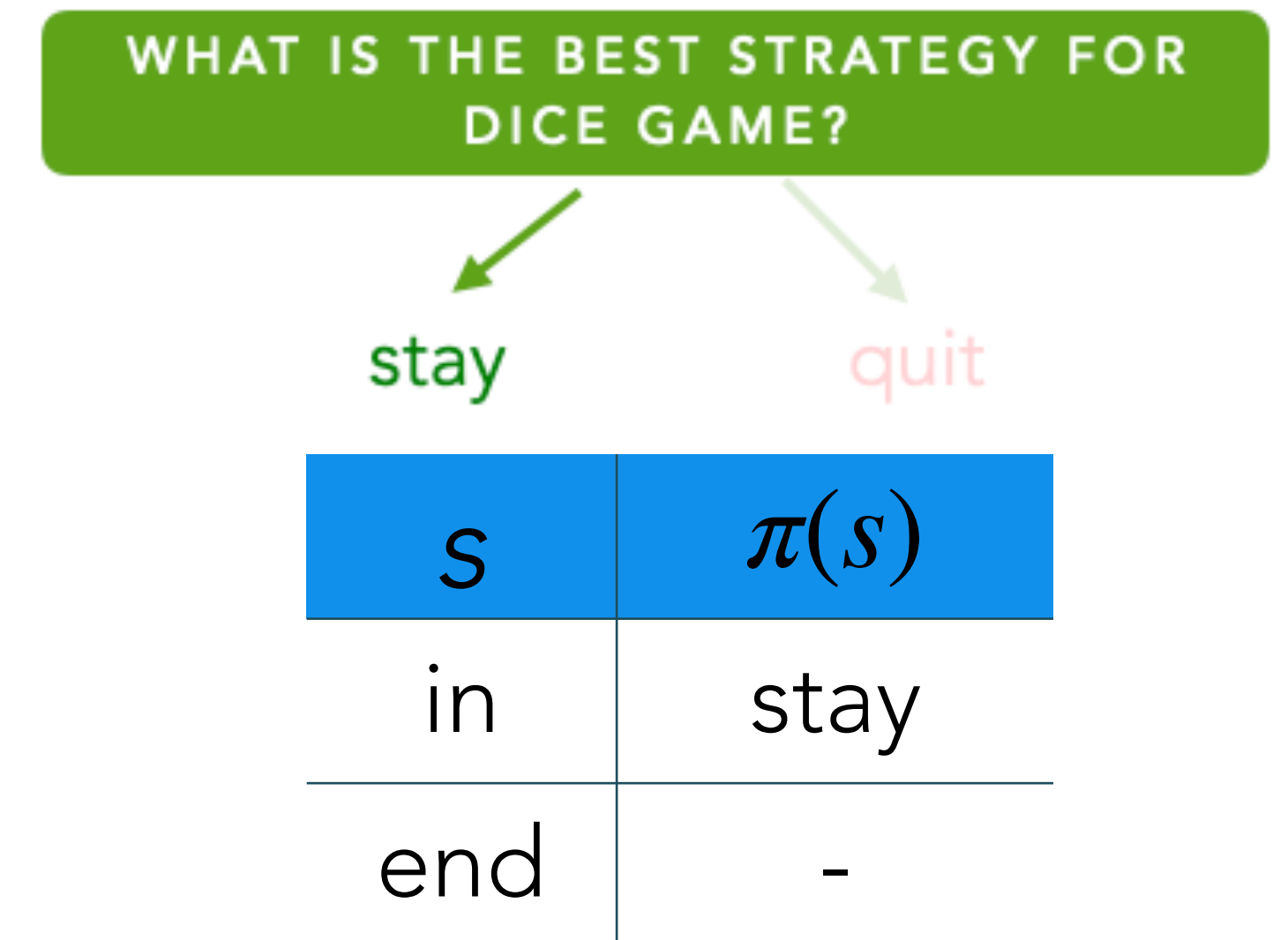
Solution to an MDP is called a 'policy'

Answer: Value Iteration Algorithm

An iterative algorithm to find the optimal policy (best solution) for an MDP

Policy: Solution to an MDP

- What is a solution to an MDP?
 - Solution to an MDP is a **policy** π which is a mapping **from each state** s to an action a
- In search problem, what was the solution?
 - Solution was the path, which was a sequence of actions from the **start state** to **goal state** only
 - Eg, In an arbitrary 8-puzzle problem the solution was *UP, LEFT, LEFT, RIGHT, ...*



Mathematically: How to calculate $V(s)$ values?

- How do we calculate $V_{opt}(in)$ (optimal V-values for each state) from the recursive Bellman equation:

$$V_{opt}(in) = \max_{a \in \text{actions}(s)} \sum_{s'} p(s' | in, a) [R(in, a, s') + \gamma V_{opt}(s')]$$

A RECURRENCE RELATION

- Turn recursive Bellman equations into update equations:

$$V_{opt}^0(in) = 0$$

$$V_{opt}^t(in) = \max_{a \in \text{actions}(s)} \sum_{s'} p(s' | in, a) [R(in, a, s') + \gamma V_{opt}^{t-1}(s')]$$

Mathematically: How to calculate $V(s)$ values?

- More precisely the update equations are as follows:

NON-TERMINAL STATE

$$V_{opt}^t(in) = \max \left\{ \begin{array}{l} Q_{opt}^t(in, quit) \\ Q_{opt}^t(in, stay) \end{array} \right\}$$
$$= \max \left\{ \begin{array}{l} \frac{1}{1}(10 + \gamma V_{opt}^{t-1}(end)) \\ \frac{2}{3}(4 + \gamma V_{opt}^{t-1}(in)) + \frac{1}{3}(4 + \gamma V_{opt}^{t-1}(end)) \end{array} \right\}$$

TERMINAL STATE

$$V_{opt}(end) = 0$$

THE RECURRENCE RELATION IS NOW EXPRESSED AS ITERATIVE UPDATES THAT CAN BE COMPUTED INCREMENTALLY

Build an Expectiminimax-like search tree

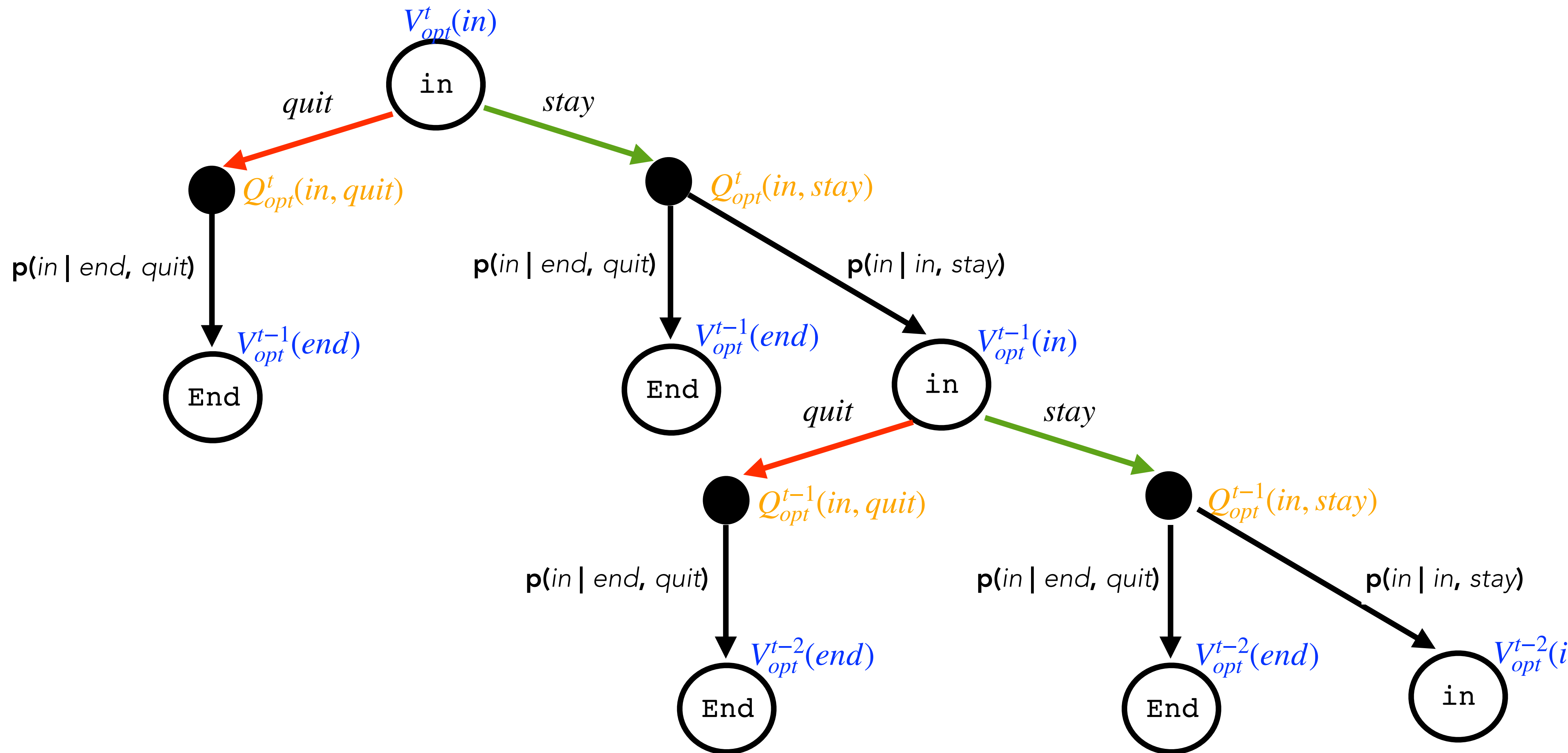
max

expectation

max

expectation

max



Algorithmically: Value iteration algorithm

INITIALIZE VALUES WITH SOME ROUGH ESTIMATES

REFINE IT OVER AND OVER AGAIN UNTIL THEIR VALUES DO NOT CHANGE

Algorithm 2: Value Iteration Algorithm

Initialize state-value $v_{opt}^0(s)$ for all the states s

//Repeat the iterative procedure for some iterations, let's say $T = 10000$

for $t = 1$ **to** T **do**

for *each state* s **do**

$v_{opt}^t(s) \leftarrow \max_{a \in \text{actions}(s)} \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v_{opt}^{t-1}(s')]$

This is because the "value iteration algorithm" is now finding the best policy out of all possible candidate policies

Value iteration algorithm on our MDP

TERMINAL STATE

$$V_{opt}(end) = 0$$

NON-TERMINAL STATE

$$V_{opt}^t(in) = \max \left\{ \begin{array}{l} Q_{opt}^t(in, quit) \\ Q_{opt}^t(in, stay) \end{array} \right\}$$

$$= \max \left\{ \begin{array}{l} \frac{1}{1}(10 + \gamma V_{opt}^{t-1}(end)) \\ \frac{2}{3}(4 + \gamma V_{opt}^{t-1}(in)) + \frac{1}{3}(4 + \gamma V_{opt}^{t-1}(end)) \end{array} \right\}$$

Algorithm 2: Value Iteration Algorithm

Initialize state-value $v_{opt}^0(s)$ for all the states s

//Repeat the iterative procedure for some iterations, let's say $T = 10000$

for $t = 1$ to T do

 for each state s do

$$v_{opt}^t(s) \leftarrow \max_{a \in actions(s)} \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v_{opt}^{t-1}(s')]$$

$V_{opt}(end)$	0	0	0	0	0	0	0	0	0	0
$\pi_{opt}(end)$	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
$V_{opt}(in)$	0									
$\pi_{opt}(in)$										
	t_0	t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	t_9

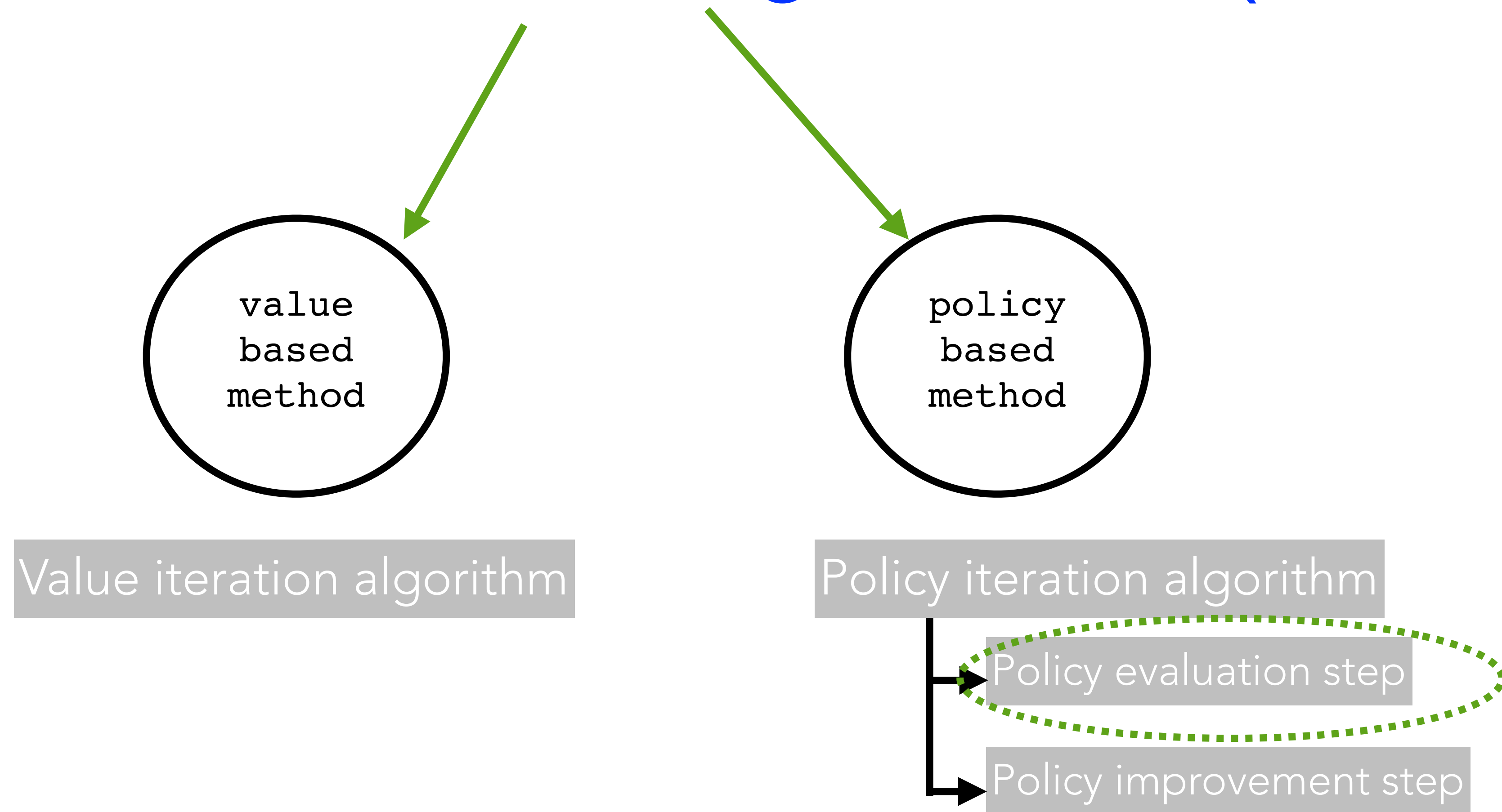
Question: Is there any other way to solve an MDP besides the Value Iteration algorithm?

Solution to an MDP is called a 'policy'

Answer (yes): Policy Iteration Algorithm

Another iterative algorithm to find the optimal policy (best solution) for an MDP

Reinforcement Learning Solvers (so far)



Let's discuss Policy Evaluation Algorithm first as it will be used by Policy Iteration Algorithm

- For an arbitrary MDP and its given policy, we can **apply** policy evaluation algorithm repeatedly until the the state-value and action-values converge for each state
- We can apply this on a very large MDP, not just the ones with 2 states (eg, our simple dice game)
- Let's say for our dice game, we are given a policy that returns action '**stay**' for each state. Since we have only one non-terminal state '**in**', it can be written as follows:

$$\pi(in) = stay$$

Policy evaluation algorithm

INITIALIZE STATE-VALUES WITH SOME ROUGH ESTIMATES

REFINE IT OVER AND OVER AGAIN UNTIL THEIR VALUES DO NOT CHANGE

Algorithm 1: Policy Evaluation Algorithm

Initialize state-value $v_{\pi}^0(s)$ for all the states s

//Repeat the iterative procedure for some iterations, let's say $T = 10000$

for $t = 1$ **to** T **do**

for *each state* s **do**

$$v_{\pi}^t(s) \leftarrow \sum_{s'} p(s'|s, \pi(s)) [r(s, \pi(s), s') + \gamma v_{\pi}^{t-1}(s')]$$

It is a specific action returned by the given policy $\pi(s)$

Policy evaluation algorithm

- We can apply **policy evaluation algorithm** repeatedly until the state-values **converge** for all the states
- How to find **convergence**? Find the difference of a state-value between two adjacent time steps until it goes down below a small threshold

$$\begin{aligned} |v_{\pi}^t(s) - v_{\pi}^{t-1}(s)| &\leq \epsilon \\ &\leq 0.001 \end{aligned}$$

Mathematically: How to calculate $V_{\pi}(s)$ values?

- How do we calculate $V_{\pi}(in)$ for a **fixed policy** $\pi(in)$ from the recursive Bellman equation:

$$V_{\pi}(in) = \sum_{s'} p(s' | in, \pi(s)) [R(in, \pi(s), s') + \gamma V_{\pi}(s')]$$

A RECURRENCE RELATION

- Turn recursive Bellman equations into update equations:

$$V_{\pi}^0(in) = 0$$

$$V_{\pi}^t(in) = \sum_{s'} p(s' | in, \pi(in)) [R(in, \pi(in), s') + \gamma V_{\pi}^{t-1}(s')]$$

Mathematically: How to calculate $V_{\pi}(s)$ values?

- More precisely the update equations are as follows:

NON-TERMINAL STATE

$$V_{\pi}^t(in) = Q_{\pi}^t(in, \pi(in))$$

$$= \begin{cases} \frac{2}{3}(4 + \gamma V_{\pi}^{t-1}(in)) + \frac{1}{3}(4 + \gamma V_{\pi}^{t-1}(end)) & \text{if } \pi(in) = \textit{stay} \\ \frac{1}{1}(10 + \gamma V_{\pi}^{t-1}(end)) & \text{if } \pi(in) = \textit{quit} \end{cases}$$

TERMINAL STATE

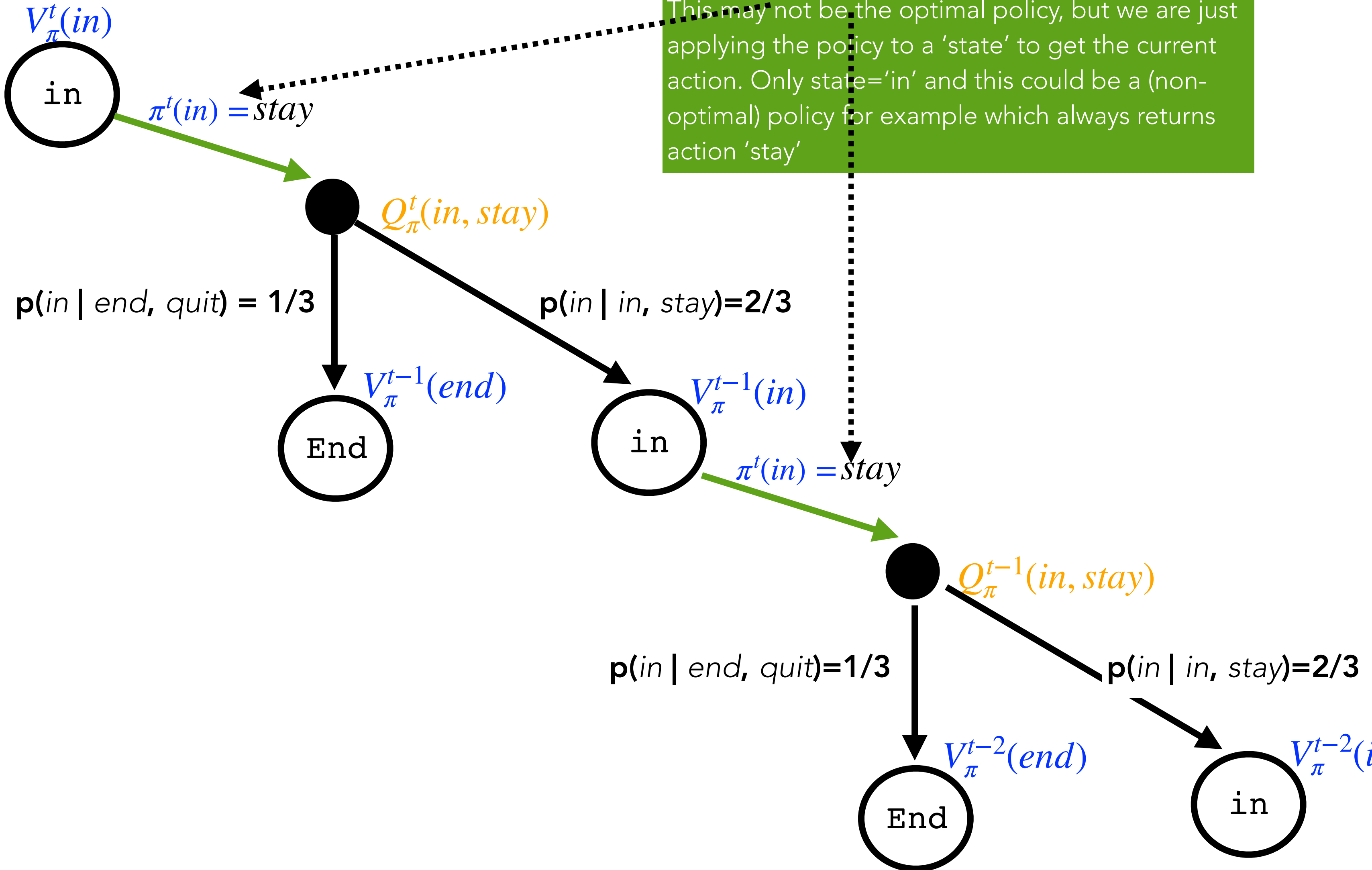
$$V_{\pi}(end) = 0$$

THE RECURRENCE RELATION IS NOW EXPRESSED AS ITERATIVE UPDATES THAT CAN BE COMPUTED INCREMENTALLY

Build an unrolled simpler tree without max computation

PLEASE NOTICE:
THE POLICY IS GIVEN

This may not be the optimal policy, but we are just applying the policy to a 'state' to get the current action. Only state='in' and this could be a (non-optimal) policy for example which always returns action 'stay'



Notice, we have excluded **max** computation

expectation

Notice, we have excluded **max** computation

expectation

Notice, we have excluded **max** computation

In contrast: value iteration with max computation

PLEASE NOTICE:
NO POLICY IS GIVEN FOR
VALUE ITERATION
ALGORITHM

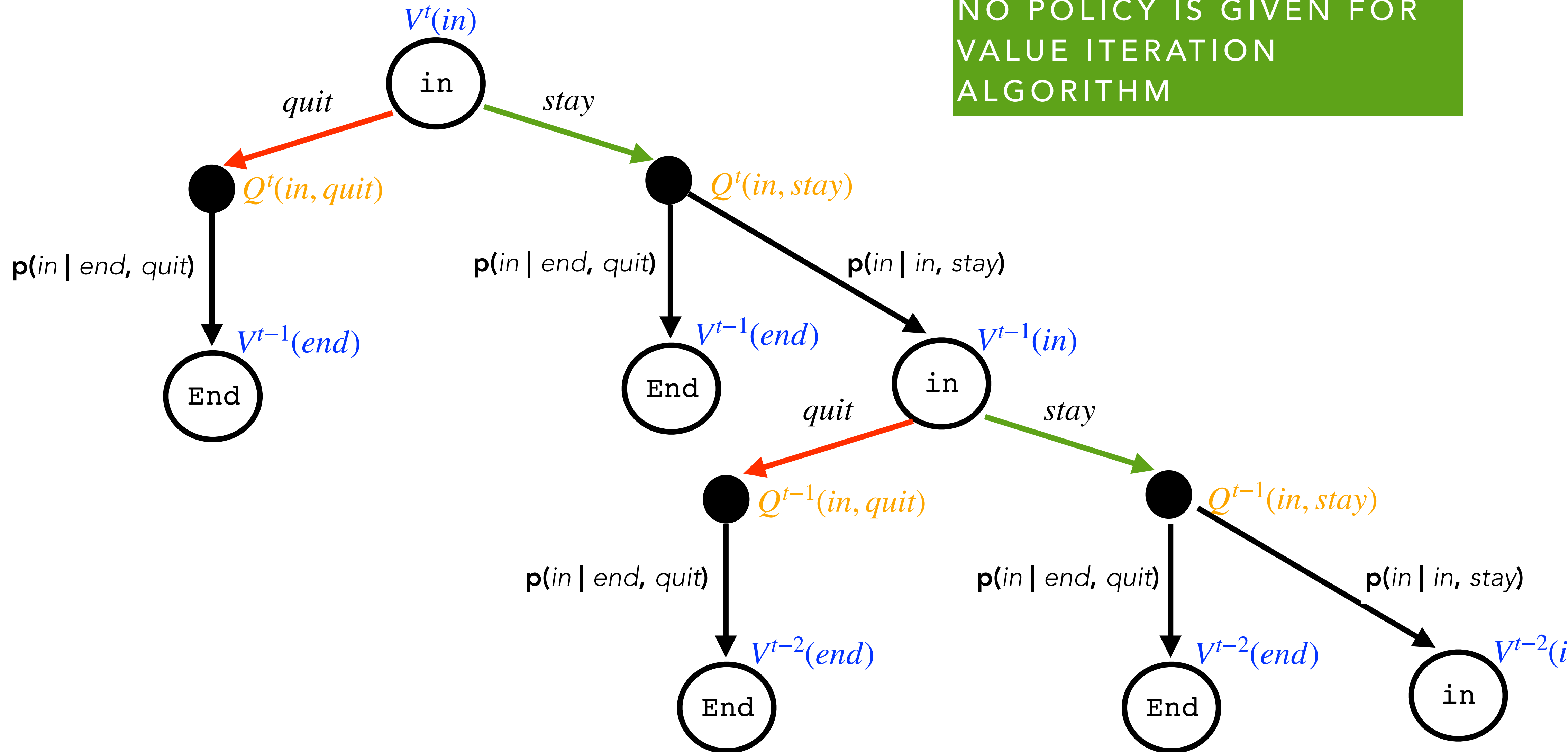
max

expectation

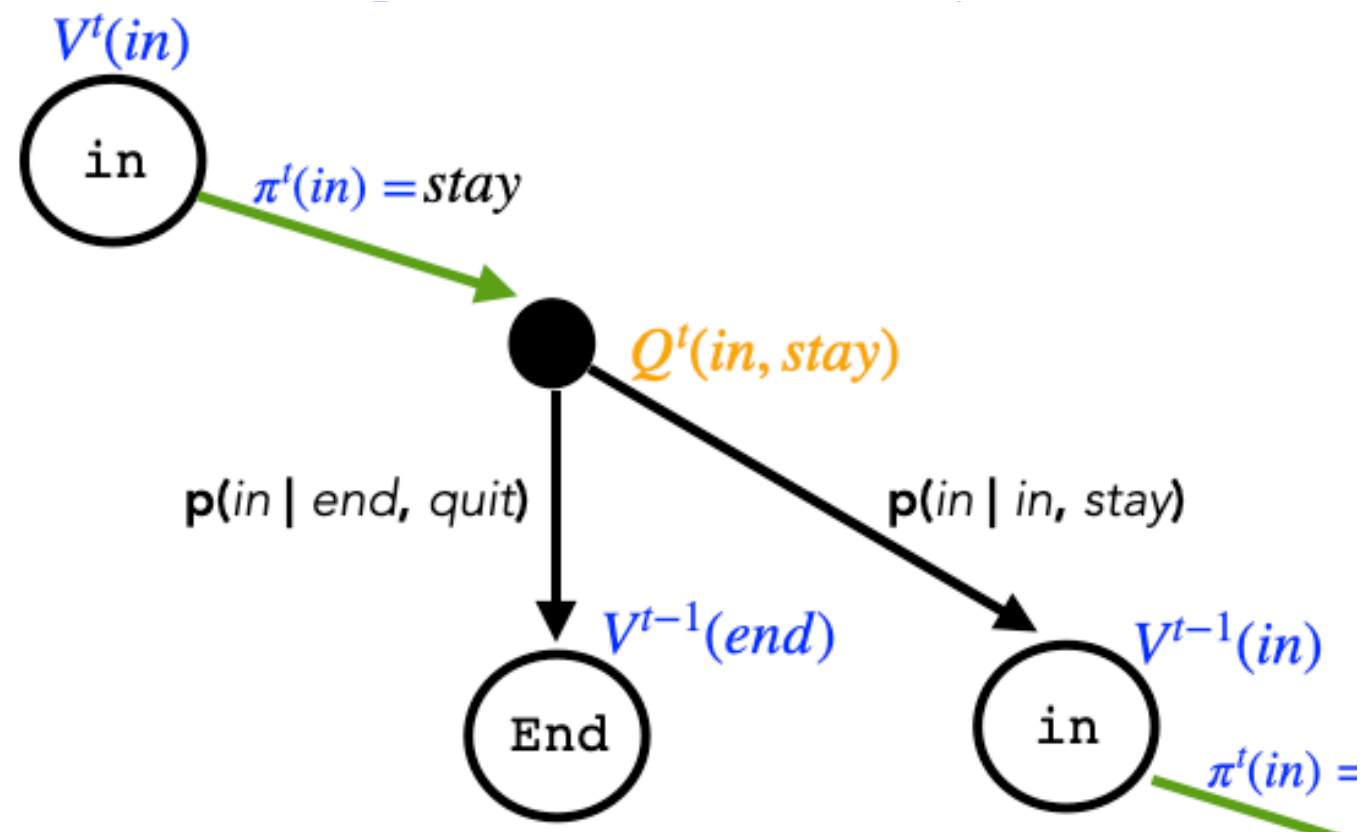
max

expectation

max



Policy evaluation on dice game



TERMINAL STATE

$$v_{\pi}(end) = 0$$

NON-TERMINAL STATE

Algorithm 1: Policy Evaluation Algorithm

Initialize state-value $v_{\pi}^0(s)$ for all the states s

//Repeat the iterative procedure for some iterations, let's s

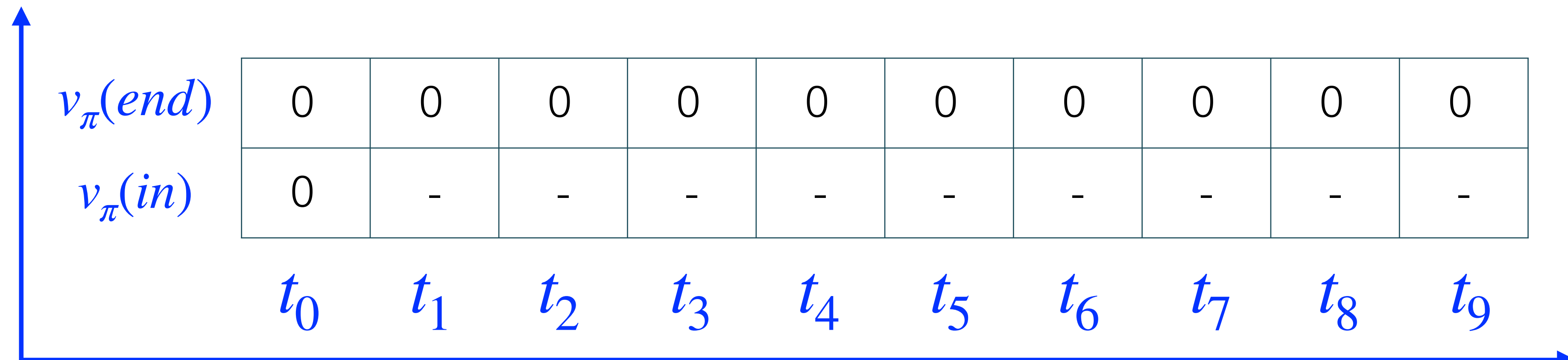
for $t = 1$ **to** T **do**

for each state s **do**

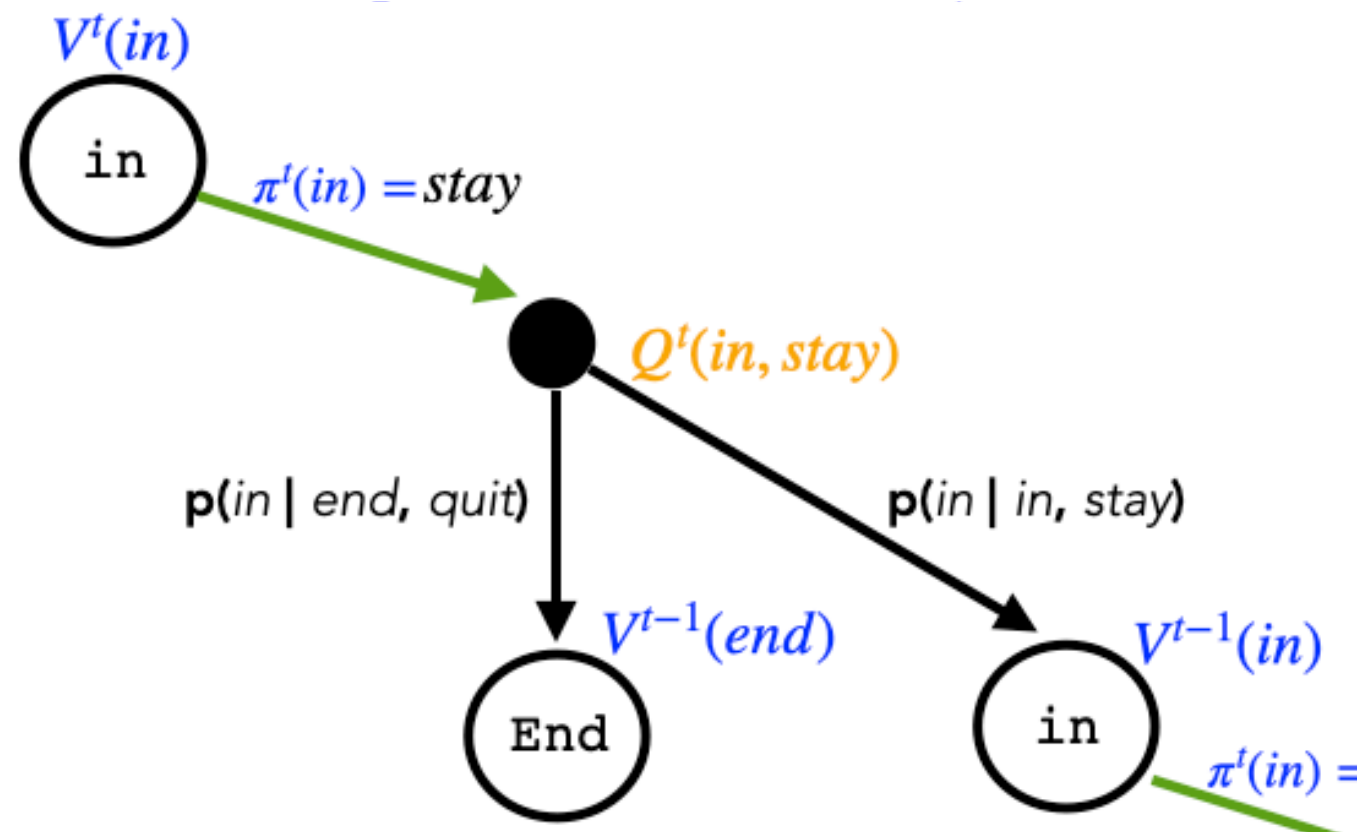
$$v_{\pi}^t(s) \leftarrow \sum_{s'} p(s'|s, \pi(s)) [r(s, \pi(s), s') + \gamma v_{\pi}^{t-1}(s')]$$

$$v_{\pi}^t(in) = q_{\pi}^t(in, a)$$

$$= \frac{2}{3}(4 + \gamma v_{\pi}^{t-1}(in)) + \frac{1}{3}(4 + \gamma v_{\pi}^{t-1}(end))$$



Solution



TERMINAL STATE

$$v_{\pi}(end) = 0$$

NON-TERMINAL STATE

Algorithm 1: Policy Evaluation Algorithm

Initialize state-value $v_{\pi}^0(s)$ for all the states s

//Repeat the iterative procedure for some iterations, let's s

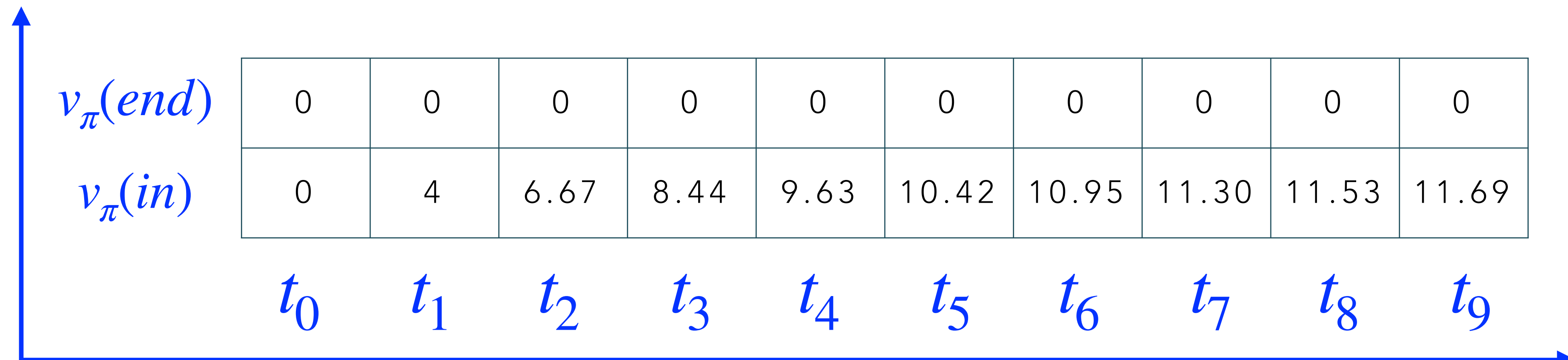
for $t = 1$ **to** T **do**

for each state s **do**

$$v_{\pi}^t(s) \leftarrow \sum_{s'} p(s'|s, \pi(s)) [r(s, \pi(s), s') + \gamma v_{\pi}^{t-1}(s')]$$

$$v_{\pi}^t(in) = q_{\pi}^t(in, a)$$

$$= \frac{2}{3}(4 + \gamma v_{\pi}^{t-1}(in)) + \frac{1}{3}(4 + \gamma v_{\pi}^{t-1}(end))$$



Contrast: value iteration vs. policy evaluation

<i>Value Iteration Algorithm</i>	<i>Policy Evaluation Algorithm</i>
No policy is given	a non-optimal policy is given
Has max computations for $V(s)$ nodes	No max computations for $V(s)$ nodes
Max operation on several action branches for $V(s)$ nodes	Just one action branch for $V(s)$ nodes
Returns: $(MDP) \rightarrow V_{opt}$ and π_{opt}	Returns: $(MDP, \pi) \rightarrow V_{\pi}$
<p>Explore the MDP and do two things: i) calculate state-value of each state ii) find the best policy</p>	<p>A (non-optimal) policy is given to you, you don't need to find it</p>

Contrast: value iteration vs. policy evaluation

Value Iteration Algorithm

Policy Evaluation Algorithm

Got my solution:
Optimal Policy



Didn't get my solution:
Optimal Policy?



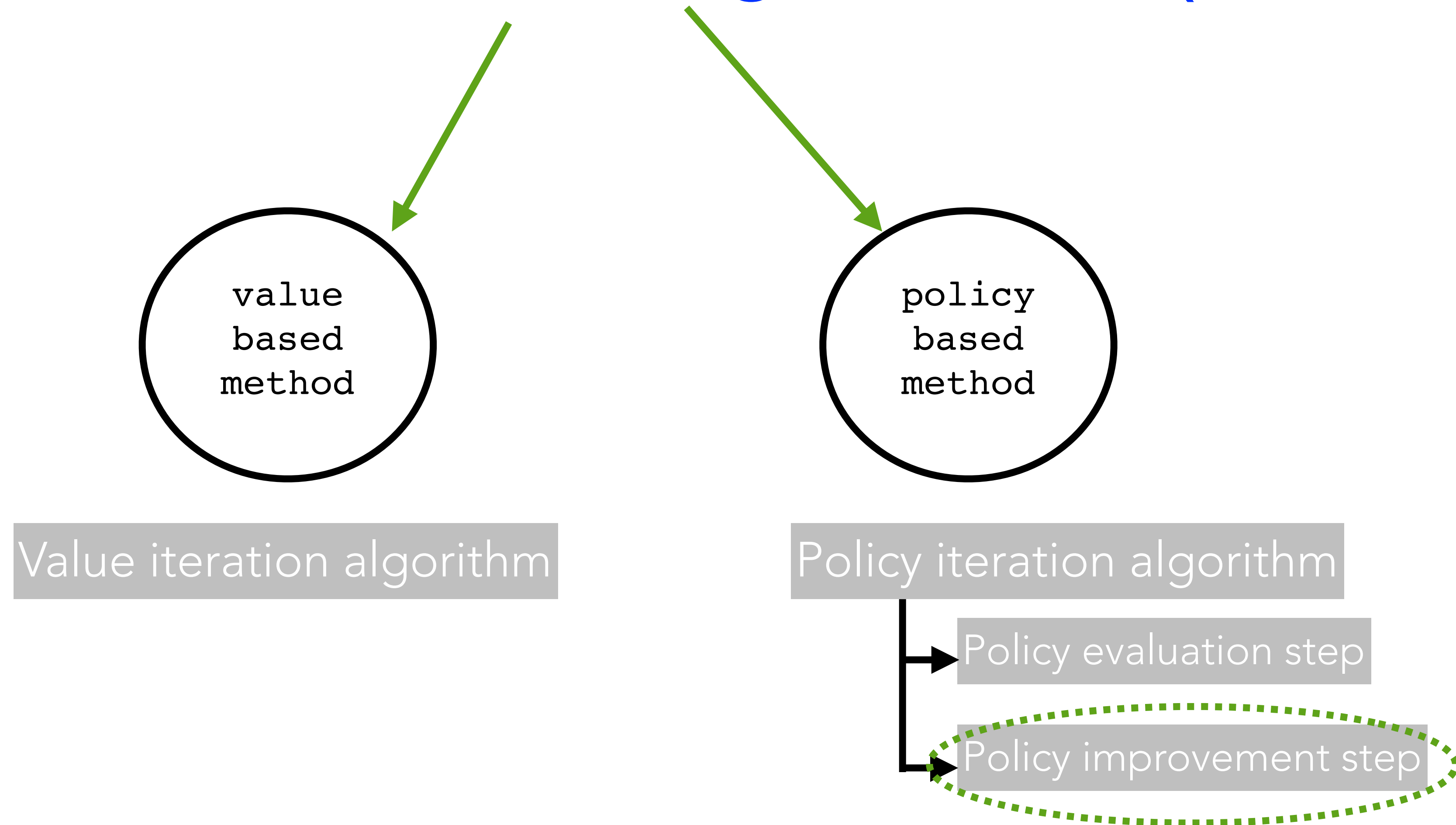
Returns: $(MDP) \rightarrow V_{opt}$ and π_{opt}

Returns: $(MDP, \pi) \rightarrow V_{\pi}$ and $?$

Explore the MDP and do two things:
i) calculate state-value of each state
ii) find the best policy

A (non-optimal) policy is
given to you, you don't need
to find it

Reinforcement learning solvers (so far)



Policy Iteration Algorithm

Iterate $j=1, 2, 3, \dots, T$ steps

- **Step 1**: For an MDP and **its given policy**, apply **policy evaluation algorithm** repeatedly until the the state-value and action-values converge for each state
- **Step 2**: **update the policy** using one step look-ahead

Policy Iteration Algorithm

Iterate $j=1, 2, 3, \dots$ T steps

- **Step 1:** For an MDP and its given policy, apply policy evaluation algorithm repeatedly until the the state-value and action-values converge for each state

$$V_{\pi_j}^t(s) = \sum_{s'} p(s' | s, \pi_j(s)) [R(s, \pi_j(s), s') + \gamma V_{\pi_j}^{t-1}(s')]$$

- **Step 2:** update the policy using one step look-ahead

$$\pi_{j+1}(s) = \arg \max_a \sum_{s'} p(s' | s, a) [R(s, a, s') + \gamma V_{\pi_j}(s')]$$