# CSE 150A-250A AI: Probabilistic Models

## Lecture 15

Fall 2025

Trevor Bonjour
Department of Computer Science and Engineering
University of California, San Diego

Slides adapted from previous versions of the course (Prof. Lawrence, Prof. Alvarado, Prof Berg-Kirkpatrick)

Review

Value functions

Planning in MDPs

Policy Based

    Policy Evaluation

    Policy Improvement

    Policy Iteration

# Review

# Reinforcement learning (RL)

- Learning from experience in the world



- Formalization as Markov decision process

$\mathcal{S}$      state space    *finite*

$\mathcal{A}$      action space    "

$P(s'|s, a)$      transition probabilities

$R(s)$      reward function

MDP      $\{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s)\}$

- Definition

  $\pi(s)$ .

  A **policy** $\pi : \mathcal{S} \rightarrow \mathcal{A}$ is a mapping of states to actions.
  In this class we will only consider deterministic policies.

- **Definition**

  A **policy** $\pi : \mathcal{S} \rightarrow \mathcal{A}$ is a mapping of states to actions.
  In this class we will only consider deterministic policies.

- **Number of policies**

  If there are $|\mathcal{A}|$ possible actions in each of $|\mathcal{S}|$ states,
  then there are *combinatorially* many policies:

  $$\text{\# policies} \ = \ |\mathcal{A}|^{|\mathcal{S}|}$$

- Definition

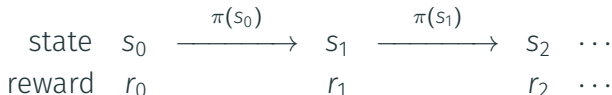  A **policy** $\pi : \mathcal{S} \to \mathcal{A}$ is a mapping of states to actions.
  In this class we will only consider deterministic policies.

- Number of policies

  If there are $|\mathcal{A}|$ possible actions in each of $|\mathcal{S}|$ states,
  then there are *combinatorially* many policies:

  $$\# \text{ policies } = |\mathcal{A}|^{|\mathcal{S}|}$$

- Experience under policy $\pi$

$$\begin{array}{cccccc}
\text{state} & s_0 & \xrightarrow{\pi(s_0)} & s_1 & \xrightarrow{\pi(s_1)} & s_2 & \cdots \\
\text{reward} & r_0 & & r_1 & & r_2 & \cdots
\end{array}$$

Transitions occur with probabilities $P(s'|s, \pi(s))$.

_deterministic_

A policy $\pi$ completely determines the next state $s'$ that an agent will end up in after taking an action from state $s$.

True (A) or False (B)?

# How to measure long-term return?

# How to measure long-term return?

1. Finite-horizon return

1. **Finite-horizon return**

$$\text{return} = \frac{1}{T}(r_0 + r_1 + \cdots + r_{T-1}) \quad \text{for a } T\text{-step horizon}$$

## How to measure long-term return?

1. Finite-horizon return

$$\text{return} = \frac{1}{T}(r_0 + r_1 + \cdots + r_{T-1}) \quad \text{for a } T\text{-step horizon}$$

2. Undiscounted return with infinite horizon

## How to measure long-term return?

1. **Finite-horizon return**

$$\text{return} = \frac{1}{T}(r_0 + r_1 + \cdots + r_{T-1}) \quad \text{for a } T\text{-step horizon}$$

2. **Undiscounted return with infinite horizon**

$$\text{return} = \lim_{T \to \infty} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_t \right]$$

1. **Finite-horizon return**

$$\text{return} \ = \ \frac{1}{T}(r_0 + r_1 + \cdots + r_{T-1}) \quad \text{for a } T\text{-step horizon}$$

2. **Undiscounted return with infinite horizon**

$$\text{return} \ = \ \lim_{T \to \infty} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_t \right]$$

These are the most obvious ways to accumulate rewards.

1. **Finite-horizon return**

$$\text{return} = \frac{1}{T}(r_0 + r_1 + \cdots + r_{T-1}) \quad \text{for a } T\text{-step horizon}$$

2. **Undiscounted return with infinite horizon**

$$\text{return} = \lim_{T \to \infty} \left[ \frac{1}{T} \sum_{t=0}^{T-1} r_t \right]$$

These are the most obvious ways to accumulate rewards.
But they are **not** the most commonly used in practice …

3. Discounted return with infinite horizon

3. Discounted return with infinite horizon

   Let $\gamma \in [0, 1)$ denote the so-called **discount factor**.

3. **Discounted return with infinite horizon**

Let $\gamma \in [0, 1)$ denote the so-called **discount factor**.
Then define

$$\text{return} = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \cdots = \sum_{t=0}^{\infty} \gamma^t r_t$$

3. **Discounted return with infinite horizon**

Let $\gamma \in [0, 1)$ denote the so-called **discount factor**.
Then define

$$\text{return} = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \cdots = \sum_{t=0}^{\infty} \gamma^t r_t$$

What does it mean when the discount factor $\gamma << 1$?

- A. Immediate and future rewards are valued equally.

- B. Future rewards are heavily discounted compared to immediate.

- C. Future rewards are lightly discounted compared to immediate.

- D. Only future rewards are considered.

3. **Discounted return with infinite horizon**

Let $\gamma \in [0, 1)$ denote the so-called **discount factor**.
Then define

$$\text{return} = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \cdots = \sum_{t=0}^{\infty} \gamma^t r_t$$

What does it mean when the discount factor $\gamma \sim 1$?

A. Immediate and future rewards are valued equally.

B. Future rewards are heavily discounted compared to immediate.

C. Future rewards are lightly discounted compared to immediate.

D. Only future rewards are considered.

3. **Discounted return with infinite horizon** Let $\gamma \in [0, 1)$
   denote the so-called **discount factor**.
   Then define

   $$\text{return} = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 + \cdots = \sum_{t=0}^{\infty} \gamma^t r_t$$

   When $\gamma \ll 1$, future rewards are heavily discounted.
   These returns can be optimized by short-sighted agents.

   When $\gamma$ is close to 1, future rewards are lightly discounted.
   These returns can only be optimized by far-sighted agents.

Psychologist:    *Why discount rewards from the distant future?*

Psychologist: *Why discount rewards from the distant future?*

Economist: *Why favor investments with short-term payoffs?*

Psychologist: *Why discount rewards from the distant future?*
   Economist: *Why favor investments with short-term payoffs?*

1. Intuition

Psychologist: *Why discount rewards from the distant future?*
  Economist: *Why favor investments with short-term payoffs?*

1. Intuition

   Many models are only approximations to the real world;

Psychologist: *Why discount rewards from the distant future?*
Economist: *Why favor investments with short-term payoffs?*

1. Intuition

   Many models are only approximations to the real world;
   we should not attempt to extrapolate them indefinitely.

Psychologist: *Why discount rewards from the distant future?*
  Economist: *Why favor investments with short-term payoffs?*

1. Intuition

   Many models are only approximations to the real world;
   we should not attempt to extrapolate them indefinitely.

2. Mathematical convenience

Psychologist: *Why discount rewards from the distant future?*
Economist: *Why favor investments with short-term payoffs?*

1. Intuition

   Many models are only approximations to the real world;
   we should not attempt to extrapolate them indefinitely.

2. Mathematical convenience

   Discounted returns lead to simple iterative algorithms
   with strong guarantees of convergence.

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
But we can try to optimize its expected value:

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
But we can try to optimize its expected value:

$$\mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

*the expected value of the*
*discounted infinite-horizon return,*

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
But we can try to optimize its expected value:

$$\mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty}\gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

*the expected value of the*
*discounted infinite-horizon return,*
*starting in state s at time $t=0$,*

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
But we can try to optimize its expected value:

$$\mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

*the expected value of the*
*discounted infinite-horizon return,*
*starting in state s at time $t=0$,*
*and following policy $\pi$.*

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
But we can try to optimize its expected value:

$$\mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

*the expected value of the discounted infinite-horizon return, starting in state s at time $t = 0$, and following policy $\pi$.*

Maximizing the expected return is:

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
**But we can try to optimize its expected value:**

$$\mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

*the expected value of the*
*discounted infinite-horizon return,*
*starting in state s at time $t = 0$,*
*and following policy $\pi$.*

**Maximizing the expected return is:**

    – generally wiser than maximizing the best-case return,

The discounted return $\sum_{t=0}^{\infty} \gamma^t r_t$ is a random variable.
**But we can try to optimize its expected value:**

$$\mathrm{E}^\pi\left[\sum_{t=0}^{\infty}\gamma^t R(s_t) \,\bigg|\, s_0 = s\right]$$

*the expected value of the discounted infinite-horizon return, starting in state s at time $t = 0$, and following policy $\pi$.*

**Maximizing the expected return is:**

– generally wiser than maximizing the best-case return,
– but not as robust as minimizing the worst-case return.

# Value functions

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^{\infty}\gamma^t R(s_t)\;\middle|\; s_0 = s\right]$$

## State value function

$$V^{\pi}(s) = \mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

expected return,
starting in state $s$,
following policy $\pi$

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^{\infty} \gamma^t R(s_t)\,\middle|\, s_0 = s\right]$$

expected return,
starting in state $s$,
following policy $\pi$

- Values versus rewards:

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

expected return,
starting in state $s$,
following policy $\pi$

- **Values versus rewards:**

  The reward $R(s)$ give **immediate** feedback to the agent.

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^\infty \gamma^t R(s_t)\,\middle|\, s_0 = s\right]$$

expected return,
starting in state $s$,
following policy $\pi$

- **Values versus rewards:**

  The reward $R(s)$ give **immediate** feedback to the agent.
  The value $V^\pi(s)$ computes the expected **long-term** return.

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$
expected return,
starting in state $s$,
following policy $\pi$

- **Values versus rewards:**

  The reward $R(s)$ give **immediate** feedback to the agent.
  The value $V^\pi(s)$ computes the expected **long-term** return.

- **Types of behaviors:**

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\Big|\, s_0 = s\right]$$

expected return,
starting in state $s$,
following policy $\pi$

- Values versus rewards:

  The reward $R(s)$ give **immediate** feedback to the agent.
  The value $V^\pi(s)$ computes the expected **long-term** return.

- Types of behaviors:

  Sacrifice now for long-term gain: $R(s) < 0$, $V^\pi(s) > 0$.

$$V^\pi(s) \;=\; \mathrm{E}^\pi\!\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

expected return,
starting in state $s$,
following policy $\pi$

- **Values versus rewards:**

  The reward $R(s)$ give **immediate** feedback to the agent.
  The value $V^\pi(s)$ computes the expected **long-term** return.

- **Types of behaviors:**

  Sacrifice now for long-term gain: $R(s) < 0$, $V^\pi(s) > 0$.
  Win now at the expense of later: $R(s) > 0$, $V^\pi(s) < 0$.

- Experience under policy $\pi$

- Experience under policy $\pi$

$$\text{state} \quad s_0 \quad \xrightarrow{\pi(s_0)} \quad s_1 \quad \xrightarrow{\pi(s_1)} \quad s_2 \quad \cdots$$

- Experience under policy $\pi$

$$\text{state} \quad s_0 \xrightarrow{\pi(s_0)} s_1 \xrightarrow{\pi(s_1)} s_2 \cdots$$
$$\text{reward} \quad r_0 \qquad\qquad r_1 \qquad\qquad r_2 \cdots$$

- Experience under policy $\pi$

$$\text{state} \quad s_0 \xrightarrow{\pi(s_0)} s_1 \xrightarrow{\pi(s_1)} s_2 \cdots$$
$$\text{reward} \quad r_0 \qquad\qquad r_1 \qquad\qquad r_2 \cdots$$

- Adjacent states

- Experience under policy $\pi$

$$\text{state} \quad s_0 \xrightarrow{\pi(s_0)} s_1 \xrightarrow{\pi(s_1)} s_2 \cdots$$
$$\text{reward} \quad r_0 \qquad\qquad r_1 \qquad\qquad r_2 \cdots$$

- Adjacent states

  States $(s, s')$ can be visited in succession if
  $P(s'|s, \pi(s)) > 0$.

- Experience under policy $\pi$

$$\begin{array}{cccccccc}
\text{state} & s_0 & \xrightarrow{\pi(s_0)} & s_1 & \xrightarrow{\pi(s_1)} & s_2 & \cdots \\
\text{reward} & r_0 & & r_1 & & r_2 & \cdots
\end{array}$$

- Adjacent states

States $(s, s')$ can be visited in succession if
$P(s'|s, \pi(s)) > 0$.

The values $V^\pi(s)$ and $V^\pi(s')$ should be related, but how?

- Experience under policy $\pi$

$$\text{state} \quad s_0 \xrightarrow{\pi(s_0)} s_1 \xrightarrow{\pi(s_1)} s_2 \cdots$$
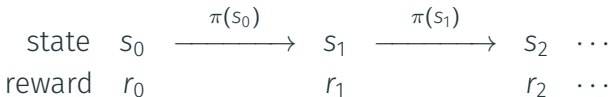$$\text{reward} \quad r_0 \qquad\qquad r_1 \qquad\qquad r_2 \cdots$$

- Adjacent states

States $(s, s')$ can be visited in succession if $P(s'|s, \pi(s)) > 0$.

The values $V^\pi(s)$ and $V^\pi(s')$ should be related, but how?

The **Bellman equation** tells us how.

$$V^\pi(s) \quad =$$

$$V^\pi(s) = \mathrm{E}^\pi\left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \,\middle|\, s_0 = s\right]$$

$$=$$

$$V^\pi(s) = \mathrm{E}^\pi\left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \Big| s_0 = s\right]$$

$$= R(s) + \gamma\,\mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \Big| s_0 = s\right]$$

$$=$$

$$V^\pi(s) = \mathrm{E}^\pi\left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \middle| s_0 = s\right]$$

$$= R(s) + \gamma\,\mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \middle| s_0 = s\right]$$

$$= R(s) + \gamma \sum_{s'} P(s'|s, \pi(s))\,\mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \middle| s_1 = s'\right]$$

$$=$$

*markov prop*

*+*

*Law of*
*total expectation.*

*$s_1 = s$*

*$s_1 = s'$, $s_0 = s$*

*Value at*
*$s'$*

# Bellman equation

$$
\begin{aligned}
V^\pi(s) &= \mathrm{E}^\pi\left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \,\middle|\, s_0 = s\right] \\[2ex]
&= R(s) + \gamma \, \mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \,\middle|\, s_0 = s\right] \\[2ex]
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) \, \mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \,\middle|\, s_1 = s'\right] \\[2ex]
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) \, V^\pi(s')
\end{aligned}
$$

$$V^\pi(s) = \mathrm{E}^\pi\left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \,\Big|\, s_0 = s\right]$$

$$= R(s) + \gamma\,\mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \,\Big|\, s_0 = s\right]$$

$$= R(s) + \gamma \sum_{s'} P(s'|s, \pi(s))\,\mathrm{E}^\pi\left[R(s_1) + \gamma R(s_2) + \cdots \,\Big|\, s_1 = s'\right]$$

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

*(handwritten annotations:)*

$V^\pi_{\in}(s')$

$= R(s')$

$+ \gamma \sum_{s'} P(s'|s'\pi(s))$

$V(s'')$

The Bellman equation is the basis for much that will follow:

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

*(handwritten annotations:)* $s$    $s' \to$ all adjacent states.

$Q^\pi(s, a) =$

$$Q^\pi(s, a) = \mathrm{E}^\pi \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \middle| s_0 = s, a_0 = a \right]$$

$$Q^\pi(s, a) = \mathrm{E}^\pi \left[ \sum_{t=0}^\infty \gamma^t R(s_t) \middle| s_0 = s, a_0 = a \right]$$

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

$$Q^\pi(s, a) = \mathrm{E}^\pi\left[\sum_{t=0}^{\infty} \gamma^t R(s_t)\,\middle|\, s_0 = s, a_0 = a\right]$$

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

- Motivation

$$Q^\pi(s, a) = \mathrm{E}^\pi\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s, a_0 = a\right]$$

expected return,
starting from state $s$,
**taking action $a$,**
then following policy $\pi$

- Motivation

  *not ness true*

  Useful to imagine how small changes affect expected outcomes.

$\pi(s) \rightarrow a$

$\rightarrow a \neq \pi(s)$

$after \downarrow$

$\pi(s)$

$$Q^\pi(s, a) \;=\; \mathrm{E}^\pi\left[\sum_{t=0}^{\infty} \gamma^t R(s_t)\,\middle|\, s_0 = s, a_0 = a\right]$$

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

- Motivation

  Useful to imagine how small changes affect expected outcomes.
  What if (just once) the agent acted differently in state $s$?

$$Q^\pi(s, a) = \mathrm{E}^\pi\left[\sum_{t=0}^{\infty} \gamma^t R(s_t)\,\middle|\, s_0 = s, a_0 = a\right]$$

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

· Motivation

Useful to imagine how small changes affect expected outcomes.
What if (just once) the agent acted differently in state $s$?

· Analogous to the Bellman equation:

$$Q^\pi(s, a) \;=\; \mathrm{E}^\pi\left[\sum_{t=0}^\infty \gamma^t R(s_t)\,\middle|\, s_0 = s, a_0 = a\right]$$

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

- Motivation

  Useful to imagine how small changes affect expected outcomes.
  What if (just once) the agent acted differently in state $s$?

- Analogous to the Bellman equation:

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

$$Q^\pi(s, a) \;=\; \mathrm{E}^\pi\left[\sum_{t=0}^\infty \gamma^t R(s_t) \,\middle|\, s_0 = s, a_0 = a\right]$$

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

· Motivation

  Useful to imagine how small changes affect expected outcomes.
  What if (just once) the agent acted differently in state $s$?

· Analogous to the Bellman equation:

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

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

- Goal

- Goal

  Find the optimal policy given the environment that the agent is in.

- Goal

  Find the optimal policy given the environment that the agent is in.

- Planning

- **Goal**

  Find the optimal policy given the environment that the agent is in.

- **Planning**

  If reward function and transition probabilities are known.

- Goal

  Find the optimal policy given the environment that the agent is in.

- Planning

  If reward function and transition probabilities are known.

- Reinforcement Learning

# Optimality

- **Goal**

  Find the optimal policy given the environment that the agent is in.

- **Planning**

  If reward function and transition probabilities are known.

- **Reinforcement Learning**
  If reward function and transition probabilities are unknown.

- Theorem

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

$$V^*(s) \quad =$$

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

$$V^*(s) \quad = \quad V^{\pi^*}(s)$$

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

$$V^*(s) = V^{\pi^*}(s)$$
$$Q^*(s, a) =$$

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

$$
\begin{aligned}
V^*(s) &= V^{\pi^*}(s) \\
Q^*(s, a) &= Q^{\pi^*}(s, a)
\end{aligned}
$$

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

$$
\begin{aligned}
V^*(s) &= V^{\pi^*}(s) \\
Q^*(s, a) &= Q^{\pi^*}(s, a)
\end{aligned}
$$

These optimal value functions are **unique**.

- Theorem

  There exists at least one policy $\pi^*$ (and perhaps many) such that $V^{\pi^*}(s) \geq V^\pi(s)$ for all policies $\pi$ and states $s$ of the MDP.

- Notation

$$
\begin{aligned}
V^*(s) &= V^{\pi^*}(s) \\
Q^*(s, a) &= Q^{\pi^*}(s, a)
\end{aligned}
$$

  These optimal value functions are **unique**.
  (All optimal policies share the same value functions.)

- From the optimal action value function:

- From the optimal action value function:

$$V^*(s) \quad =$$

- From the optimal action value function:

$$V^*(s) \quad = \quad \max_a \left[ Q^*(s, a) \right]$$

- From the optimal action value function:

$$V^*(s) \quad = \quad \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) \quad =$$

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s, a) \right] \rightarrow \textit{value}$$

$$\pi^*(s) = \operatorname*{argmax}_a \left[ Q^*(s, a) \right] \quad \textit{Single action}$$

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) = \arg\max_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

- From the optimal action value function:

$$V^*(s) \quad = \quad \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) \quad = \quad \arg\max_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

$$Q^*(s, a) \quad =$$

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) = \arg\max_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

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

- From the optimal action value function:

$$V^*(s) \quad = \quad \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) \quad = \quad \operatorname{argmax}_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

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

$$\pi^*(s) \quad =$$

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s,a) \right]$$

$$\pi^*(s) = \arg\max_a \left[ Q^*(s,a) \right]$$

- From the optimal state value function:

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

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

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) = \operatorname*{argmax}_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

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

$$\pi^*(s) = \operatorname*{argmax}_a \left[ R(s) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right]$$

- Why are these relations useful?

- From the optimal action value function:

$$V^*(s) \;=\; \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) \;=\; \operatorname*{argmax}_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

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

$$\pi^*(s) \;=\; \operatorname*{argmax}_a \left[ R(s) \;+\; \gamma \sum_{s'} P(s'|s, a) V^*(s') \right]$$

- Why are these relations useful?

  Sometimes it can be easier to estimate $Q^*(s, a)$ or $V^*(s)$

# Relations at optimality

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) = \arg\max_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

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

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

- Why are these relations useful?

  Sometimes it can be easier to estimate $Q^*(s, a)$ or $V^*(s)$ (which are continuous)

- From the optimal action value function:

$$V^*(s) = \max_a \left[ Q^*(s, a) \right]$$

$$\pi^*(s) = \arg\max_a \left[ Q^*(s, a) \right]$$

- From the optimal state value function:

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

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

- Why are these relations useful?

  Sometimes it can be easier to estimate $Q^*(s, a)$ or $V^*(s)$ (which are continuous) than to learn $\pi^*(s)$ (which is discrete).

# Planning in MDPs

Given a complete model of the agent and its environment as a Markov decision process,

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute (i.e., in time *polynomial in the number of states*)

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} \; = \; \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute (i.e., in time *polynomial in the
number of states*) any of the following:

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute (i.e., in time *polynomial in the number of states*) any of the following:

1. an optimal policy $\pi^*(s)$?

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute (i.e., in time *polynomial in the number of states*) any of the following:

1. an optimal policy $\pi^*(s)$?
2. the optimal state value function $V^*(s)$?

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$MDP = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute (i.e., in time *polynomial in the number of states*) any of the following:

1. an optimal policy $\pi^*(s)$?
2. the optimal state value function $V^*(s)$?
3. the optimal action value function $Q^*(s, a)$?

Given a complete model of the agent and its environment
as a Markov decision process, namely

$$\text{MDP} = \{\mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma\},$$

how can we *efficiently* compute (i.e., in time *polynomial in the number of states*) any of the following:

1. an optimal policy $\pi^*(s)$?
2. the optimal state value function $V^*(s)$?
3. the optimal action value function $Q^*(s, a)$?

This is the problem of **planning** in MDPs.

# Policy Based

1. Policy evaluation

   How to compute $V^\pi(s)$ for some fixed policy $\pi$?

1. **Policy evaluation**

   How to compute $V^\pi(s)$ for some fixed policy $\pi$?

2. **Policy improvement**

   How to compute a policy $\pi'$ such that $V^{\pi'}(s) \geq V^\pi(s)$?

1. **Policy evaluation**

   How to compute $V^\pi(s)$ for some fixed policy $\pi$?

2. **Policy improvement**

   How to compute a policy $\pi'$ such that $V^{\pi'}(s) \geq V^\pi(s)$?

3. **Policy iteration**

   How to compute an optimal policy $\pi^*(s)$?

- How to compute the state value function?

- How to compute the state value function?

$$V^{\pi}(s) \;=\; \mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

- How to compute the state value function?

$$V^\pi(s) = \mathrm{E}^\pi\left[\sum_{t=0}^{\infty}\gamma^t R(s_t)\,\middle|\,s_0 = s\right]$$

- Bellman equation:

- How to compute the state value function?

$$V^{\pi}(s) = \mathrm{E}^{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

- Bellman equation:

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

- How to compute the state value function?

$$V^\pi(s) = \mathrm{E}^\pi\left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \,\middle|\, s_0 = s\right]$$

- Bellman equation:

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

- Solve linear system:

- How to compute the state value function?

$$V^\pi(s) = \mathrm{E}^\pi \left[ \sum_{t=0}^\infty \gamma^t R(s_t) \,\middle|\, s_0 = s \right]$$

- Bellman equation:

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

- **Solve linear system:** There are $n$ equations for $n$ unknowns (where $s = 1, 2, \ldots, n$).

- From the Bellman equation:

- From the Bellman equation:

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

- From the Bellman equation:

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

- Rearranging terms:

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

$$= \sum_{s'} \Bigg[$$

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

$$= \sum_{s'} \left[ \phantom{xx} I(s, s') \right.$$

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

$$= \sum_{s'} \Bigg[ \underbrace{I(s,s')}_{\text{identity matrix}}$$

# Solving the linear system

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^{\pi}(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^{\pi}(s')$$

$$= \sum_{s'} \left[ \underbrace{I(s, s')}_{\text{identity matrix}} - \gamma P(s'|s, \pi(s)) \right]$$

- From the Bellman equation:

$$V^\pi(s) \;=\; R(s) \,+\, \gamma \sum_{s'} P(s'|s, \pi(s))\, V^\pi(s').$$

- Rearranging terms:

$$R(s) \;=\; V^\pi(s) \,-\, \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

$$= \; \sum_{s'} \left[ \underbrace{I(s, s')}_{\text{identity matrix}} \,-\, \gamma P(s'|s, \pi(s)) \right] V^\pi(s')$$

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^{\pi}(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^{\pi}(s')$$

$$= \sum_{s'} \left[ \underbrace{I(s, s')}_{\text{identity matrix}} - \gamma P(s'|s, \pi(s)) \right] V^{\pi}(s')$$

- In matrix-vector form:

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s))V^\pi(s')$$

$$= \sum_{s'} \left[ \underbrace{I(s, s')}_{\text{identity matrix}} - \gamma P(s'|s, \pi(s)) \right] V^\pi(s')$$

- In matrix-vector form:

$$R = \left[ I - \gamma P^\pi \right] V^\pi$$

- From the Bellman equation:

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

- Rearranging terms:

$$R(s) = V^\pi(s) - \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

$$= \sum_{s'} \left[ \underbrace{I(s, s')}_{\text{identity matrix}} - \gamma P(s'|s, \pi(s)) \right] V^\pi(s')$$

- In matrix-vector form:

$$R = \left[ I - \gamma P^\pi \right] V^\pi$$

$$\begin{bmatrix} \text{column vector of} \\ n \text{ known rewards} \end{bmatrix} = \begin{bmatrix} n \times n \text{ matrix} \\ \text{(known)} \end{bmatrix} \begin{bmatrix} \text{column vector of} \\ n \text{ unknown values} \end{bmatrix}$$

- Solution

- Solution
$$R = \left[I - \gamma P^\pi\right] V^\pi \implies$$

- Solution

$$R = \left[ I - \gamma P^\pi \right] V^\pi \quad \implies \quad V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Solution
$$R = \left[ I - \gamma P^\pi \right] V^\pi \implies V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Complexity

- Solution
$$R = \left[I - \gamma P^\pi\right] V^\pi \quad \implies \quad V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Complexity

  It takes $O(n^3)$ operations to solve this system of equations.

- Solution

$$R = \left[ I - \gamma P^\pi \right] V^\pi \quad \implies \quad V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Complexity

  It takes $O(n^3)$ operations to solve this system of equations.

- Example

- **Solution**

$$R = \left[I - \gamma P^\pi\right] V^\pi \quad \Longrightarrow \quad V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- **Complexity**

  It takes $O(n^3)$ operations to solve this system of equations.

- **Example**

  Let $\mathcal{S} = \{1, 2\}$ and $P(s'|s, \pi(s)) = 0.5$ for all $(s, s')$.

- Solution

$$R = \left[I - \gamma P^\pi\right] V^\pi \implies V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Complexity

  It takes $O(n^3)$ operations to solve this system of equations.

- Example

  Let $\mathcal{S} = \{1, 2\}$ and $P(s'|s, \pi(s)) = 0.5$ for all $(s, s')$.

$$\left[\begin{array}{c} V^\pi(1) \\ V^\pi(2) \end{array}\right] =$$

- Solution
$$R = \left[I - \gamma P^\pi\right] V^\pi \implies V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Complexity

  It takes $O(n^3)$ operations to solve this system of equations.

- Example

  Let $\mathcal{S} = \{1, 2\}$ and $P(s'|s, \pi(s)) = 0.5$ for all $(s, s')$.

$$\left[\begin{array}{c} V^\pi(1) \\ V^\pi(2) \end{array}\right] = \left(\left[\begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array}\right]\right.$$

- **Solution**
  $$R = \left[I - \gamma P^\pi\right] V^\pi \implies V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- **Complexity**

  It takes $O(n^3)$ operations to solve this system of equations.

- **Example**

  Let $\mathcal{S} = \{1, 2\}$ and $P(s'|s, \pi(s)) = 0.5$ for all $(s, s')$.

  $$\left[\begin{array}{c} V^\pi(1) \\ V^\pi(2) \end{array}\right] = \left(\left[\begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array}\right] - \gamma \left[\begin{array}{cc} 0.5 & 0.5 \\ 0.5 & 0.5 \end{array}\right]\right)^{-1}$$

- Solution

$$R = \left[ I - \gamma P^\pi \right] V^\pi \implies V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\text{matrix inverse}} R$$

- Complexity

  It takes $O(n^3)$ operations to solve this system of equations.

- Example

  Let $\mathcal{S} = \{1, 2\}$ and $P(s'|s, \pi(s)) = 0.5$ for all $(s, s')$.

$$\begin{bmatrix} V^\pi(1) \\ V^\pi(2) \end{bmatrix} = \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \gamma \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} \right)^{-1} \begin{bmatrix} R(1) \\ R(2) \end{bmatrix}.$$

# Policy improvement

- Problem statement

- Problem statement

  Given a policy $\pi$ and its state value function $V^\pi(s)$,

## Policy improvement

- Problem statement

  Given a policy $\pi$ and its state value function $V^\pi(s)$,
  how to compute a policy $\pi'$ such that

## Policy improvement

- Problem statement

  Given a policy $\pi$ and its state value function $V^\pi(s)$,
  how to compute a policy $\pi'$ such that

  $$V^{\pi'}(s) \geq V^\pi(s) \quad \text{for all states } s?$$

- Problem statement

  Given a policy $\pi$ and its state value function $V^\pi(s)$,
  how to compute a policy $\pi'$ such that

  $$V^{\pi'}(s) \geq V^\pi(s) \quad \text{for all states } s?$$

- Definition

- Problem statement

  Given a policy $\pi$ and its state value function $V^{\pi}(s)$,
  how to compute a policy $\pi'$ such that

  $$V^{\pi'}(s) \geq V^{\pi}(s) \quad \text{for all states } s?$$

- Definition

  Given the action value function $Q^{\pi}(s, a)$ for policy $\pi$, we
  define the **greedy policy** $\pi'$ by

- Problem statement

  Given a policy $\pi$ and its state value function $V^\pi(s)$, how to compute a policy $\pi'$ such that

  $$V^{\pi'}(s) \geq V^\pi(s) \quad \text{for all states } s?$$

- Definition

  Given the action value function $Q^\pi(s, a)$ for policy $\pi$, we define the **greedy policy** $\pi'$ by

  $$\pi'(s) = \operatorname*{argmax}_a \left[ Q^\pi(s, a) \right].$$

- Problem statement

  Given a policy $\pi$ and its state value function $V^\pi(s)$, how to compute a policy $\pi'$ such that

  $$V^{\pi'}(s) \geq V^\pi(s) \quad \text{for all states } s?$$

- Definition

  Given the action value function $Q^\pi(s, a)$ for policy $\pi$, we define the **greedy policy** $\pi'$ by

  $$\pi'(s) = \underset{a}{\arg\max} \left[ Q^\pi(s, a) \right].$$

  Why *greedy*?

- **Problem statement**

  Given a policy $\pi$ and its state value function $V^\pi(s)$,
  how to compute a policy $\pi'$ such that

  $$V^{\pi'}(s) \geq V^\pi(s) \quad \text{for all states } s?$$

- **Definition**

  Given the action value function $Q^\pi(s, a)$ for policy $\pi$, we
  define the **greedy policy** $\pi'$ by

  $$\pi'(s) = \underset{a}{\arg\max} \left[ Q^\pi(s, a) \right].$$

  Why *greedy*? Because we change the action in state $s$ to
  whatever appears to improve the expected return.

- In terms of the state value function:

- In terms of the state value function:

$$\pi'(s) = \underset{a}{\arg\max} \left[ Q^\pi(s, a) \right]$$

- In terms of the state value function:

$$\pi'(s) = \operatorname*{argmax}_a \left[ Q^\pi(s, a) \right]$$

$$= \operatorname*{argmax}_a \left[ R(s) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s') \right]$$

*following*

*$\pi$ after action*

*$a$*

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \operatorname*{argmax}_a \left[ Q^\pi(s, a) \right] \\
&= \operatorname*{argmax}_a \left[ R(s) + \gamma \sum_{s'} P(s'|s, a)\, V^\pi(s') \right] \\
&= \operatorname*{argmax}_a \left[ \sum_{s'} P(s'|s, a)\, V^\pi(s') \right]
\end{aligned}
$$

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \underset{a}{\text{argmax}} \left[ Q^\pi(s, a) \right] \\
&= \underset{a}{\text{argmax}} \left[ R(s) + \gamma \sum_{s'} P(s'|s, a) \, V^\pi(s') \right] \\
&= \underset{a}{\text{argmax}} \left[ \sum_{s'} P(s'|s, a) \, V^\pi(s') \right]
\end{aligned}
$$

- Test your understanding:

- In terms of the state value function:

$$\pi'(s) = \underset{a}{\text{argmax}} \left[ Q^\pi(s, a) \right]$$

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

$$= \underset{a}{\text{argmax}} \left[ \sum_{s'} P(s'|s, a) V^\pi(s') \right]$$

- Test your understanding:

  $\pi'(s) = \pi(s)$ for some $s \in \mathcal{S}$?

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \operatorname*{argmax}_{a}\left[Q^{\pi}(s, a)\right] \\
&= \operatorname*{argmax}_{a}\left[R(s) + \gamma\sum_{s'}P(s'|s, a)\, V^{\pi}(s')\right] \\
&= \operatorname*{argmax}_{a}\left[\sum_{s'}P(s'|s, a)\, V^{\pi}(s')\right]
\end{aligned}
$$

- Test your understanding:

    $\pi'(s) = \pi(s)$ for some $s \in \mathcal{S}$?        not necessarily

## Greedy policies

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \underset{a}{\arg\max}\left[Q^\pi(s, a)\right] \\
&= \underset{a}{\arg\max}\left[R(s) + \gamma\sum_{s'}P(s'|s, a)\,V^\pi(s')\right] \\
&= \underset{a}{\arg\max}\left[\sum_{s'}P(s'|s, a)\,V^\pi(s')\right]
\end{aligned}
$$

- Test your understanding:

  $\pi'(s) = \pi(s)$ for some $s \in \mathcal{S}$?     <span style="color:orange">not necessarily</span>

  $\pi'(s) \neq \pi(s)$ for some $s \in \mathcal{S}$?

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \operatorname*{argmax}_a \left[ Q^\pi(s, a) \right] \\
&= \operatorname*{argmax}_a \left[ R(s) + \gamma \sum_{s'} P(s'|s, a)\, V^\pi(s') \right] \\
&= \operatorname*{argmax}_a \left[ \sum_{s'} P(s'|s, a)\, V^\pi(s') \right]
\end{aligned}
$$

- Test your understanding:

$\pi'(s) = \pi(s)$ for some $s \in \mathcal{S}$?　　　　not necessarily

$\pi'(s) \neq \pi(s)$ for some $s \in \mathcal{S}$?　　　　not necessarily

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \arg\max_a \left[ Q^\pi(s, a) \right] \\
&= \arg\max_a \left[ R(s) + \gamma \sum_{s'} P(s'|s, a) \, V^\pi(s') \right] \\
&= \arg\max_a \left[ \sum_{s'} P(s'|s, a) \, V^\pi(s') \right]
\end{aligned}
$$

- Test your understanding:

    $\pi'(s) = \pi(s)$ for some $s \in \mathcal{S}$?      not necessarily

    $\pi'(s) \neq \pi(s)$ for some $s \in \mathcal{S}$?      not necessarily

    $Q^\pi(s, \pi'(s)) \geq Q^\pi(s, \pi(s))$ for all $s \in \mathcal{S}$?

- In terms of the state value function:

$$
\begin{aligned}
\pi'(s) &= \underset{a}{\arg\max}\left[Q^\pi(s,a)\right] \\
&= \underset{a}{\arg\max}\left[R(s) + \gamma\sum_{s'}P(s'|s,a)\,V^\pi(s')\right] \\
&= \underset{a}{\arg\max}\left[\sum_{s'}P(s'|s,a)\,V^\pi(s')\right]
\end{aligned}
$$

- Test your understanding:

| | |
|---|---|
| $\pi'(s) = \pi(s)$ for some $s \in \mathcal{S}$? | not necessarily |
| $\pi'(s) \neq \pi(s)$ for some $s \in \mathcal{S}$? | not necessarily |
| $Q^\pi(s,\pi'(s)) \geq Q^\pi(s,\pi(s))$ for all $s \in \mathcal{S}$? | TRUE |

- Greedy policy:

# Policy improvement

- Greedy policy:

$$\pi'(s) = \underset{a}{\arg\max} \; Q^\pi(s, a)$$

## Policy improvement

- Greedy policy:

$$\pi'(s) \ = \ \operatorname*{argmax}_a Q^\pi(s, a)$$

- Theorem:
  The greedy policy $\pi'(s) = \arg\max_a Q^\pi(s, a)$ improves
  everywhere on the policy $\pi$ from which it was derived:

## Policy improvement

- Greedy policy:

$$\pi'(s) \;=\; \arg\max_a Q^\pi(s, a)$$

- Theorem:
  The greedy policy $\pi'(s) = \arg\max_a Q^\pi(s, a)$ improves
  everywhere on the policy $\pi$ from which it was derived:

$$V^{\pi'}(s) \;\geq\; V^\pi(s) \quad \text{for all states } s \in \mathcal{S}$$

## Policy improvement

- Greedy policy:

$$\pi'(s) \;=\; \arg\max_a Q^\pi(s, a)$$

- Theorem:
  The greedy policy $\pi'(s) = \arg\max_a Q^\pi(s, a)$ improves
  everywhere on the policy $\pi$ from which it was derived:

$$V^{\pi'}(s) \;\geq\; V^\pi(s) \quad \text{for all states } s \in \mathcal{S}$$

- Intuition:

## Policy improvement

- **Greedy policy:**

$$\pi'(s) \;=\; \underset{a}{\arg\max}\; Q^\pi(s, a)$$

- **Theorem:**
  The greedy policy $\pi'(s) = \arg\max_a Q^\pi(s, a)$ improves everywhere on the policy $\pi$ from which it was derived:

$$V^{\pi'}(s) \;\geq\; V^\pi(s) \quad \text{for all states } s \in \mathcal{S}$$

- **Intuition:**
  If it's better to choose action $a$ in state $s$ before following $\pi$, then it's always better to make this choice.

## Policy improvement

- Greedy policy:

$$\pi'(s) \;=\; \arg\max_a Q^\pi(s, a)$$

- Theorem:
  The greedy policy $\pi'(s) = \arg\max_a Q^\pi(s, a)$ improves everywhere on the policy $\pi$ from which it was derived:

$$V^{\pi'}(s) \;\geq V^\pi(s) \quad \text{for all states } s \in \mathcal{S}$$

- Intuition:
  If it's better to choose action $a$ in state $s$ before following $\pi$, then it's always better to make this choice.

- Proof idea:
  We'll prove a key inequality for *one-step deviations* from $\pi$, then we'll extend this inequality by an iterative argument.

- Comparing value functions:

- Comparing value functions:

$$V^\pi(s) \;=\; Q^\pi(s, \pi(s))$$

- Comparing value functions:

$$
\begin{aligned}
V^\pi(s) &= Q^\pi(s, \pi(s)) \\
&\leq \max_a Q^\pi(s, a)
\end{aligned}
$$

- Comparing value functions:

$$
\begin{aligned}
V^\pi(s) &= Q^\pi(s, \pi(s)) \\
&\leq \max_a Q^\pi(s, a) \\
&= Q^\pi(s, \pi'(s))
\end{aligned}
$$

- Comparing value functions:

$$
\begin{aligned}
V^{\pi}(s) &= Q^{\pi}(s, \pi(s)) \\
&\leq \max_a Q^{\pi}(s, a) \\
&= Q^{\pi}(s, \pi'(s)) \\
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s')
\end{aligned}
$$

- Comparing value functions:

$$
\begin{aligned}
V^{\pi}(s) &= Q^{\pi}(s, \pi(s)) \\
&\leq \max_{a} Q^{\pi}(s, a) \\
&= Q^{\pi}(s, \pi'(s)) \\
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s))V^{\pi}(s')
\end{aligned}
$$

- Combining these steps:

- Comparing value functions:

$$
\begin{aligned}
V^{\pi}(s) &= Q^{\pi}(s, \pi(s)) \\
&\leq \max_{a} Q^{\pi}(s, a) \\
&= Q^{\pi}(s, \pi'(s)) \\
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s')
\end{aligned}
$$

- Combining these steps:

$$
V^{\pi}(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s')
$$

- Comparing value functions:

$$
\begin{aligned}
V^\pi(s) &= Q^\pi(s, \pi(s)) \\
&\leq \max_a Q^\pi(s, a) \\
&= Q^\pi(s, \pi'(s)) \\
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')
\end{aligned}
$$

- Combining these steps:

$$
V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')
$$

- Intuition:

- Comparing value functions:

$$
\begin{aligned}
V^\pi(s) &= Q^\pi(s, \pi(s)) \\
&\leq \max_a Q^\pi(s, a) \\
&= Q^\pi(s, \pi'(s)) \\
&= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')
\end{aligned}
$$

- Combining these steps:

$$
V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')
$$

- Intuition:

It is better to take one step under $\pi'$, then revert to $\pi$, than to always follow $\pi$.

- One-step inequality:

- One-step inequality:

$$V^\pi(s) \ \leq \ R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

What happens if we plug this inequality into itself?

- One-step inequality:

$$V^\pi(s) \ \leq \ R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s))V^\pi(s')$$

What happens if we plug this inequality into itself?
Then we obtain …

- One-step inequality:

$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

What happens if we plug this inequality into itself?
Then we obtain …

- Two-step inequality:

- One-step inequality:

$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

What happens if we plug this inequality into itself?
Then we obtain ...

- Two-step inequality:

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

- One-step inequality:

$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

What happens if we plug this inequality into itself?
Then we obtain …

- Two-step inequality:

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

- Intuition:

- One-step inequality:

$$V^\pi(s) \;\leq\; R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

What happens if we plug this inequality into itself?
Then we obtain …

- Two-step inequality:

$$V^\pi(s) \;\leq\; R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) \left[ R(s') \;+\; \gamma \sum_{s''} P(s''|s', \pi'(s')) V^\pi(s'') \right]$$

- Intuition:

It is better to take **two** steps under $\pi'$, then revert to $\pi$, than to always follow $\pi$.

- Two-step inequality:

- Two-step inequality:

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

- Two-step inequality:

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

- Apply the inequality $t$ times:

- Two-step inequality:

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

- Apply the inequality $t$ times:

  It is better to take $t$ steps under $\pi'$, then revert to $\pi$, than to always follow $\pi$.

- Two-step inequality:

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

- Apply the inequality $t$ times:

It is better to take $t$ steps under $\pi'$, then revert to $\pi$, than to always follow $\pi$. Last term is of order $O(\gamma^t)$.

- Two-step inequality:

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

- Apply the inequality $t$ times:

  It is better to take $t$ steps under $\pi'$, then revert to $\pi$, than to always follow $\pi$. Last term is of order $O(\gamma^t)$.

- Take the limit $t \to \infty$:

- Two-step inequality:

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

- Apply the inequality $t$ times:

  It is better to take $t$ steps under $\pi'$, then revert to $\pi$, than to always follow $\pi$. Last term is of order $O(\gamma^t)$.

- Take the limit $t \to \infty$:

  It is better to follow $\pi'$ (always) than to follow $\pi$ (always).

- Two-step inequality:

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

- Apply the inequality $t$ times:

  It is better to take $t$ steps under $\pi'$, then revert to $\pi$, than to always follow $\pi$. Last term is of order $O(\gamma^t)$.

- Take the limit $t \to \infty$:

  It is better to follow $\pi'$ (always) than to follow $\pi$ (always). Conclude that $V^\pi(s) \leq V^{\pi'}(s)$ for all states $s \in \mathcal{S}$.

# Policy iteration

How to compute $\pi^*$?

# Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

## Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

# Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.

# Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.

## Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a \, Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

   $\pi_0$

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

   $$\pi_0 \xrightarrow{\text{evaluate}}$$

# Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

   $$\pi_0 \quad \xrightarrow{\text{evaluate}} \quad V^{\pi_0}(s)$$

## Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{l} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{array}$$

# Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \mathrm{argmax}_a \, Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{array} \xrightarrow{\text{improve}}$$

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

   $$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}}$$

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

$$\pi_0 \quad \xrightarrow{\text{evaluate}} \quad \begin{matrix} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{matrix} \quad \xrightarrow{\text{improve}} \quad \pi_1 \quad \xrightarrow{\text{evaluate}} \quad V^{\pi_1}(s)$$

## Policy iteration

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{matrix} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{matrix} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \begin{matrix} V^{\pi_1}(s) \\ Q^{\pi_1}(s, a) \end{matrix}$$

How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \to \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_1}(s) \\ Q^{\pi_1}(s, a) \end{array} \xrightarrow{\text{improve}}$$
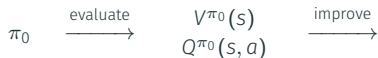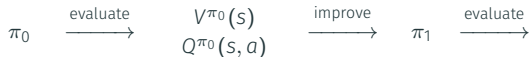
How to compute $\pi^*$?

1. Choose an initial policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$.

2. Repeat until convergence:

   Compute the action value function $Q^\pi(s, a)$.
   Compute the greedy policy $\pi'(s) = \text{argmax}_a Q^\pi(s, a)$.
   Replace $\pi$ by $\pi'$.

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{matrix} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{matrix} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \begin{matrix} V^{\pi_1}(s) \\ Q^{\pi_1}(s, a) \end{matrix} \xrightarrow{\text{improve}} \cdots$$

Policy iteration is guaranteed to terminate.

True (A) or False (B)?

# Policy iteration

- How to compute $\pi^*$?

## Policy iteration

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s, a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

## Policy iteration

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.

## Policy iteration

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.
But does it converge to an optimal policy?

## Policy iteration

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.
But does it converge to an optimal policy?

- Theorem

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.
But does it converge to an optimal policy?

- Theorem

If $\pi'(s) = \arg\max_a Q^\pi(s,a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all $s \in \mathcal{S}$,

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.
But does it converge to an optimal policy?

- Theorem

If $\pi'(s) = \arg\max_a Q^\pi(s,a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all $s \in \mathcal{S}$,
then $V^\pi(s) = V^*(s)$ for all $s \in \mathcal{S}$.

## Policy iteration

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.
But does it converge to an optimal policy?

- Theorem

If $\pi'(s) = \arg\max_a Q^\pi(s,a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all $s \in \mathcal{S}$,
then $V^\pi(s) = V^*(s)$ for all $s \in \mathcal{S}$.

- Proof idea

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

This process is guaranteed to terminate.
But does it converge to an optimal policy?

- Theorem

If $\pi'(s) = \arg\max_a Q^\pi(s,a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all $s \in \mathcal{S}$,
then $V^\pi(s) = V^*(s)$ for all $s \in \mathcal{S}$.

- Proof idea

Prove a key equality/inequality for terminal/non-terminal
policies;

- How to compute $\pi^*$?

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{array}{c} V^{\pi_0}(s) \\ Q^{\pi_0}(s,a) \end{array} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} \cdots$$

  This process is guaranteed to terminate.
  But does it converge to an optimal policy?

- Theorem

  If $\pi'(s) = \arg\max_a Q^\pi(s,a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all $s \in \mathcal{S}$,
  then $V^\pi(s) = V^*(s)$ for all $s \in \mathcal{S}$.

- Proof idea

  Prove a key equality/inequality for terminal/non-terminal
  policies; iterate $t$ times, then compare the limits as $t \to \infty$.

- Suppose policy iteration converges to $\pi'$.

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) \;=\; R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s') \qquad \boxed{\text{Bellman equation}}$$

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) \quad = \quad R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\pi}(s) \quad = \quad R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s') \qquad \boxed{\text{at convergence}}$$

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) \;=\; R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\pi}(s) \;=\; R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s') \qquad \boxed{\text{at convergence}}$$

Now exploit that $\pi'$ is greedy with respect to $\pi$ …

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s))V^{\pi'}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s))V^{\pi}(s') \qquad \boxed{\text{at convergence}}$$

Now exploit that $\pi'$ is greedy with respect to $\pi$ …

- Bellman optimality equation

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s')$$   Bellman equation

$$V^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s')$$   at convergence

Now exploit that $\pi'$ is greedy with respect to $\pi$ …

- Bellman optimality equation

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

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) \;=\; R(s) \;+\; \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s')$$  Bellman equation

$$V^{\pi}(s) \;=\; R(s) \;+\; \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s')$$  at convergence

Now exploit that $\pi'$ is greedy with respect to $\pi$ …

- Bellman optimality equation

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

These equations are nonlinear due to the max operation.

- Suppose policy iteration converges to $\pi'$.

$$V^{\pi'}(s) \;=\; R(s) \,+\, \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi'}(s')$$  | Bellman equation |

$$V^{\pi}(s) \;=\; R(s) \,+\, \gamma \sum_{s'} P(s'|s, \pi'(s)) V^{\pi}(s')$$  | at convergence |

Now exploit that $\pi'$ is greedy with respect to $\pi$ …

- Bellman optimality equation

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

These equations are nonlinear due to the max operation.
There are $n$ equations for $n$ unknowns (where $s = 1, 2, \ldots, n$).

- Let $\tilde{\pi}$ be any policy of the MDP:

- Let $\tilde{\pi}$ be any policy of the MDP:

$$V^{\tilde{\pi}}(s) \;=\; R(s) \,+\, \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}}$$

- Let $\tilde{\pi}$ be any policy of the MDP:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) &= R(s) + \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}} \\
V^{\tilde{\pi}}(s) &\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \qquad \boxed{\text{greedy}}
\end{aligned}
$$

- Let $\tilde{\pi}$ be any policy of the MDP:

$$V^{\tilde{\pi}}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\tilde{\pi}}(s) \leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \qquad \boxed{\text{greedy}}$$

- Compare to Bellman optimality equation (BOE):

- Let $\tilde{\pi}$ be any policy of the MDP:

$$V^{\tilde{\pi}}(s) \;=\; R(s) \;+\; \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\tilde{\pi}}(s) \;\leq\; R(s) \;+\; \gamma \max_{a} \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \qquad \boxed{\text{greedy}}$$

- Compare to Bellman optimality equation (BOE):

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

- Let $\tilde{\pi}$ be any policy of the MDP:

$$V^{\tilde{\pi}}(s) \;=\; R(s) \,+\, \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\tilde{\pi}}(s) \;\leq\; R(s) \,+\, \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \qquad \boxed{\text{greedy}}$$

- Compare to Bellman optimality equation (BOE):

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

- Understanding the difference:

- Let $\tilde{\pi}$ be any policy of the MDP:

$$V^{\tilde{\pi}}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\tilde{\pi}}(s) \leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \qquad \boxed{\text{greedy}}$$

- Compare to Bellman optimality equation (BOE):

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

- Understanding the difference:

  The inequality holds for any policy $\tilde{\pi}$ of the MDP.

# Proof — 2. Inequality

- Let $\tilde{\pi}$ be any policy of the MDP:

$$V^{\tilde{\pi}}(s) \;=\; R(s) \,+\, \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s') \qquad \boxed{\text{Bellman equation}}$$

$$V^{\tilde{\pi}}(s) \;\leq\; R(s) \,+\, \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \qquad \boxed{\text{greedy}}$$

- Compare to Bellman optimality equation (BOE):

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

- Understanding the difference:

  The inequality holds for any policy $\tilde{\pi}$ of the MDP.
  The **BOE** only holds for a solution $\pi$ from policy iteration.

- Iterating the inequality:

- Iterating the inequality:

$$V^{\tilde{\pi}}(s) \quad \leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s')$$

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) \quad &\leq \quad R(s) + \gamma \max_{a} \sum\nolimits_{s'} P(s'|s,a) V^{\tilde{\pi}}(s') \\
&\leq \quad R(s) + \gamma \max_{a} \sum\nolimits_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum\nolimits_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) \quad &\leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\tilde{\pi}}(s') \\
&\leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the BOE:

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) &\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\tilde{\pi}}(s') \\
&\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the BOE:

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

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) &\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\tilde{\pi}}(s') \\
&\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the BOE:

$$
\begin{aligned}
V^{\pi}(s) &= R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\pi}(s') \\
&= R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\pi}(s'') \right]
\end{aligned}
$$

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) &\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\tilde{\pi}}(s') \\
&\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the BOE:

$$
\begin{aligned}
V^{\pi}(s) &= R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\pi}(s') \\
&= R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\pi}(s'') \right]
\end{aligned}
$$

- Iterating $t$ times:

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) \quad &\leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\tilde{\pi}}(s') \\
&\leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the BOE:

$$
\begin{aligned}
V^{\pi}(s) \quad &= \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\pi}(s') \\
&= \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\pi}(s'') \right]
\end{aligned}
$$

- Iterating $t$ times:

  Both right sides agree up to term of order $\gamma^t$.

## Proof — 3. Taking the limit

- Iterating the inequality:

$$
\begin{aligned}
V^{\tilde{\pi}}(s) \quad &\leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s') \\
&\leq \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s', a') V^{\tilde{\pi}}(s'') \right]
\end{aligned}
$$

- Iterating the BOE:

$$
\begin{aligned}
V^{\pi}(s) \quad &= \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\pi}(s') \\
&= \quad R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s', a') V^{\pi}(s'') \right]
\end{aligned}
$$

- Iterating $t$ times:

  Both right sides agree up to term of order $\gamma^t$.
  Taking the limit $t \to \infty$, we find $V^{\tilde{\pi}}(s) \leq V^{\pi}(s)$ for all $s \in \mathcal{S}$.

- Iterating the inequality:

$$V^{\tilde{\pi}}(s) \leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^{\tilde{\pi}}(s')$$

$$\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s',a') V^{\tilde{\pi}}(s'') \right]$$

- Iterating the BOE:

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

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

- Iterating $t$ times:

  Both right sides agree up to term of order $\gamma^t$.
  Taking the limit $t \to \infty$, we find $V^{\tilde{\pi}}(s) \leq V^{\pi}(s)$ for all $s \in \mathcal{S}$.

  Since $\tilde{\pi}$ is arbitrary, we conclude that $\pi$ is optimal .

That's all folks!