

# Reinforcement Learning IV

Lecture 22

# Learning strategy

Model-based  
(planning)

Model-free

## Reinforcement Learning

Knowledge of **Environment**

**No knowledge**  
Must learn from  
experience

**Monte Carlo Learning**

**Perfect  
knowledge**  
Known MDP

**Dynamic Programming**  
Policy iteration  
Value iteration

# Reinforcement Learning

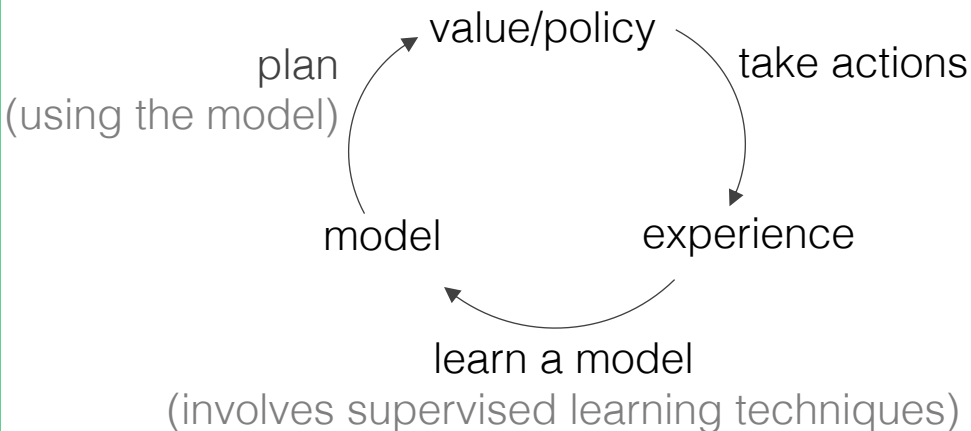
## Learning strategy

Model-based  
(planning)

Model-free

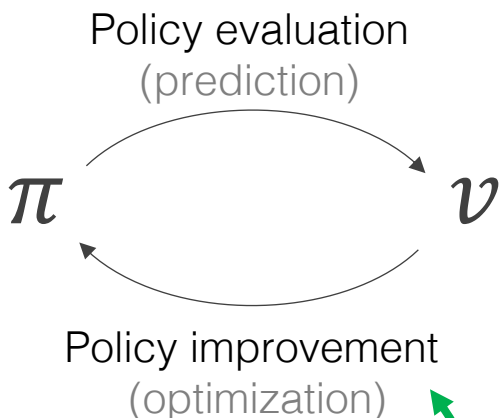
### Reinforcement Learning

Simulation-based search



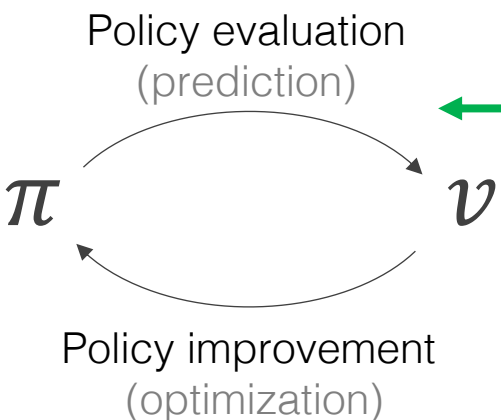
Monte Carlo Learning

Temporal Difference Learning



**Dynamic Programming**

Policy iteration  
Value iteration



**This class:**

1. How to **compute optimal policies** for known MDPs?
2. How to extend this to the case **without full knowledge** of MDPs (model-free learning)

Knowledge of **Environment**

**No knowledge**  
Must learn from experience

**Perfect knowledge**  
Known MDP

# From Dynamic Programming to true RL

We assume a **fully known MDP environment**

(Markov Decision Process)

- |   |                                     |
|---|-------------------------------------|
| 1. What returns will a policy yield?      | <b>Policy evaluation</b>            |
| 2. How can we find a better policy?       | <b>Policy improvement</b>           |
| 3. How do we find the best policy?        | <b>Policy iteration</b>             |
| 4. How do we find the best policy faster? | <b>Value iteration</b>              |
| 5. Are there other approaches?            | <b>Generalized Policy Iteration</b> |
| What if we don't have a fully known MDP?  | <b>Monte Carlo Methods</b>          |

# Markov Decision Process

## Components:

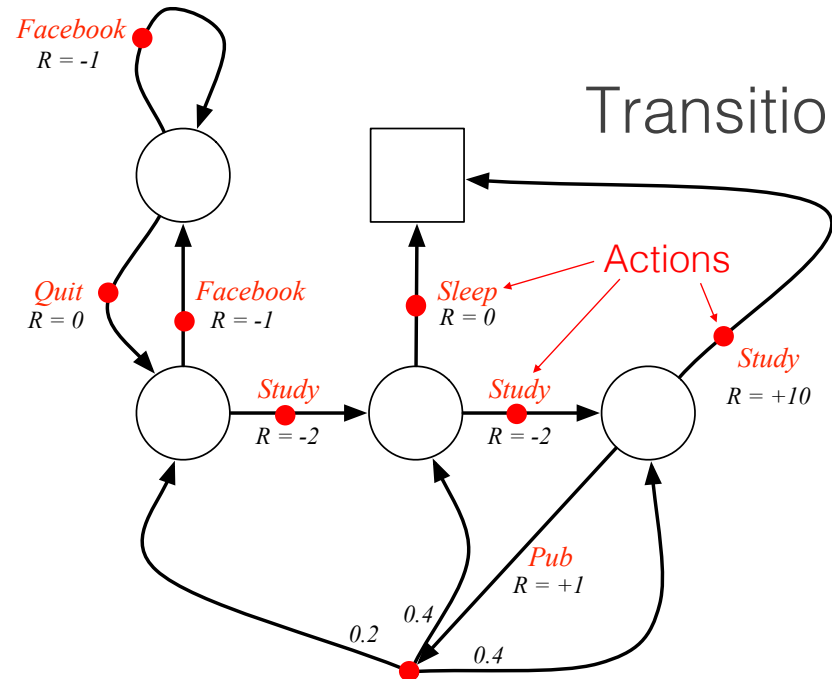
State space  $\mathcal{S}$

Transition probabilities,  $P$

Rewards,  $R$

Discount rate,  $\gamma$

Actions,  $\mathcal{A}$



## Returns (Expected future rewards)

(discount factor weights the the future)

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots$$

Policy (how we choose actions)

(can be stochastic or deterministic)

$$\pi(a|s) = P(a|s)$$

## State value function

(expected return from state  $s$ , and following policy  $\pi$ )

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

$$v_{\pi}(s) = E[R_s^a + \gamma v_{\pi}(s')|s]$$

## Action value function

(expected return from state  $s$ , taking action  $a$ , and following policy  $\pi$ )

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

$$q_{\pi}(s, a) = E[R_s^a + \gamma q_{\pi}(s', a')|s, a]$$

David Silver, UCL, 2015

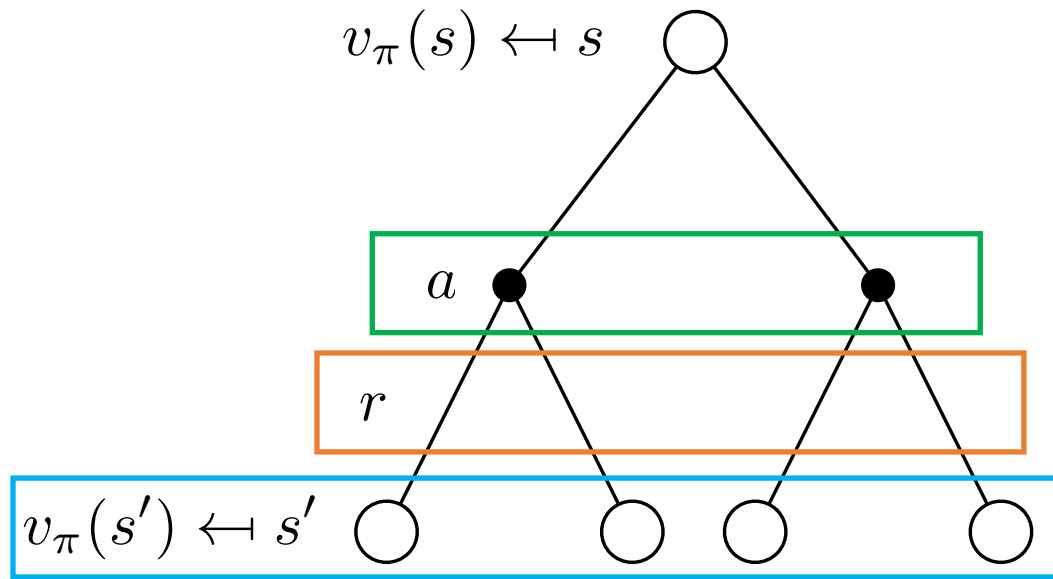
# Bellman Expectation Equations for the **state value** function

(expected return from state  $s$ , and following policy  $\pi$ )

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

$$v_{\pi}(s) = E[R_s^a + \gamma v_{\pi}(s') | s]$$

$$R_s^a = E[r_{t+1} | S_t = s, A_t = a]$$



Expectation over the possible actions

Expectation over the rewards

(based on state and choice of action)

Expectation over the next possible states

$$v_{\pi}(s) = \underbrace{\sum_a}_{\text{Expectation over actions}} \underbrace{\pi(a|s)}_{\text{Expectation over rewards}} \left( \underbrace{R_s^a + \gamma \sum_{s'} P_{ss'}^a v_{\pi}(s')}_{\text{Expectation over next states}} \right)$$

# Bellman Expectation Equations for the **action value** function

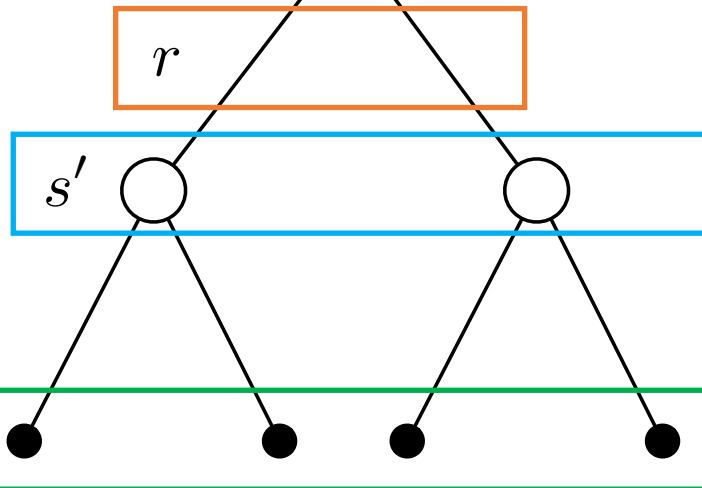
(expected return from state  $s$ , taking action  $a$ , then following policy  $\pi$ )

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

$$q_{\pi}(s, a) = E[R_s^a + \gamma q_{\pi}(s', a') | s, a]$$

$$R_s^a = E[r_{t+1} | S_t = s, A_t = a]$$

$$q_{\pi}(s, a) \leftarrow s, a$$



Expectation over the rewards

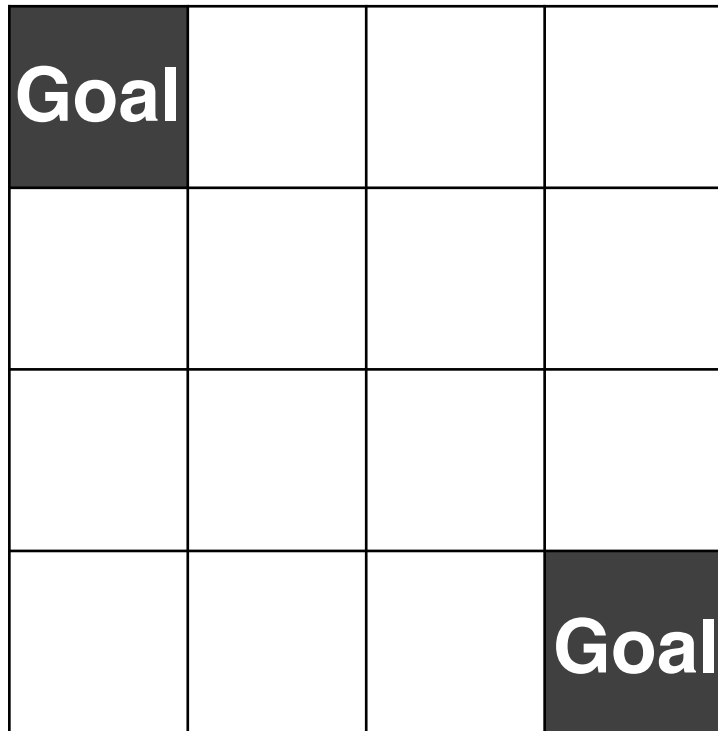
(based on state and choice of action)

Expectation over the next possible states

Expectation over the possible actions

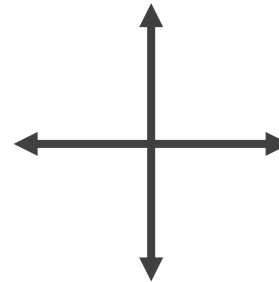
$$q_{\pi}(s, a) = \underbrace{R_s^a}_{\text{orange}} + \gamma \underbrace{\sum_{s'} P_{ss'}^a}_{\text{blue}} \underbrace{\sum_{a'} \pi(a' | s') q_{\pi}(s', a')}_{\text{green}}$$

# Running example: Gridworld



16 states and 2 terminal states labeled “goal”

Valid actions:  
(unless there is a wall)



Reward:

-1 for all transitions  
(until the terminal state has been reached)

Note: actions that would take the agent off the board are not allowed

Sutton and Barto, 2018



# 1. Policy Evaluation

Input: policy  $\pi(a|s)$

Output: value function  $v_\pi(s)$   
(unknown)

- 1 Select a policy function to evaluate (estimate the value function)
- 2 Start with a guess of the value function,  $v_0$  (often all zeros)
- 3 **Iteratively** apply the Bellman Expectation Equation to “backup” the values until they converge on the actual value function for the policy,  $v_\pi$

$$v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_\pi$$

Adapted from David Silver, 2015

# 1. Policy Evaluation

in Gridworld

Policy:

$$\pi(a|s) = \frac{1}{N_{\text{valid\_actions}}}$$

Randomly go in any valid direction

Value function initialization:

$$v_0(s) = 0 \text{ (all zeros)}$$

$v_0(s)$   
(initialization)

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

We estimate the value function that corresponds to the policy:  $v_\pi(s)$

# 1. Policy Evaluation in Gridworld

Policy:  $\pi(a|s) = 1/N_{\text{valid\_actions}}$   
(randomly go in any direction)

$v_0(s)$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

Bellman Expectation Equation:

$$v_{k+1}(s) = \sum_a \pi(a|s) \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a v_k(s') \right)$$

In Gridworld:  $\frac{1}{N_a}$   $-1$   $1$  (once you pick an action there's no uncertainty as to which state you'll transition to)

$$v_{k+1}(s) = \sum_a \frac{1}{N_a} \left( -1 + \sum_{s'} v_k(s') \right) = -1 + \sum_a \frac{1}{N_a} \sum_{s'} v_k(s')$$

Each action leads to only one state, so the sum over states is not needed

$$= -1 + \sum_a \frac{1}{N_a} v_k(s')$$

Average of the value of the  $N_a$  neighboring states

# 1. Policy Evaluation

in Gridworld

$$v_{k+1}(s) = -1 + \sum_a \frac{1}{N_a} v_k(s')$$

$$v_1 = -1 + \sum_a \frac{1}{4} v_k(s') = -1$$

$v_0(s)$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

One neighborhood  
in  $v_0(s)$

	0		
0		0	
	0		

$v_1(s)$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

$v_0(s)$ 

<b>0</b>	0	0	0
0	0	0	0
0	0	0	0
0	0	0	<b>0</b>

 $v_1(s)$ 

<b>0</b>	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	<b>0</b>

 $v_2(s)$ 

<b>0</b>	-1.7	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	<b>0</b>

 $v_3(s)$ 

<b>0</b>	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	<b>0</b>

 $v_{10}(s)$ 

<b>0</b>	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	<b>0</b>

 $v_{\infty}(s) = v_{\pi}(s)$ 

<b>0</b>	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	<b>0</b>

We've found the value function  
(expected returns) from our random  
movement policy

# 1. Policy Evaluation

in Gridworld

## 2. Policy Improvement

Input: policy  $\pi(a|s)$   
Output: better policy  $\pi'(a|s)$

Definition of better: has greater or equal expected return in all states:  
 $v_{\pi'}(s) \geq v_{\pi}(s)$  for all states

- 1 Select a policy function to improve
- 2 Evaluate the value function (our last discussion)
- 3 **Greedily** select a new policy,  $\pi'$ , that chooses actions that maximize value

$$\pi'(s) = \text{greedy}(s)$$

(i.e. pick the action that brings us to the state with highest value)

Adapted from David Silver, 2015

## 2. Policy Improvement

Input: policy  $\pi(a|s)$   
Output: better policy  $\pi'(a|s)$

**How do we do this:**  $\pi'(s) = \text{greedy}(s)$

i.e. pick the **action** that brings us to the state with **highest value**

We can use the state value function to help us choose the right action:

$$\pi'(s) = \arg \max_a q_{\pi}(s, a)$$

Reminder:

Action value function

(expected return from state  $s$ , taking action  $a$ , and following policy  $\pi$ )

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

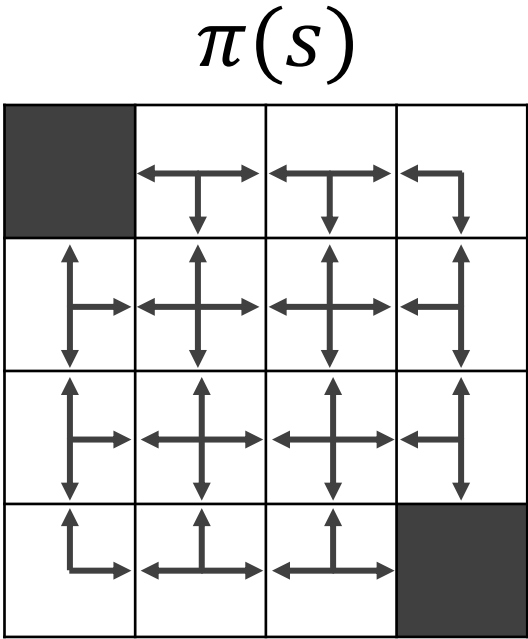
# 2. Policy Improvement

in Gridworld

Value function:

In this case,  $q_{\pi}(s, \pi(s)) = v_{\pi}(s)$  since each action leads to only one state

Initial policy:  
 $\pi(a|s)$  = randomly go in any valid direction



$$v_0(s)$$

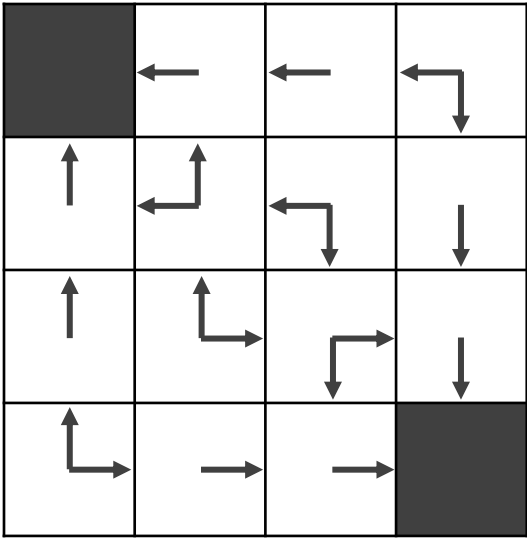
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$$v_{\infty}(s) = v_{\pi}(s)$$

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

Improved policy:  $\pi'(s)$

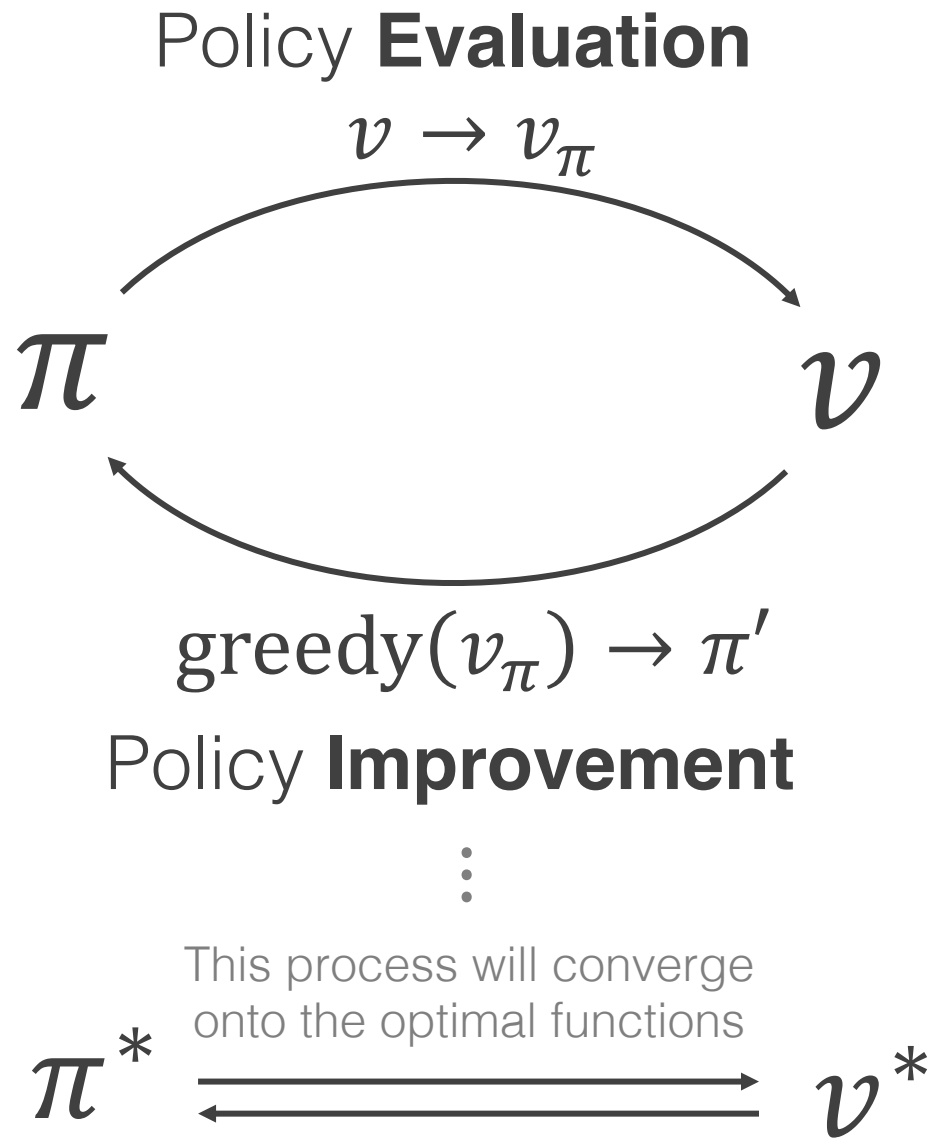
(in this case this is an optimal policy)





# 3. Policy Iteration

Input: policy  $\pi(a|s)$   
Output: **best** policy  $\pi^*(a|s)$   
Best in the sense that:  $v_{\pi^*}(s) \geq v_{\pi}(s)$  for all states and for all policies

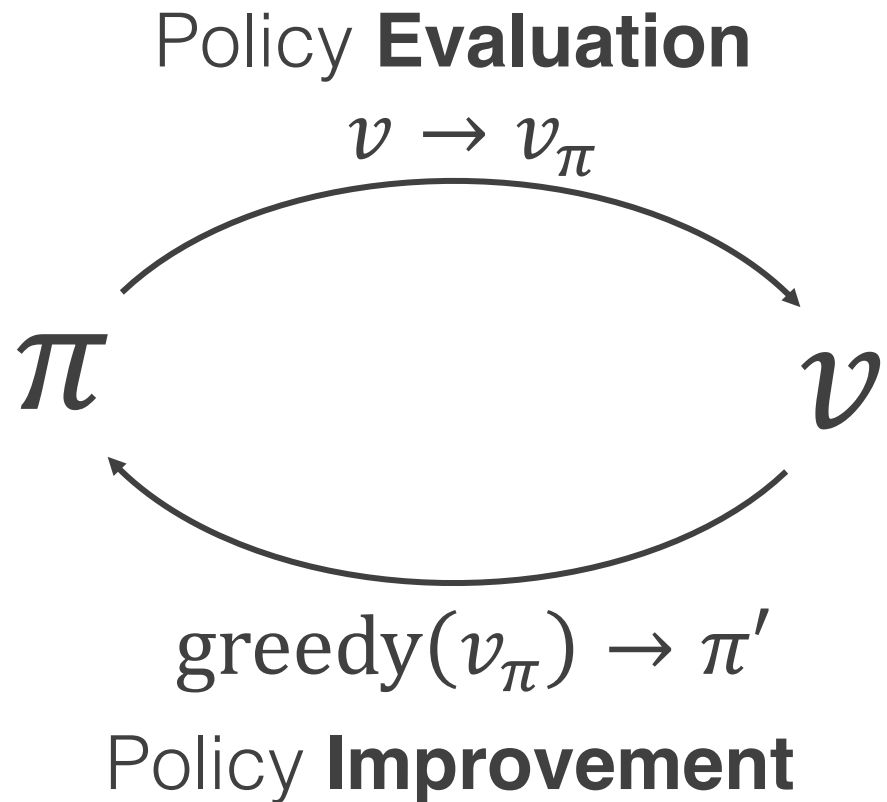


$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \rightarrow \dots \xrightarrow{I} \pi^* \xrightarrow{E} v^*$$

Adapted from David Silver, 2015 and Sutton and Barto, 1998

### 3. Policy Iteration

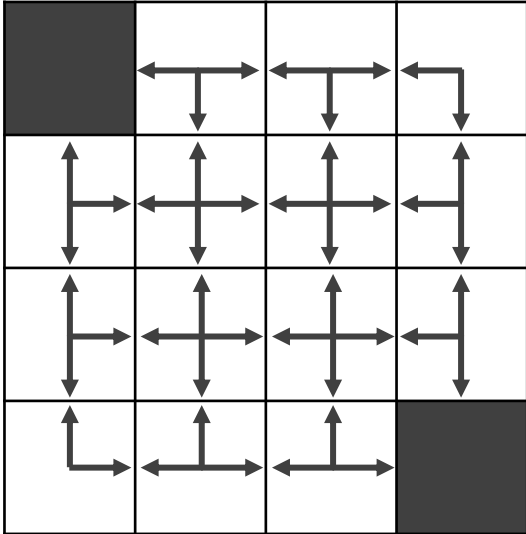
Input: policy  $\pi(a|s)$   
Output: **best** policy  $\pi^*(a|s)$



- 1 Policy Evaluation:** estimate  $v_\pi$   
Iterative policy evaluation  
Note: This is VERY slow
- 2 Policy Improvement:** generate  $\pi' \geq \pi$   
Greedy policy improvement
- 3** Iterate 1 and 2 until convergence

Adapted from David Silver, 2015 and Sutton and Barto, 1998

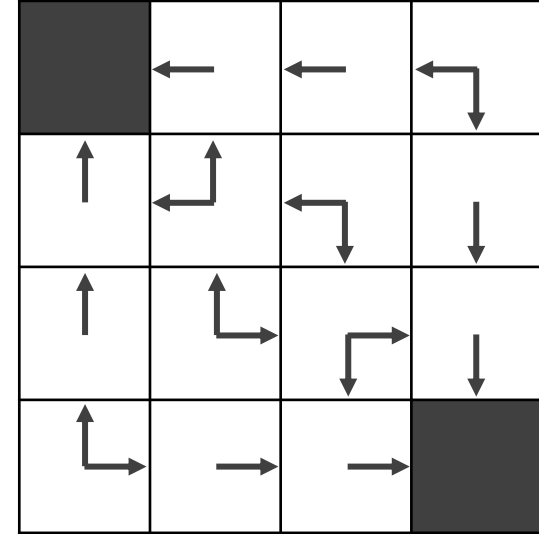
1  $\pi_0(s)$



$v_0(s)$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

3  $\pi_1(s) = \pi^*(s)$



$v_0(s)$

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

Policy Evaluation  
Policy Improvement

2  $v_\infty(s) \rightarrow v_{\pi_0}(s)$

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

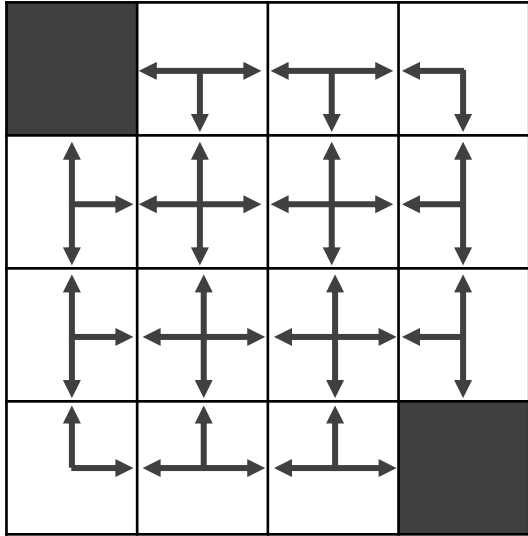
Policy Evaluation

4  $v_\infty(s) \rightarrow v_{\pi_0}(s)$

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

$v_{\pi^*}(s)$

$$\pi_0(s)$$



$$v_0(s)$$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$$v_1(s)$$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

What if we stopped after one sweep. This is...

## 4. Value Iteration

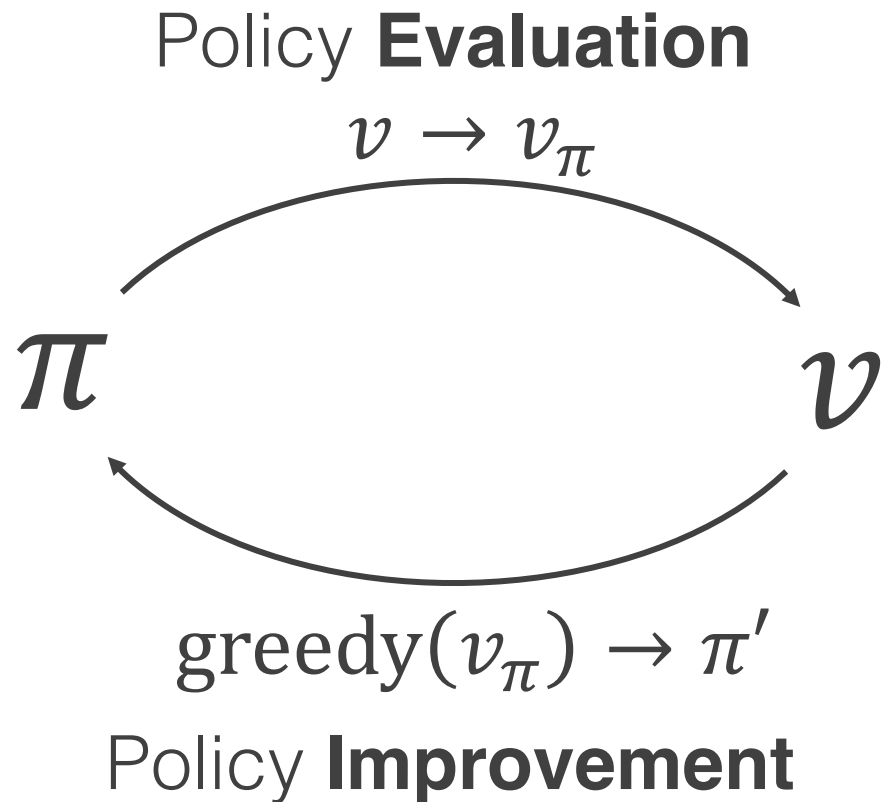
# 4. Value Iteration

Input: policy

$\pi(a|s)$

Output: **best** policy

$\pi^*(a|s)$



- 1 Policy Evaluation:** estimate  $v_\pi$   
**One-sweep** of policy evaluation
- 2 Policy Improvement:** generate  $\pi' \geq \pi$   
Greedy policy improvement
- 3** Iterate 1 and 2 until convergence

Adapted from David Silver, 2015 and Sutton and Barto, 1998

$v_0(s)$ 

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

 $v_1(s)$ 

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

 $v_2(s)$ 

0	-1.7	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	0

 $v_3(s)$ 

0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0

 $v_{10}(s)$ 

0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0

 $v_\infty(s) = v_\pi(s)$ 

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

So far, we've run policy evaluation all the way to convergence (**this is slow**)

1  $v_0(s)$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

3  $v_1(s)$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

5  $v_2(s)$

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

7  $v_3(s) = v_{\pi^*}(s)$

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

2  $\pi_0(s)$

	↔	↔	↖
↑	↕	↕	↕
↑	↕	↕	↕
↙	↕	↕	

4  $\pi_1(s)$

	←	↔	↖
↑	↕	↕	↕
↑	↕	↕	↓
↙	↕	→	

6  $\pi_2(s)$

	←	←	↖
↑	↖	↕	↓
↑	↕	↘	↓
↙	→	→	

8  $\pi_3(s) = \pi^*(s)$

	←	←	↖
↑	↖	↕	↓
↑	↕	↘	↓
↙	→	→	

# Demo

[https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld\\_dp.html](https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html)



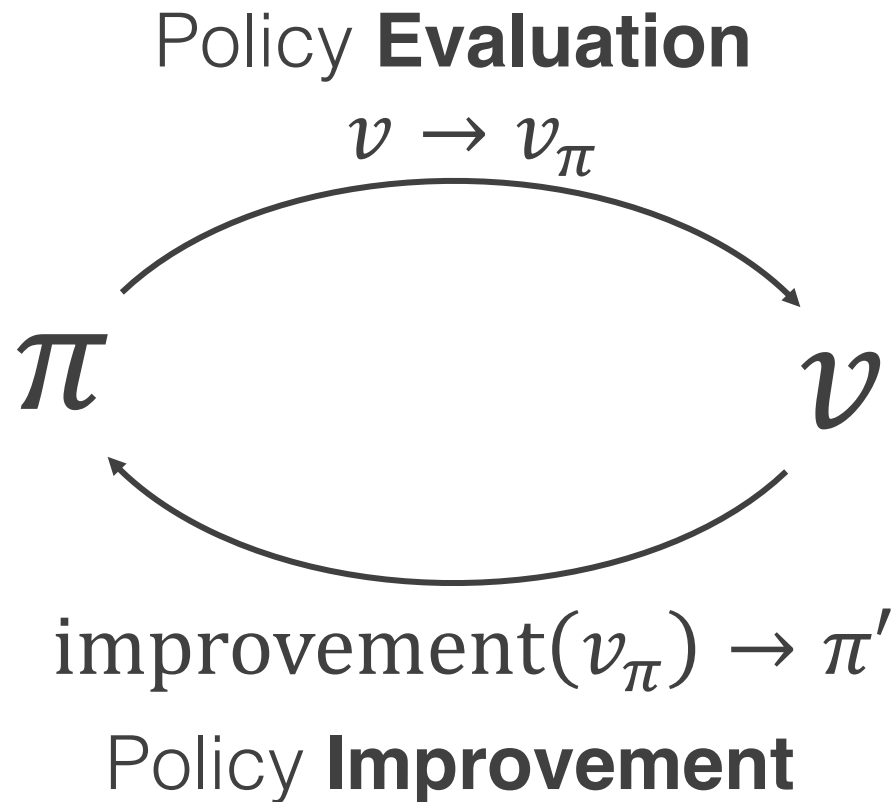
# 5. Generalized Policy Iteration

Input: policy

$\pi(a|s)$

Output: **best** policy

$\pi^*(a|s)$



- 1 Policy Evaluation:** estimate  $v_\pi$   
**Any** policy evaluation algorithm
- 2 Policy Improvement:** generate  $\pi' \geq \pi$   
**Any** policy improvement algorithm
- 3** Iterate 1 and 2 until convergence

Adapted from David Silver, 2015 and Sutton and Barto, 1998

So far, we've assumed full knowledge of the environment (MDP)

What if we DO NOT assume full knowledge of the environment (MDP)

This means we have to learn by experience:  
**true reinforcement learning**

## 6. Monte Carlo Policy Evaluation

For **state** values

Input: policy  $\pi(a|s)$   
Output: state value  $v_\pi(s)$

- 1 Select a policy function to evaluate (estimate the value function)
- 2 Start with a guess of the value function,  $v_0$  (often all zeros)
- 3 Repeat forever:
  - A Generate an episode (take actions until a terminal state)
  - B Save the returns following the first occurrence of each state
  - C Assign  $\text{AVG}(\text{Returns}(s)) \rightarrow \hat{v}_\pi(s)$

Sutton and Barto, 1998

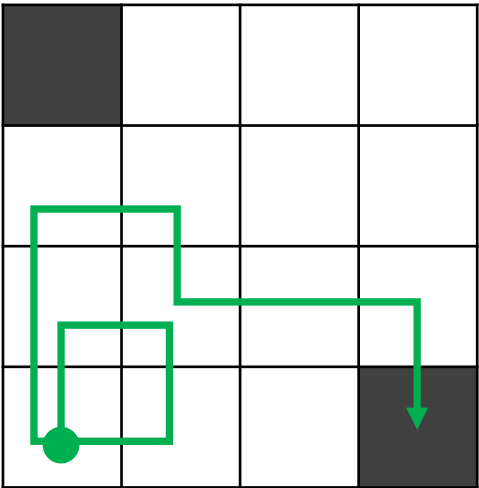
# 6. Monte Carlo Policy Evaluation

For **state** values  
"First Visit"

For each state, we store the running returns seen **after** the first visit to that state

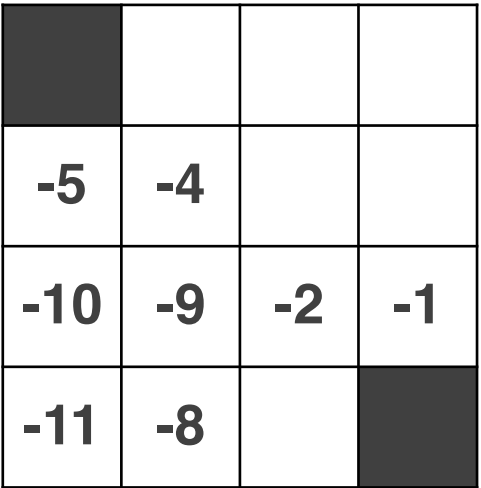
## Episode 1

Total Reward: -11



Episode 1 **returns** after the first visit of each state

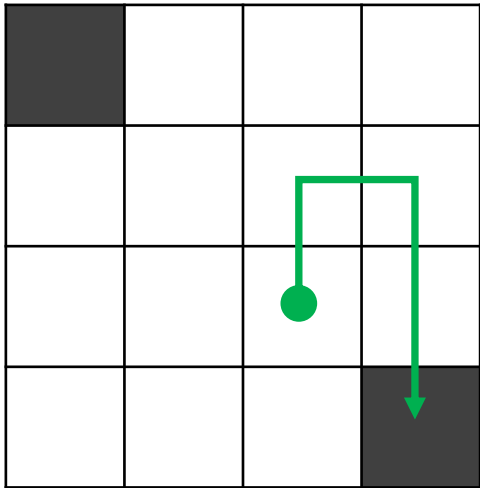
$G^{(1)}$



Discount rate:  $\gamma = 1$

## Episode 2

Total Reward: -4



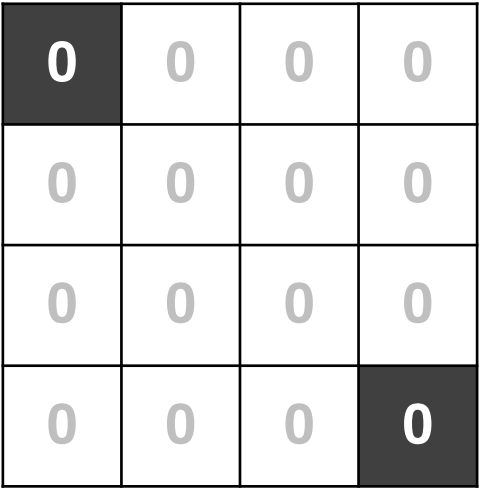
Episode 2 **returns** from the first visit of each state

$G^{(2)}$



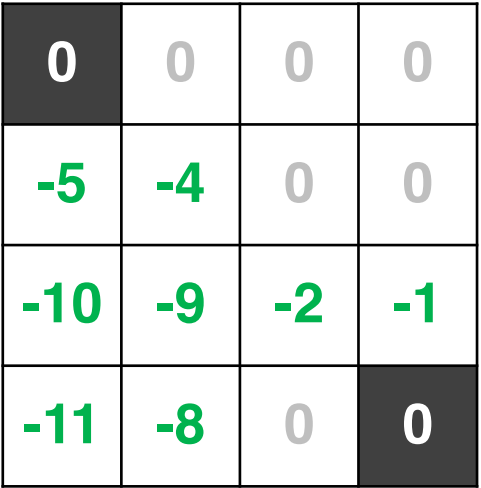
Discount rate:  $\gamma = 1$

$$v_0(s)$$



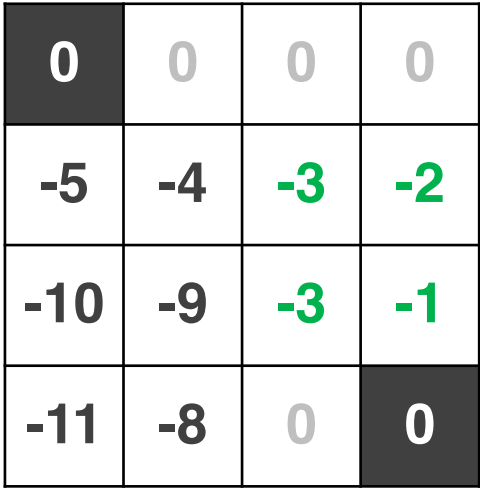
The value function is the **running average** of the returns after the visit to that state, averaged over episodes (only average over episodes when state is visited)

$$v_1(s)$$



$v_1$  is just the first visit returns,  $G^{(1)}$

$$v_2(s)$$



$v_2$  is the average first visit returns,  $G^{(1)}$  and  $G^{(2)}$ , for those states visited

# Which value function?

The **state value function** doesn't tell us directly about actions

If we don't have a model, to pick a policy we need **action values**

# Which value function?

Greedy policy improvement over  $v(s)$  **requires a model of the MDP**

$$\pi'(s) = \operatorname{argmax}_a R_s^a + P_{ss'}^a v(s')$$

Greedy policy improvement over  $q(s, a)$  is **model-free**

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

And the two value functions are related:  $v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$

## 6. Monte Carlo Policy Evaluation

Input: policy  $\pi(a|s)$   
Output: action value  $q_\pi(s, a)$

For **action** values

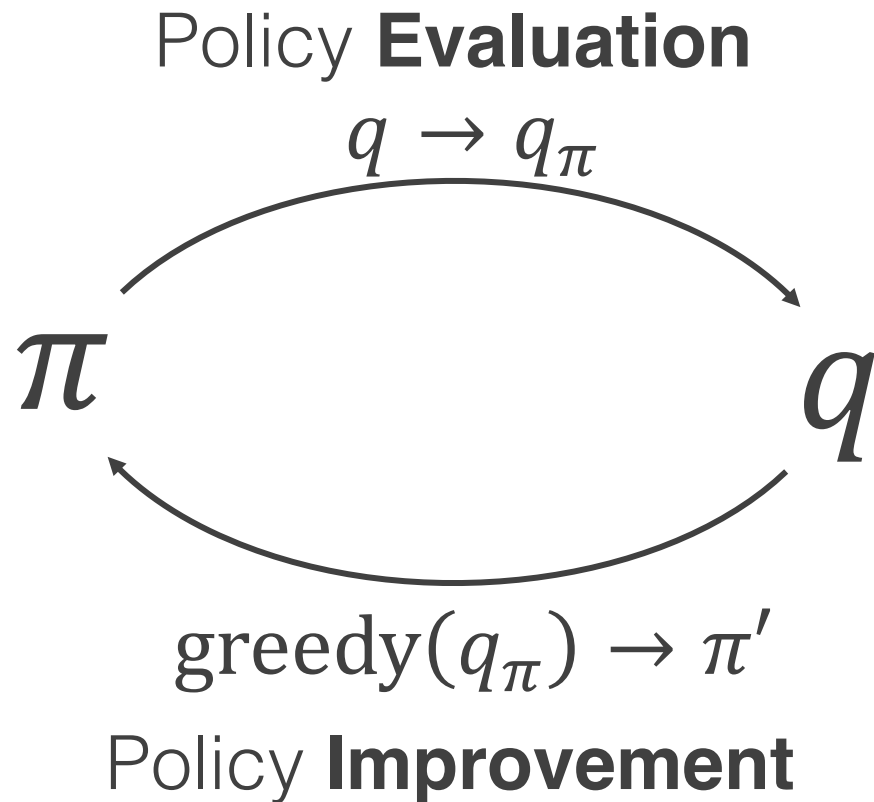
- 1 Select a policy function to evaluate (estimate its value function)
- 2 Start with a guess of the action value function,  $q_0$  (often all zeros)
- 3 Repeat forever:
  - A Generate an episode (take actions until a terminal state)
  - B Save returns following first occurrence of each state **& action**
  - C Assign  $\text{AVG}(\text{Returns}(s, a)) \rightarrow \hat{q}_\pi(s, a)$

Sutton and Barto, 1998

# 7. Monte Carlo Control

(policy iteration)

Input: policy  $\pi(a|s)$   
Output: **best policy**  $\pi^*(a|s)$



- 1 Policy Evaluation:** estimate  $q_\pi$   
**Monte Carlo action policy evaluation**
- 2 Policy Improvement:** generate  $\pi' \geq \pi$   
Greedy policy improvement
- 3** Iterate 1 and 2 until convergence

Sutton and Barto, 1998



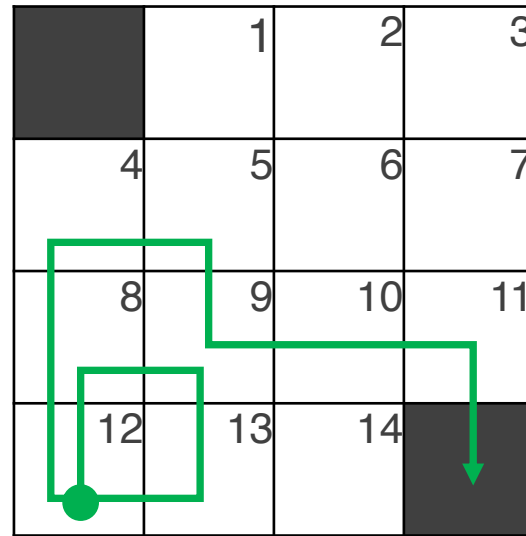
# 7. Monte Carlo Control

“First Visit” (of state AND action) is recorded

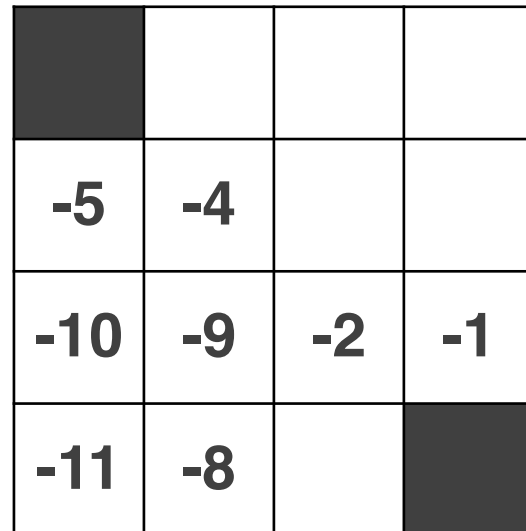
**Episode 1**  
Total Reward: -11

Episode 1 **returns** after  
the first visit of each state

Discount rate:  $\gamma = 1$



**1** MC Policy Evaluation



$q_{\pi}(s, a)$

Action ( $a$ ):  $\uparrow$   $\rightarrow$   $\leftarrow$   $\downarrow$

State ( $s$ )

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

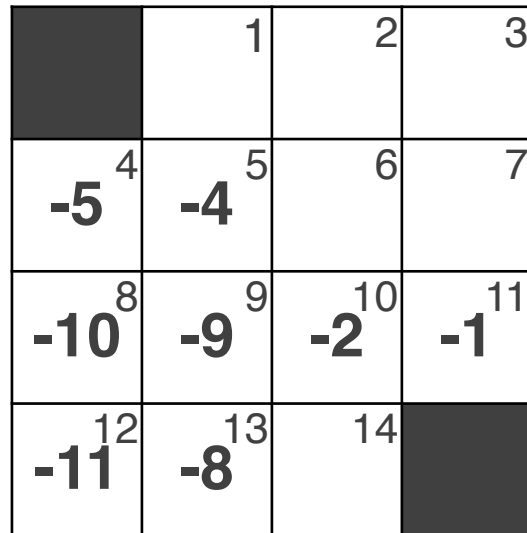
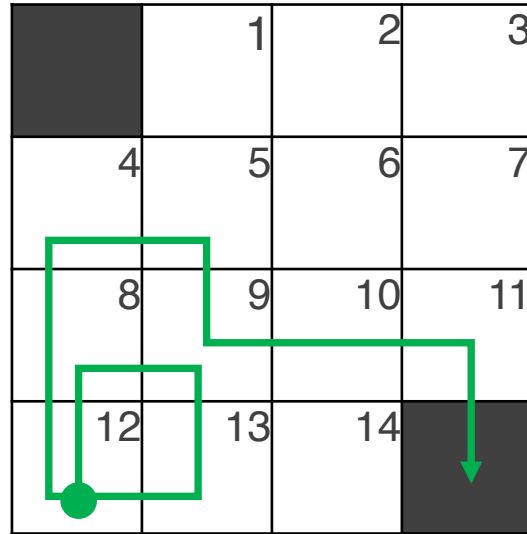
# 7. Monte Carlo Control

“First Visit” (of state AND action) is recorded

**Episode 1**  
Total Reward: -11

Episode 1 **returns** after  
the first visit of each state

Discount rate:  $\gamma = 1$



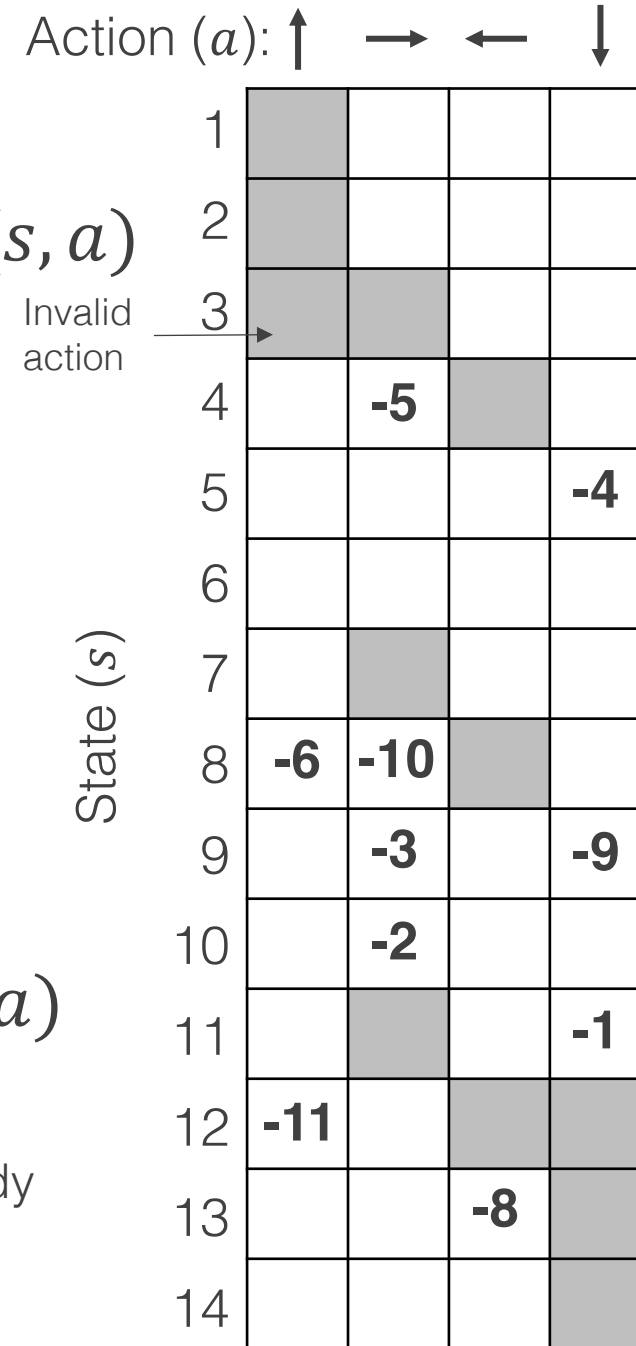
- 1 MC Policy Evaluation
- 2 MC Policy Improvement

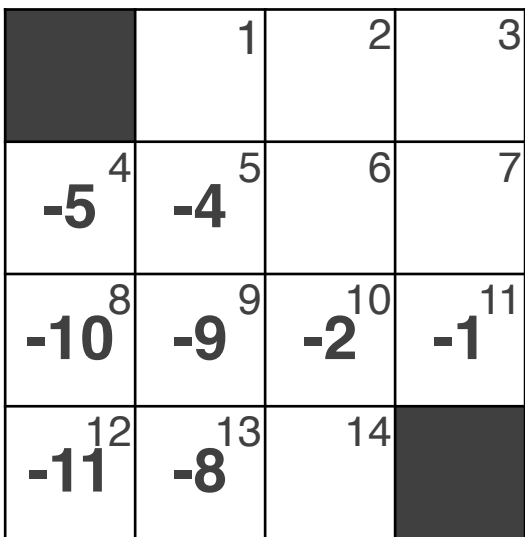
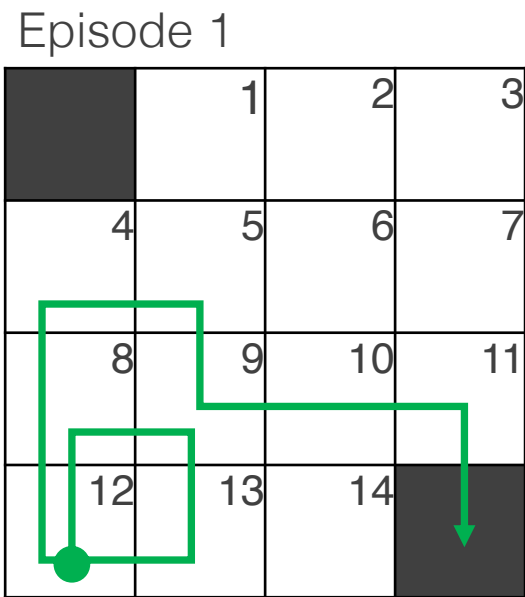
$$\pi'(s) = \operatorname{argmax}_a q_{\pi}(s, a)$$

Typically this is set to be  $\epsilon$ -greedy  
to better learn  $q(s, a)$

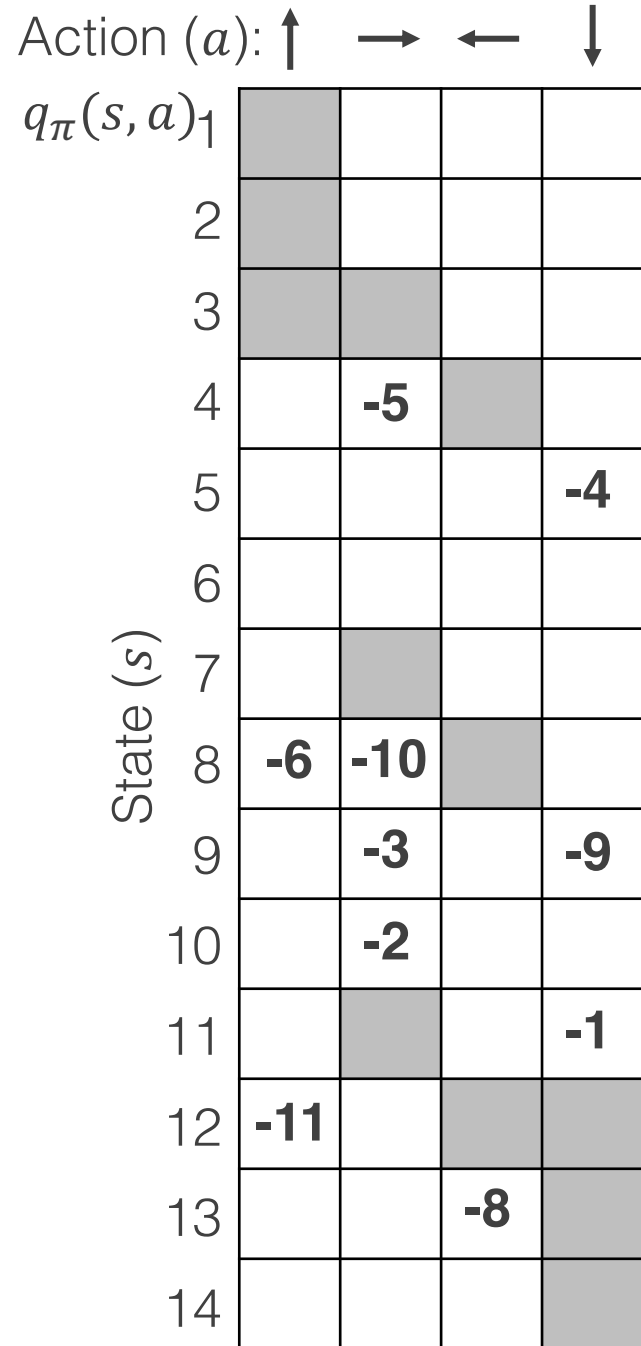
$$q_{\pi}(s, a)$$

Invalid  
action

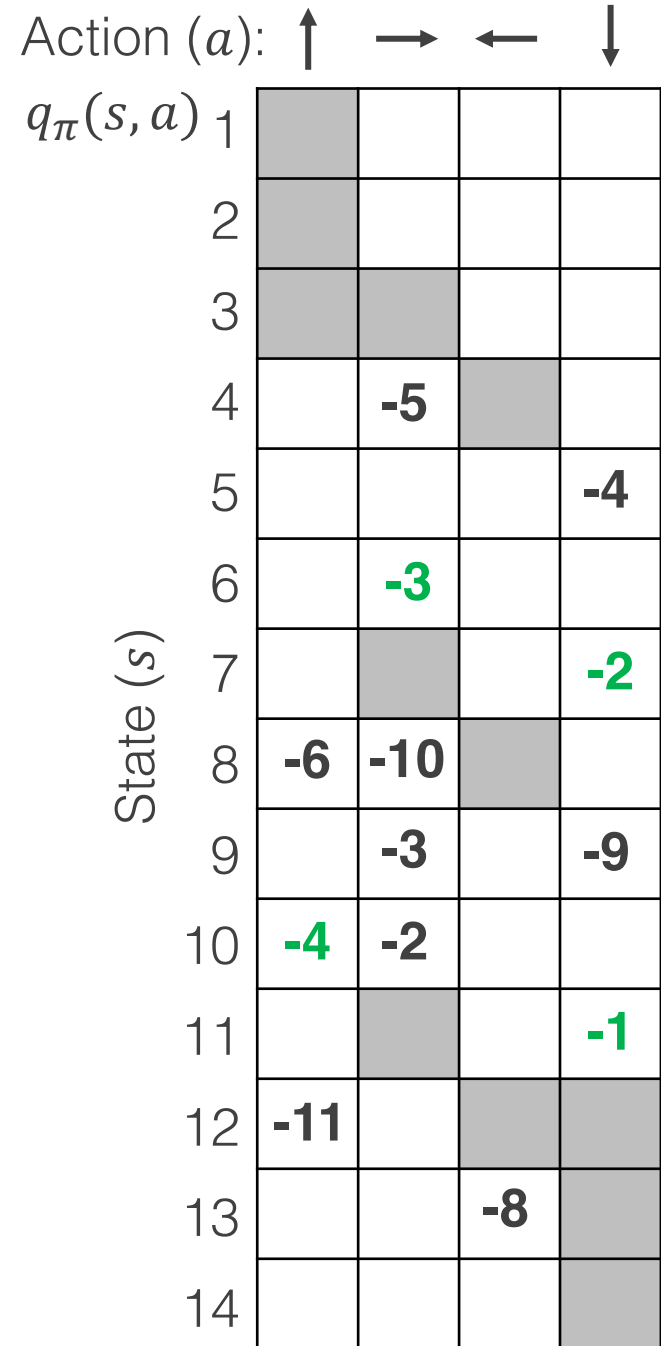




Returns from episode 1



Returns from episode 2



Action ( $a$ ):  $\uparrow \rightarrow \leftarrow \downarrow$

$q_{\pi^*}(s, a)$

State ( $s$ )

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

If we're in state 4,  
take the up action

If we know the  
optimal action  
value function,  
we also have our  
optimal policy



$v_{\pi^*}(s)$

<b>0</b>	-1 <sup>1</sup>	-2 <sup>2</sup>	-3 <sup>3</sup>
-1 <sup>4</sup>	-2 <sup>5</sup>	-3 <sup>6</sup>	-2 <sup>7</sup>
-2 <sup>8</sup>	-3 <sup>9</sup>	-2 <sup>10</sup>	-1 <sup>11</sup>
-3 <sup>12</sup>	-2 <sup>13</sup>	-1 <sup>14</sup>	<b>0</b>

$\pi^*(s)$

	$\leftarrow$ <sup>1</sup>	$\leftarrow$ <sup>2</sup>	$\swarrow$ <sup>3</sup>
$\uparrow$ <sup>4</sup>	$\swarrow$ <sup>5</sup>	$\updownarrow$ <sup>6</sup>	$\downarrow$ <sup>7</sup>
$\uparrow$ <sup>8</sup>	$\left\langle\right\rangle$ <sup>9</sup>	$\searrow$ <sup>10</sup>	$\downarrow$ <sup>11</sup>
$\swarrow$ <sup>12</sup>	$\rightarrow$ <sup>13</sup>	$\rightarrow$ <sup>14</sup>	

# Extensions

Monte Carlo methods require that we finish each episode before updating

**Solution:** **Temporal Difference** (TD) methods

What if we want to learn about one policy while following or observing another?

**Solution:** **Off-policy learning** instead of on-policy learning

What if our state space has too many states that we can't build a table of values?

**Solution:** **Value function approximation** (involving supervised learning techniques)

How can we simulate what the environment might output for next states and rewards?

**Solution:** **Model-based learning**: simulate the environment and plan ahead

# From Dynamic Programming to true RL

We assume a **fully known MDP environment**

(Markov Decision Process) **Sutton & Barto, Chapter 3**

- |   |                                     |
|---|-------------------------------------|
| 1. What returns will a policy yield?      | <b>Policy evaluation</b>            |
| 2. How can we find a better policy?       | <b>Policy improvement</b>           |
| 3. How do we find the best policy?        | <b>Policy iteration</b>             |
| 4. How do we find the best policy faster? | <b>Value iteration</b>              |
| 5. Are there other approaches?            | <b>Generalized Policy Iteration</b> |

What if we don't have a fully known MDP?

**Monte Carlo Methods**

**Sutton & Barto, Chapter 5**

# Reinforcement Learning

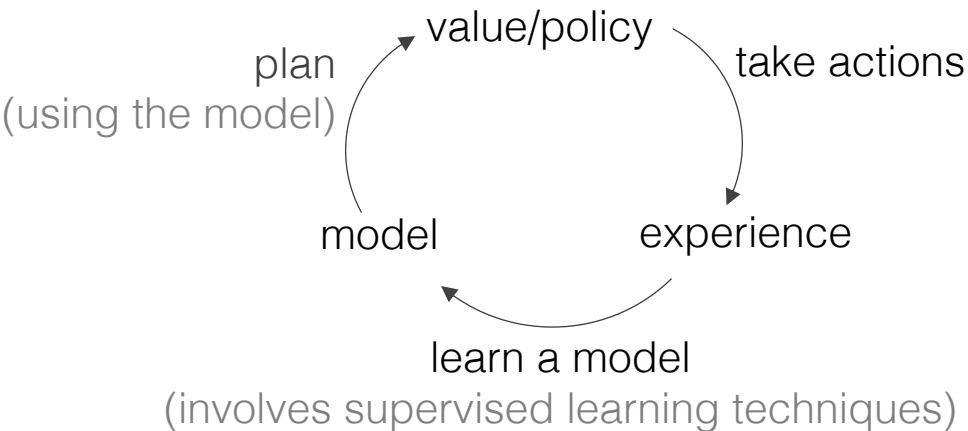
## Learning strategy

Model-based  
(planning)

Model-free

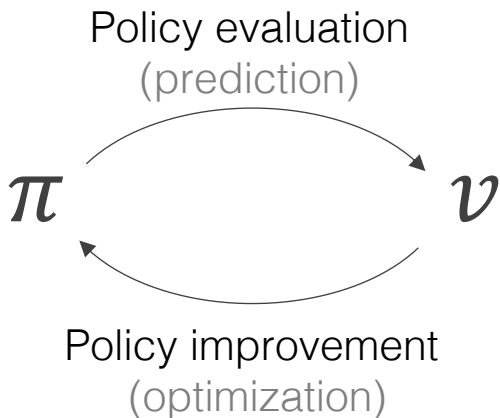
## Reinforcement Learning

Simulation-based search



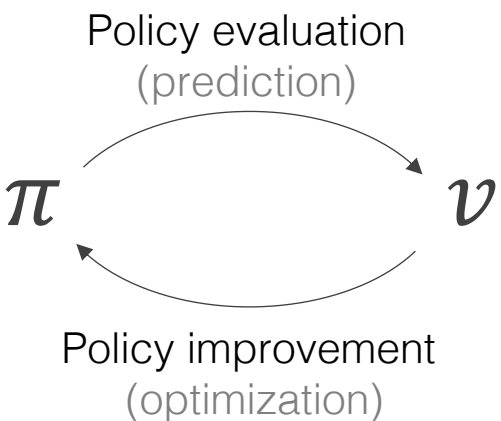
### Monte Carlo Learning

Temporal Difference Learning



### Dynamic Programming

Policy iteration  
Value iteration



Knowledge of Environment

**No knowledge**  
Must learn from experience

**Perfect knowledge**  
Known MDP