

# Chapter 3: Markov Decision Process (MDP)

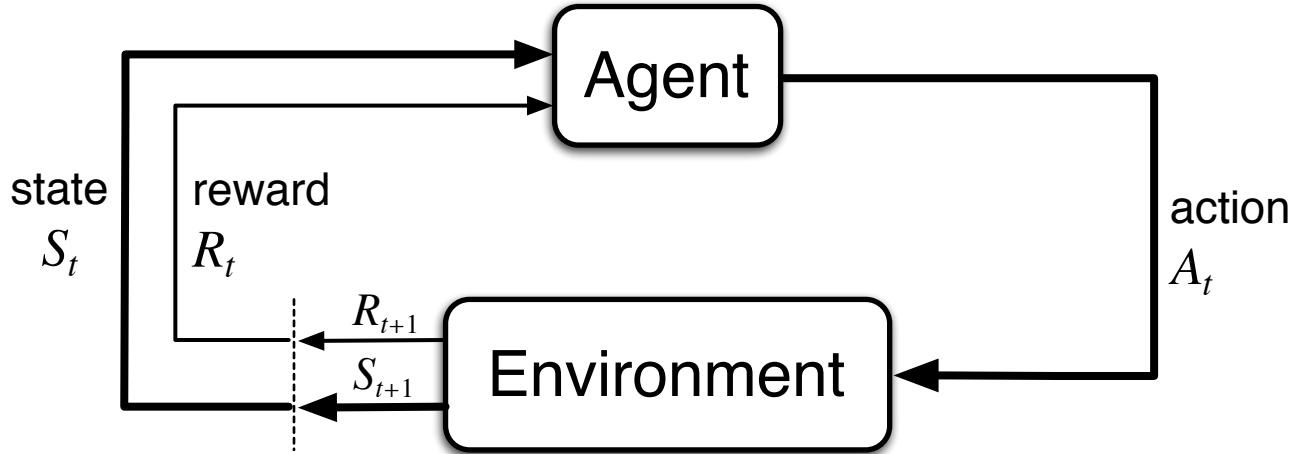
## The Reinforcement Learning Problem

Objectives of this chapter:

- present Markov decision processes—an idealized form of the AI problem for which we have precise theoretical results
- introduce key components of the mathematics: value functions and Bellman equations

# The Agent-Environment Interface

---



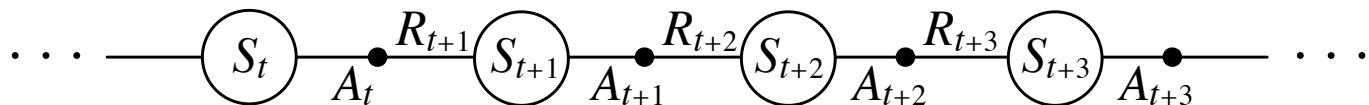
Agent and environment interact at discrete time steps:  $t = 0, 1, 2, 3, \dots$

Agent observes state at step  $t$ :  $S_t \in \mathcal{S}$

produces action at step  $t$ :  $A_t \in \mathcal{A}(S_t)$

gets resulting reward:  $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$

and resulting next state:  $S_{t+1} \in \mathcal{S}^+$



# Goals and Rewards

---

- Is a scalar reward signal an adequate notion of a goal?— maybe not, but it is surprisingly flexible.
- A goal should specify **what** we want to achieve, not **how** we want to achieve it.
- A goal must be outside the agent’s direct control—thus outside the agent.
- The agent must be able to measure success:
  - explicitly;
  - frequently during its lifespan.

# Rewards and returns

- The objective in RL is to maximize long-term future reward
- That is, to choose  $A_t$  so as to maximize  $R_{t+1}, R_{t+2}, R_{t+3}, \dots$
- But what exactly should be maximized?
- Example: the discounted return at time  $t$ :

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \quad \text{the discount rate}$$

$\gamma$	Reward sequence	Return
0.5(or any)	1 0 0 0...	
0.5	0 0 2 0 0 0...	
0.9	0 0 2 0 0 0...	
0.5	-1 2 6 3 2 0 0 0...	

# Return

---

Suppose the sequence of rewards after step  $t$  is:

$$R_{t+1}, R_{t+2}, R_{t+3}, \dots$$

What do we want to maximize?

We seek to **maximize the expected return**,  $E\{G_t\}$ , on each step t:

- Total reward,  $G_t$  = sum of all future reward in the episode
- Discounted reward,  $G_t$  = sum of all future *discounted* reward

# Episodic Tasks

---

**Episodic tasks:** interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze

In episodic tasks, we almost always use simple *total reward*:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T ,$$

where  $T$  is a final time step at which a **terminal state** is reached, ending an episode.

# Continuing Tasks

---

**Continuing tasks:** interaction does not have natural episodes, but just goes on and on...

In this class, for continuing tasks we will always use *discounted return*:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$$

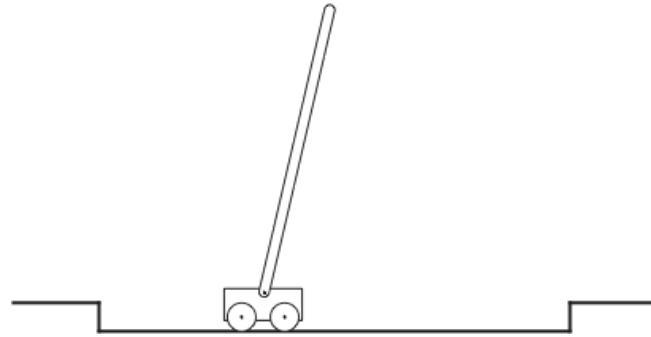
where  $\gamma, 0 \leq \gamma \leq 1$ , is the **discount rate**.

shortsighted  $0 \leftarrow \gamma \rightarrow 1$  farsighted

Typically,  $\gamma = 0.9$

# An Example: Pole Balancing

---



Avoid **failure**: the pole falling beyond a critical angle or the cart hitting end of track

As an **episodic task** where episode ends upon failure:

reward = +1 for each step before failure

⇒ return = number of steps before failure

As a **continuing task** with discounted return:

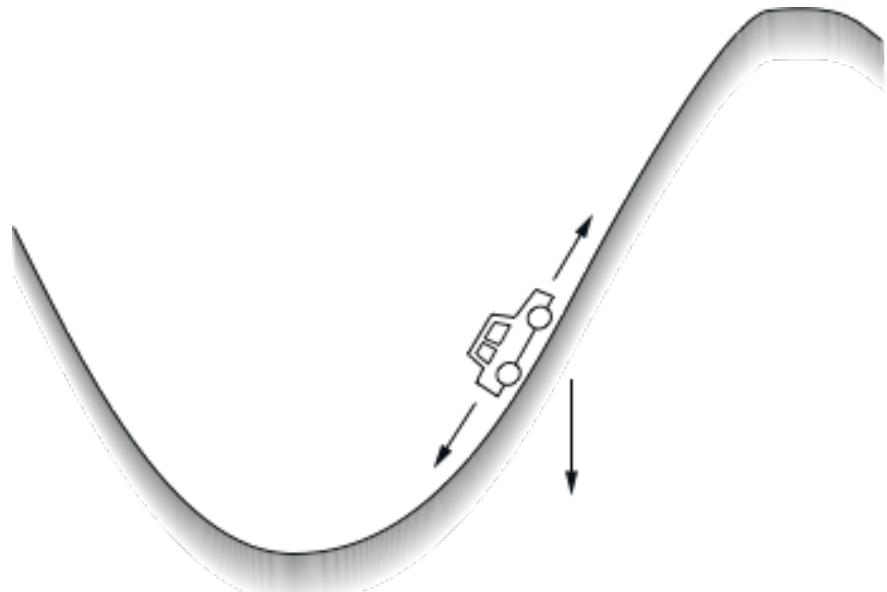
reward = -1 upon failure; 0 otherwise

⇒ return =  $-\gamma^k$ , for  $k$  steps before failure

In either case, return is maximized by avoiding failure for as long as possible.

# Another Example: Mountain Car

---



Get to the top of the hill  
as quickly as possible.

reward =  $-1$  for each step where **not** at top of hill

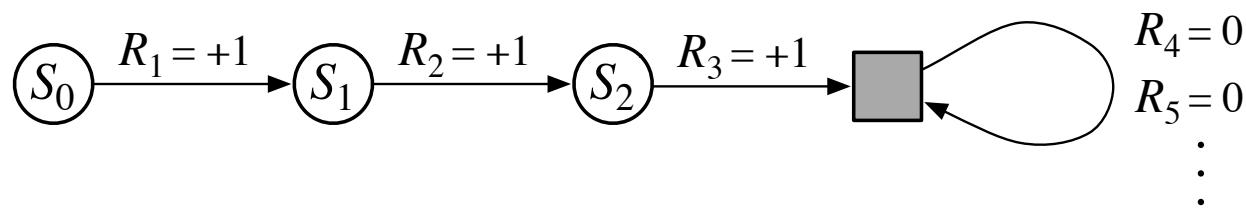
$\Rightarrow$  return = - number of steps before reaching top of hill

Return is maximized by minimizing  
number of steps to reach the top of the hill.

# A Trick to Unify Notation for Returns

---

- In episodic tasks, we number the time steps of each episode starting from zero.
- We usually do not have to distinguish between episodes, so instead of writing  $S_{t,j}$  for states in episode  $j$ , we write just  $S_t$
- Think of each episode as ending in an absorbing state that always produces reward of zero:



- We can cover all cases by writing  $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ ,
- where  $\gamma$  can be 1 only if a zero reward absorbing state is always reached.

# The Markov Property

---

- By “the state” at step  $t$ , the book means whatever information is available to the agent at step  $t$  about its environment.
- The state can include immediate “sensations,” highly processed sensations, and structures built up over time from sequences of sensations.
- Ideally, a state should summarize past sensations so as to retain all “essential” information, i.e., it should have the **Markov Property**:

$$\Pr\{R_{t+1} = r, S_{t+1} = s' \mid S_0, A_0, R_1, \dots, S_{t-1}, A_{t-1}, R_t, S_t, A_t\} = \\ p(s', r | s, a) = \Pr\{R_{t+1} = r, S_{t+1} = s' \mid S_t, A_t\}$$

for all  $s' \in \mathcal{S}^+, r \in \mathcal{R}$ , and all histories  $S_0, A_0, R_1, \dots, S_{t-1}, A_{t-1}, R_t, S_t, A_t$ .

# Markov Decision Processes

---

- If a reinforcement learning task has the Markov Property, it is basically a **Markov Decision Process (MDP)**.
- If state and action sets are finite, it is a **finite MDP**.
- To define a finite MDP, you need to give:
  - **state and action sets**
  - one-step “dynamics”

$$p(s', r | s, a) = \Pr\{S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a\}$$

$$p(s' | s, a) \doteq \Pr\{S_{t+1} = s' \mid S_t = s, A_t = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a)$$

$$r(s, a) \doteq \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a)$$

# An Example Finite MDP

---

## Recycling Robot

- At each step, robot has to decide whether it should (1) actively search for a can, (2) wait for someone to bring it a can, or (3) go to home base and recharge.
- Searching is better but runs down the battery; if runs out of power while searching, has to be rescued (which is bad).
- Decisions made on basis of current energy level: high, low.
- Reward = number of cans collected

# Recycling Robot MDP

$$\mathcal{S} = \{\text{high}, \text{low}\}$$

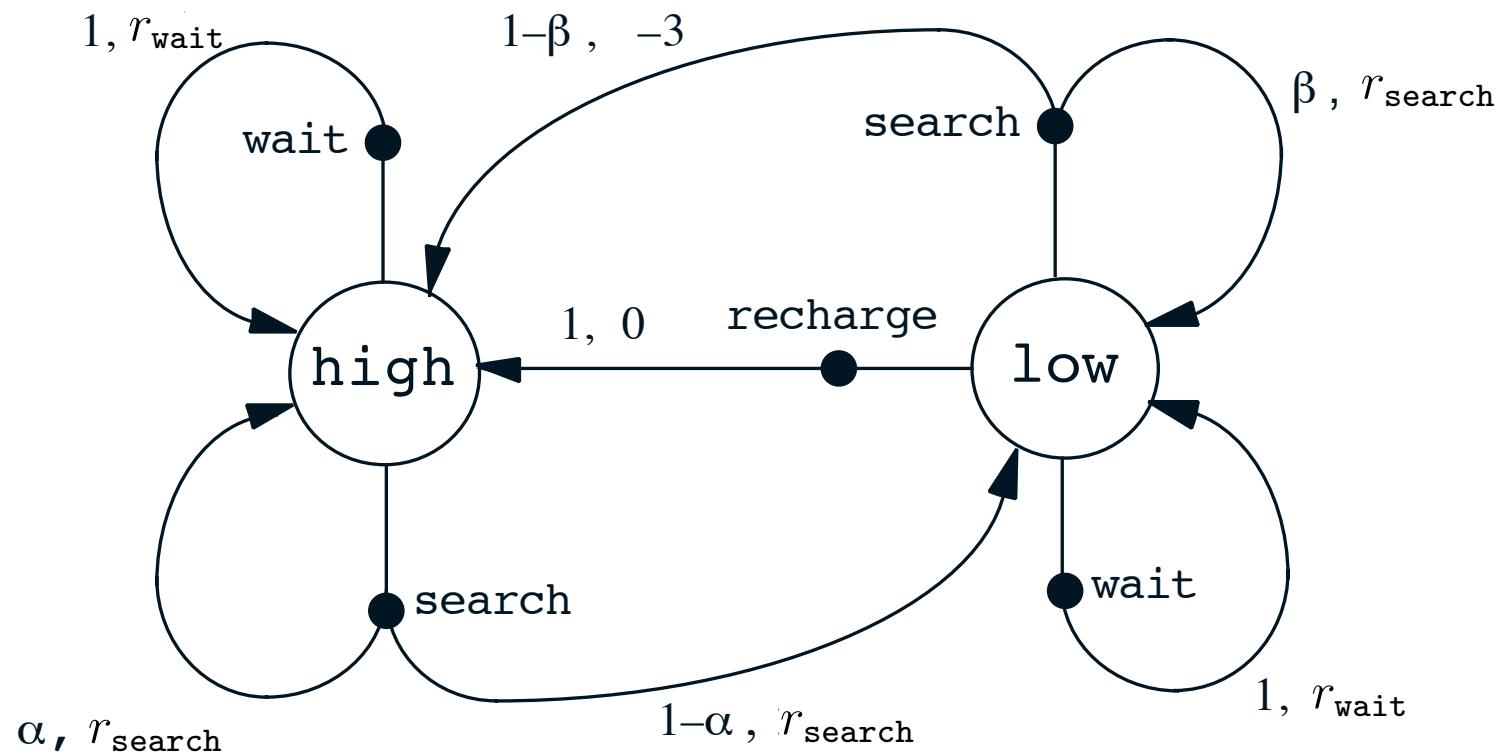
$$\mathcal{A}(\text{high}) = \{\text{search}, \text{wait}\}$$

$$\mathcal{A}(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}$$

$r_{\text{search}}$  = expected no. of cans while searching

$r_{\text{wait}}$  = expected no. of cans while waiting

$$r_{\text{search}} > r_{\text{wait}}$$



# The Agent Learns a Policy

---

**Policy at step  $t$**  =  $\pi_t$  =

a mapping from states to action probabilities

$\pi_t(a \mid s)$  = probability that  $A_t = a$  when  $S_t = s$

Special case - *deterministic policies*:

$\pi_t(s)$  = the action taken with prob=1 when  $S_t = s$

- Reinforcement learning methods specify how the agent changes its policy as a result of experience.
- Roughly, the agent's goal is to get as much reward as it can over the long run.

# Value Functions

---

- The **value of a state** is the expected return starting from that state; depends on the agent's policy:

**State - value function for policy  $\pi$  :**

$$v_{\pi}(s) = E_{\pi} \left\{ G_t \mid S_t = s \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right\}$$

- The **value of an action (in a state)** is the expected return starting after taking that action from that state; depends on the agent's policy:

**Action - value function for policy  $\pi$  :**

$$q_{\pi}(s, a) = E_{\pi} \left\{ G_t \mid S_t = s, A_t = a \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right\}$$

# Bellman Equation for a Policy $\pi$

---

The basic idea:

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

So:

$$\begin{aligned} v_\pi(s) &= E_\pi \left\{ G_t \mid S_t = s \right\} \\ &= E_\pi \left\{ R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s \right\} \end{aligned}$$

Or, without the expectation operator:

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

# More on the Bellman Equation

---

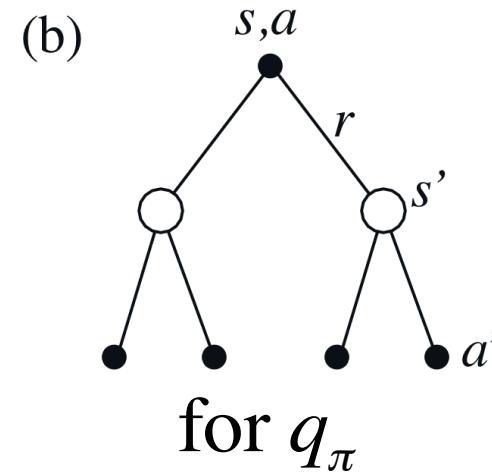
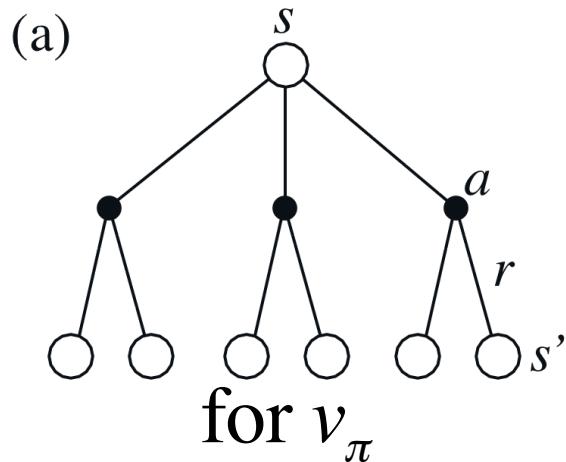
$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')]$$

This is a set of equations (in fact, linear), one for each state.  
The value function for  $\pi$  is its unique solution.

$$\begin{aligned}\text{Action Value: } q_\pi(s, a) &= \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t=s, A_t=a] \\ &= \sum_{s',r} p(s',r|s,a) [r + \gamma v_\pi(s')]\end{aligned}$$

$$\begin{aligned}v_\pi(s) &= \mathbb{E}_\pi[q_\pi(S_t, A_t) \mid S_t=s] \\ &= \sum_a \pi(a|s) q_\pi(s, a)\end{aligned}$$

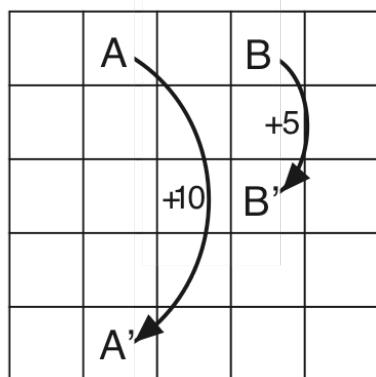
**Backup diagrams:**



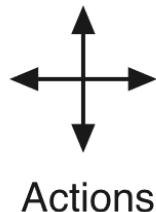
# Gridworld

---

- ❑ Actions: north, south, east, west; deterministic.
- ❑ If would take agent off the grid: no move but reward =  $-1$
- ❑ Other actions produce reward = 0, except actions that move agent out of special states A and B as shown.



(a)



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

(b)

State-value function  
for equiprobable  
random policy;  
 $\gamma = 0.9$

# Optimal Value Functions

---

- For finite MDPs, policies can be **partially ordered**:
$$\pi \geq \pi' \quad \text{if and only if } v_\pi(s) \geq v_{\pi'}(s) \text{ for all } s \in \mathcal{S}$$
- There are always one or more policies that are better than or equal to all the others. These are the **optimal policies**. We denote them all  $\pi_*$ .
- Optimal policies share the same **optimal state-value function**:
$$v_*(s) = \max_{\pi} v_{\pi}(s) \quad \text{for all } s \in \mathcal{S}$$
- Optimal policies also share the same **optimal action-value function**:
$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad \text{for all } s \in \mathcal{S} \text{ and } a \in \mathcal{A}$$

This is the expected return for taking action  $a$  in state  $s$  and thereafter following an optimal policy.

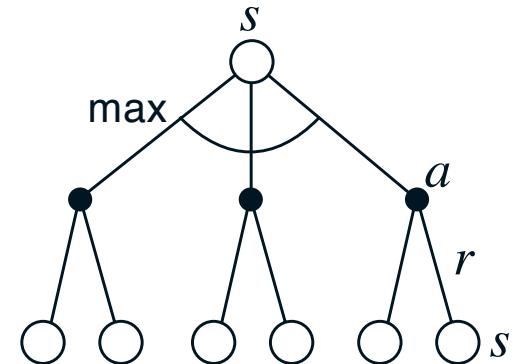
# Bellman Optimality Equation for $v_*$

---

The value of a state under an optimal policy must equal the expected return for the best action from that state:

$$\begin{aligned} v_*(s) &= \max_a q_{\pi_*}(s, a) \\ &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]. \end{aligned}$$

The relevant backup diagram:



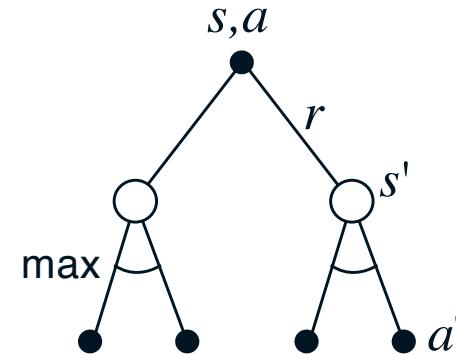
$v_*$  is the unique solution of this system of nonlinear equations.

# Bellman Optimality Equation for $q_*$

---

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \max_{a'} q_*(s', a') \right]. \end{aligned}$$

The relevant backup diagram:



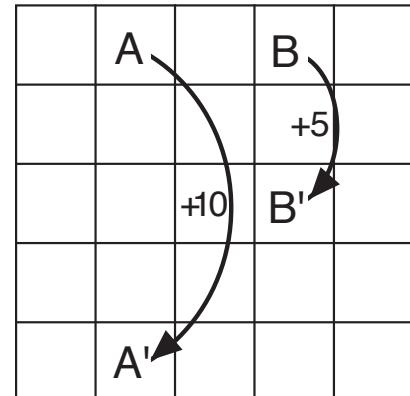
$q_*$  is the unique solution of this system of nonlinear equations.

# Why Optimal State-Value Functions are Useful

Any policy that is greedy with respect to  $v_*$  is an optimal policy.

Therefore, given  $v_*$ , one-step-ahead search produces the long-term optimal actions.

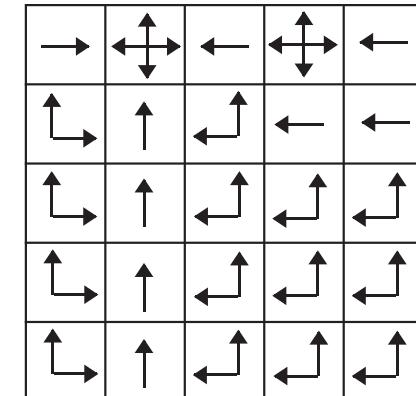
E.g., back to the gridworld:



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

b)  $v_*$



c)  $\pi_*$

# What About Optimal Action-Value Functions?

---

Given  $q_*$ , the agent does not even have to do a one-step-ahead search:

$$\pi_*(s) = \arg \max_a q_*(s, a)$$

# Solving the Bellman Optimality Equation

---

- Finding an optimal policy by solving the Bellman Optimality Equation requires the following:
  - accurate knowledge of environment dynamics;
  - we have enough space and time to do the computation;
  - the Markov Property.
- How much space and time do we need?
  - polynomial in number of states (via dynamic programming methods; Chapter 4),
  - BUT, number of states is often huge (e.g., backgammon has about  $10^{20}$  states).
- We usually have to settle for approximations.
- Many RL methods can be understood as approximately solving the Bellman Optimality Equation.

# Four value functions

	state values	action values
prediction	$v_\pi$	$q_\pi$
control	$v_*$	$q_*$

- All theoretical objects, mathematical ideals (expected values)
- Distinct from their estimates:

$$V_t(s) \quad Q_t(s, a)$$

# Values are *expected* returns

- The value of a state, given a policy:

$$v_\pi(s) = \mathbb{E}\{G_t \mid S_t = s, A_{t:\infty} \sim \pi\} \quad v_\pi : \mathcal{S} \rightarrow \mathbb{R}$$

- The value of a state-action pair, given a policy:

$$q_\pi(s, a) = \mathbb{E}\{G_t \mid S_t = s, A_t = a, A_{t+1:\infty} \sim \pi\} \quad q_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$$

- The optimal value of a state:

$$v_*(s) = \max_\pi v_\pi(s) \quad v_* : \mathcal{S} \rightarrow \mathbb{R}$$

- The optimal value of a state-action pair:

$$q_*(s, a) = \max_\pi q_\pi(s, a) \quad q_* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$$

- Optimal policy:  $\pi_*$  is an optimal policy if and only if

$$\pi_*(a|s) > 0 \text{ only where } q_*(s, a) = \max_b q_*(s, b) \quad \forall s \in \mathcal{S}$$

- in other words,  $\pi_*$  is optimal iff it is *greedy* wrt  $q_*$

# Summary

---

- ❑ Agent-environment interaction
  - States
  - Actions
  - Rewards
- ❑ Policy: stochastic rule for selecting actions
- ❑ Return: the function of future rewards agent tries to maximize
- ❑ Episodic and continuing tasks
- ❑ Markov Property
- ❑ Markov Decision Process
  - Transition probabilities
  - Expected rewards
- ❑ Value functions
  - State-value function for a policy
  - Action-value function for a policy
  - Optimal state-value function
  - Optimal action-value function
- ❑ Optimal policies
- ❑ Bellman Equations
- ❑ The need for approximation