# Planning in Markov Decision Process

## Sham Kakade and Wen Sun

**CS 6789: Foundations of Reinforcement Learning**

# Announcements

HW0: due this 9/9 11:59pm ET
Gradescope (please self-enroll)

# Announcements

HW0: due this 9/9 11:59pm ET
Gradescope (please self-enroll)

Waiting List

# Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$

# Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

# Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

Value function $V^\pi(s) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,\middle|\, s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\,\cdot\,|\,s_h, a_h)\right]$

# Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

Value function $V^{\pi}(s) = \mathbb{E}\left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,\bigg|\, s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \,|\, s_h, a_h) \right]$

Q function $Q^{\pi}(s, a) = \mathbb{E}\left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,\bigg|\, (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \,|\, s_h, a_h) \right]$

# Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$

**Theorem 1: Bellman Optimality (Q-version)**

$$Q^\star(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a' \in A} Q^\star(s', a') \right]$$

# Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \to [0,1], \quad \gamma \in [0,1)$$

**Theorem 1: Bellman Optimality (Q-version)**

$$Q^\star(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ \max_{a' \in A} Q^\star(s', a') \right]$$

**Theorem 2 (Q-version):**

For any $Q : S \times A \to \mathbb{R}$, if $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ \max_{a'} Q(s', a') \right]$,

$\forall s, a$, then $Q(s, a) = Q^\star(s, a), \forall s, a$

# Main Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ , How to find $\pi^\star$ (stationary & deterministic)

# Main Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ , How to find $\pi^{\star}$ (stationary & deterministic)

Two Approaches:

1. Value Iteration
2. Policy Iteration

# Define Bellman Operator $\mathscr{T}$:

Given a function $f : S \times A \mapsto \mathbb{R}$,

function $\hookleftarrow$ $\boxed{\mathscr{T}f}$ : $S \times A \mapsto \mathbb{R}$,

$(\mathscr{T}f)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a' \in A} f(s', a'), \forall s, a \in S \times A$

# Value Iteration Algorithm:

$$Q^0(\cdot,\cdot) \in \mathbb{R}^{|S| \times |A|}$$

$$1 + \gamma + \gamma^2 + \gamma^3 + \cdots = \frac{1}{1-\gamma}$$

1. Initialization: $Q^0 : \|Q^0\|_\infty \in (0, \frac{1}{1-\gamma})$

2. Iterate until convergence: $Q^{t+1} = \mathcal{T} Q^t$

$$\Rightarrow \forall s, a, \quad \text{see} \quad Q^{t+1}(s,a) \leftarrow r(s,a) + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \max_{a'} Q(s',a')$$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathscr{T} Q^\star$$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathscr{T}Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathscr{T}f$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathcal{T} Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathcal{T} f$

Consider the simple problem: finding fixed point solution $x^\star = \ell(x^\star)$

$\ell : [a, b] \mapsto [a, b]$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathcal{T} Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathcal{T} f$

Consider the simple problem: finding fixed point solution $x^\star = \ell(x^\star)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \ldots,$$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathcal{T}Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^\star = \ell(x^\star)$

$$x_0, x_{t+1} = \ell(x_t), t = 0,\ldots,$$

$$|x_t - x^\star| =$$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathcal{T} Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathcal{T} f$

Consider the simple problem: finding fixed point solution $x^\star = \ell(x^\star)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \ldots,$$

$$|x_t - x^\star| = |\ell(x_{t-1}) - \ell(x^\star)|$$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathscr{T} Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathscr{T} f$

Consider the simple problem: finding fixed point solution $x^\star = \ell(x^\star)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \ldots,$$

$$|x_t - x^\star| = |\ell(x_{t-1}) - \ell(x^\star)| \leq L |x_{t-1} - x^\star|$$

# Intuition:

Via Bellman optimality theorem:

$$Q^\star = \mathcal{T}Q^\star$$

i.e., $Q^\star$ is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^\star = \ell(x^\star)$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \ldots,$$

$$|x_t - x^\star| = |\ell(x_{t-1}) - \ell(x^\star)| \leq L|x_{t-1} - x^\star| \;\leq L^2 \left| x_{t-2} - x^\star \right|$$

If $L < 1$ (i.e., contraction), then it converges exponentially fast

# Convergence of Value Iteration:

*Lemma [contraction]*: Given any $Q, Q'$, we have:
$$\|\mathscr{T}Q - \mathscr{T}Q'\|_\infty \leq \gamma\|Q - Q'\|_\infty$$

**Proof:**

# Convergence of Value Iteration:

**Lemma [contraction]**: Given any $Q, Q'$, we have:

$$\|\mathscr{T}Q - \mathscr{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

**Proof:**

$$|\mathscr{T}Q(s,a) - \mathscr{T}Q'(s,a)| = \left| r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q(s',a') - \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q'(s',a') \right) \right|$$

# Convergence of Value Iteration:

**Lemma [contraction]**: Given any $Q, Q'$, we have:
$$\|\mathscr{T}Q - \mathscr{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

**Proof:**

$$|\mathscr{T}Q(s,a) - \mathscr{T}Q'(s,a)| = \left| r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q(s',a') - \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma \sum_{s'} P(s'|s,a) \left| \left( \max_{a'} Q(s',a') - \max_{a'} Q'(s',a') \right) \right|$$

# Convergence of Value Iteration:

**Lemma [contraction]**: Given any $Q, Q'$, we have:
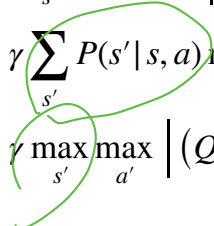$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma\|Q - Q'\|_\infty$$

**Proof:**

$$|\mathcal{T}Q(s,a) - \mathcal{T}Q'(s,a)| = \left| r(s,a) + \gamma\mathbb{E}_{s'\sim P(s,a)}\max_{a'} Q(s',a') - \left( r(s,a) + \gamma\mathbb{E}_{s'\sim P(s,a)}\max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma \sum_{s'} P(s'|s,a) \left| \left( \max_{a'} Q(s',a') - \max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma \sum_{s'} P(s'|s,a) \max_{a'} \left| \left( Q(s',a') - Q'(s',a') \right) \right|$$

# Convergence of Value Iteration:

*Lemma [contraction]*: Given any $Q, Q'$, we have:
$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma\|Q - Q'\|_\infty$$

*Proof:*

$$|\mathcal{T}Q(s,a) - \mathcal{T}Q'(s,a)| = \left| r(s,a) + \gamma\mathbb{E}_{s'\sim P(s,a)}\max_{a'} Q(s',a') - \left( r(s,a) + \gamma\mathbb{E}_{s'\sim P(s,a)}\max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma\sum_{s'} P(s'\,|\,s,a)\left| \left( \max_{a'} Q(s',a') - \max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma\sum_{s'} P(s'\,|\,s,a)\max_{a'}\left| \left( Q(s',a') - Q'(s',a') \right) \right|$$

$$\leq \gamma\max_{s'}\max_{a'}\left| \left( Q(s',a') - Q'(s',a') \right) \right|$$

# Convergence of Value Iteration:

***Lemma [contraction]***: Given any $Q, Q'$, we have:
$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma\|Q - Q'\|_\infty$$

***Proof:***

$$|\mathcal{T}Q(s,a) - \mathcal{T}Q'(s,a)| = \left| r(s,a) + \gamma\mathbb{E}_{s'\sim P(s,a)}\max_{a'} Q(s',a') - \left( r(s,a) + \gamma\mathbb{E}_{s'\sim P(s,a)}\max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma\sum_{s'} P(s'|s,a)\left| \left( \max_{a'} Q(s',a') - \max_{a'} Q'(s',a') \right) \right|$$

$$\leq \gamma\sum_{s'} P(s'|s,a)\max_{a'}\left| \left( Q(s',a') - Q'(s',a') \right) \right|$$

$$\leq \gamma\max_{s'}\max_{a'}\left| \left( Q(s',a') - Q'(s',a') \right) \right| \quad = \gamma\|Q - Q'\|_\infty$$

# Convergence of Value Iteration:

*Lemma [Convergence]*: Given $Q^0$, we have:

$$\|Q^t - Q^\star\|_\infty \leq \gamma^t \|Q^0 - Q^\star\|_\infty$$

*Proof:*

# Convergence of Value Iteration:

**Lemma [Convergence]**: Given $Q^0$, we have:
$$\|Q^t - Q^\star\|_\infty \leq \gamma^t \|Q^0 - Q^\star\|_\infty$$

**Proof:**

$$\|Q^{t+1} - Q^\star\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^\star\|_\infty \leq \gamma \|Q^t - Q^\star\|_\infty$$

# Convergence of Value Iteration:

**Lemma [Convergence]**: Given $Q^0$, we have:
$$\|Q^t - Q^\star\|_\infty \leq \gamma^t \|Q^0 - Q^\star\|_\infty$$

**Proof:**

$$\|Q^{t+1} - Q^\star\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^\star\|_\infty \leq \gamma\|Q^t - Q^\star\|_\infty$$

$$\leq \gamma^2 \|Q^{t-1} - Q^\star\|_\infty$$

$$\ldots \leq \gamma^{t+1} \|\widehat{Q}^0 - Q^\star\|_\infty$$

$$\pi^\star(s) = \arg\max_a Q^\star(s,a)$$
$$\pi(s) = \arg\max_a Q^t(s,a)$$

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s, a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^\star(s) - \dfrac{2\gamma^t}{1-\gamma}\|Q^0 - Q^\star\|_\infty \forall s \in S$

**Proof:**

$$\left| V^{\pi^t}(s) - V^\star(s) \right| \leq O\left(\frac{\gamma^t}{1-\gamma}\right)$$

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s, a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^\star(s) - \dfrac{2\gamma^t}{1 - \gamma}\|Q^0 - Q^\star\|_\infty \forall s \in S$

***Proof:***

$$V^{\pi^t}(s) - V^\star(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s,a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^\star(s) - \dfrac{2\gamma^t}{1-\gamma}\|Q^0 - Q^\star\|_\infty \forall s \in S$

### Proof:

$$V^{\pi^t}(s) - V^\star(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^t(s)) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s, a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^\star(s) - \dfrac{2\gamma^t}{1 - \gamma}\|Q^0 - Q^\star\|_\infty \forall s \in S$

## Proof:

$$V^{\pi^t}(s) - V^\star(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= \underbrace{Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^t(s))} + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= \gamma\mathbb{E}_{s'\sim P(s,\pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s, a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^{\star}(s) - \dfrac{2\gamma^t}{1-\gamma}\|Q^0 - Q^{\star}\|_\infty \, \forall s \in S$

## Proof:

$$V^{\pi^t}(s) - V^{\star}(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^{\star}(s, \pi^{\star}(s))$$

$$= Q^{\pi^t}(s, \pi^t(s)) - Q^{\star}(s, \pi^t(s)) + Q^{\star}(s, \pi^t(s)) - Q^{\star}(s, \pi^{\star}(s))$$

$$= \gamma\mathbb{E}_{s'\sim P(s,\pi^t(s))}\left(V^{\pi^t}(s') - V^{\star}(s')\right) + Q^{\star}(s, \pi^t(s)) - Q^{\star}(s, \pi^{\star}(s))$$

$$\geq \gamma\mathbb{E}_{s'\sim P(s,\pi^t(s))}\left(V^{\pi^t}(s') - V^{\star}(s')\right) + Q^{\star}(s, \pi^t(s)) \boxed{- Q^t(s, \pi^t(s)) + Q^t(s, \pi^{\star}(s))} - Q^{\star}(s, \pi^{\star}(s))$$

$$\leq 0$$

$$Q^t(s, \pi^{\star}(s)) \leq \max_a Q^t(s,a)$$
$$\geq Q^t(s, \pi^t(s))$$

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s, a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^\star(s) - \dfrac{2\gamma^t}{1-\gamma}\|Q^0 - Q^\star\|_\infty \forall s \in S$

## Proof:

$$V^{\pi^t}(s) - V^\star(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^t(s)) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= \gamma\mathbb{E}_{s'\sim P(s,\pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$\geq \gamma\mathbb{E}_{s'\sim P(s,\pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) + Q^\star(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + \boxed{Q^t(s, \pi^\star(s)) - Q^\star(s, \pi^\star(s))}$$

$$\geq \gamma\mathbb{E}_{s'\sim P(s,\pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) - 2\gamma^t\|Q^0 - Q^\star\|_\infty$$

$\geq -\gamma^t\|Q^0 - Q^\star\|_\infty$

Apply the same procedure

# Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg\max_a Q^t(s, a)$$

**Theorem:** $V^{\pi^t}(s) \geq V^\star(s) - \dfrac{2\gamma^t}{1-\gamma}\|Q^0 - Q^\star\|_\infty \forall s \in S$

***Proof:***

$$V^{\pi^t}(s) - V^\star(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^t(s)) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$= \gamma\mathbb{E}_{s' \sim P(s, \pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s))$$

$$\geq \gamma\mathbb{E}_{s' \sim P(s, \pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) + Q^\star(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^\star(s)) - Q^\star(s, \pi^\star(s))$$

$$\geq \gamma\mathbb{E}_{s' \sim P(s, \pi^t(s))}\left(V^{\pi^t}(s') - V^\star(s')\right) - 2\gamma^t\|Q^0 - Q^\star\|_\infty \quad \dots\text{Recursion}$$

# Main Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ , How to find $\pi^\star$ (stationary & deterministic)

Two Approaches:

✅ 1. Value Iteration

2. Policy Iteration

# Policy Iteration Algorithm:

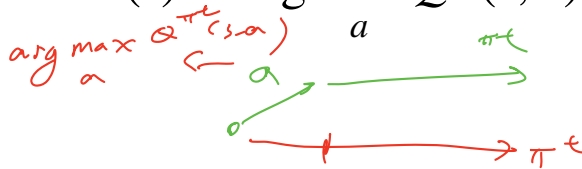1. Initialization: $\pi^0 : S \mapsto A$

# Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

# Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

3. Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

# Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

3. Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

# Policy Iteration Algorithm:

Closed-form for PE
(see 1.1.3 in Monograph)

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

3. Policy Improvement $\pi^{t+1}(s) = \arg\max_{a} Q^{\pi^t}(s, a), \forall s$

$\forall s a$

$Q^{\pi^t}(sa) = r(sa) + \gamma \underset{s' \sim p(sa)}{E} Q^{\pi^t}(s', \pi(s'))$

$|S||A|$ variable & constraints

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max\limits_{a} Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s,a)}\left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$\leftarrow$ Bell-Eqn

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$= \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$= \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$\geq 0$

$$\geq \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right]$$

$$\geq \gamma \mathbb{E}_{s'' \sim P(s', \pi^{t+1}(s'))} \left[ Q^{\pi^{t+1}}(s'', \pi^{t+1}(s'')) - Q^{\pi^t}(s'', \pi^{t+1}(s'')) \right]$$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$= \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$\geq \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right] \geq \ldots, \geq - \gamma^\infty/(1-\gamma) = 0$$

$\in \left[ -\frac{1}{1-\gamma}, \frac{1}{1-\gamma} \right]$

$\geq -\frac{\gamma}{1-\gamma}$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$= \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$\geq \gamma \mathbb{E}_{s' \sim P(s,a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right] \geq \ldots, \geq -\gamma^\infty/(1 - \gamma) = 0$$

$$V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s), \forall s$$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s,a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma\|V^{\pi^t} - V^\star\|_\infty$

$$\cdots \leq \gamma^t \|V^{\pi^0} - V^*\|_\infty$$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s,a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

$$V^\star(s) - V^{\pi^{t+1}}(s) = \max_a \left[ r(s,a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s,\pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right]$$

bell-opt

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

$$V^\star(s) - V^{\pi^{t+1}}(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s,\pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right]$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s,\pi^{t+1}(s))} V^{\pi^t}(s') \right]$$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

$$V^\star(s) - V^{\pi^{t+1}}(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right]$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right]$$

$$= \max_a (r(s, a) + \mathbb{E}_{s' \sim P(s,a)} \gamma V^\star(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^{\pi^t}(s'))$$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

$$V^\star(s) - V^{\pi^{t+1}}(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right]$$

$$\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right]$$

$:= Q^{\pi^t}(s, \pi^{t+1}(s))$

$$= \max_a (r(s, a) + \mathbb{E}_{s' \sim P(s,a)} \gamma V^\star(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^{\pi^t}(s'))$$

$$\leq \max_a \left( r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') - \left( r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^{\pi^t}(s') \right) \right)$$

$\max_x f(x) - \max_x g(x) \leq \max_x \left[ f(x) - g(x) \right]$

# Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg\max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma\|V^{\pi^t} - V^\star\|_\infty$

$\forall s$

$$V^\star(s) - V^{\pi^{t+1}}(s) = \max_a \left[ r(s, a) + \gamma\mathbb{E}_{s'\sim P(s,a)}V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma\mathbb{E}_{s'\sim P(s,\pi^{t+1}(s))}V^{\pi^{t+1}}(s') \right]$$

$$\leq \max_a \left[ r(s, a) + \gamma\mathbb{E}_{s'\sim P(s,a)}V^\star(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma\mathbb{E}_{s'\sim P(s,\pi^{t+1}(s))}V^{\pi^t}(s') \right]$$

$$= \max_a(r(s, a) + \mathbb{E}_{s'\sim P(s,a)}\gamma V^\star(s')) - \max_a(r(s, a) + \gamma\mathbb{E}_{s'\sim P(s,a)}V^{\pi^t}(s'))$$

$$\leq \max_a \left( r(s, a) + \gamma\mathbb{E}_{s'\sim P(s,a)}V^\star(s') - \left( r(s, a) + \gamma\mathbb{E}_{s'\sim P(s,a)}V^{\pi^t}(s') \right) \right)$$

$$\leq \gamma\|V^\star - V^{\pi^t}\|_\infty \qquad \Rightarrow \|V^\star - V^{\pi^{t+1}}\|_\infty \leq \gamma\|V^\star - V^{\pi^t}\|_\infty$$

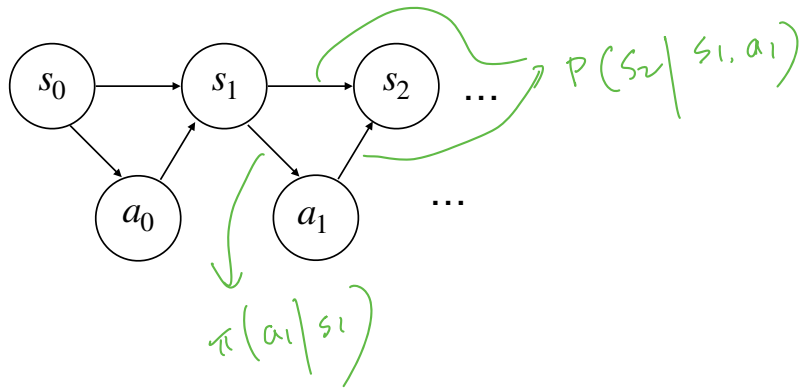# Value Iteration vs Policy Iteration?

Which one is faster?

How many iterations (computation complexity) need to find the EXACT optimal policy?

# Trajectory distribution and state-action distribution

Q: what is the probability of $\pi$ generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \ldots, s_h, a_h\}$?

# Trajectory distribution and state-action distribution

Q: what is the probability of $\pi$ generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \ldots, s_h, a_h\}$?
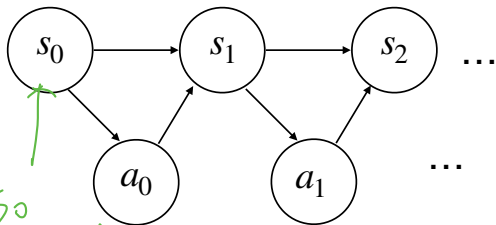


$P(s_2 \mid s_1, a_1)$

$\pi(a_1 \mid s_1)$

# Trajectory distribution and state-action distribution

Q: what is the probability of $\pi$ generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \ldots, s_h, a_h\}$?



$\mathbb{P}^\pi(s_0, a_0, \ldots, s_h, a_h)$

$= \pi(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\pi(a_1 \,|\, s_1)P(s_2 \,|\, s_1, a_1)\ldots P(s_h \,|\, s_{h-1}, a_{h-1})\pi(a_h \,|\, s_h)$

So
is fixed

# Trajectory distribution and state-action distribution

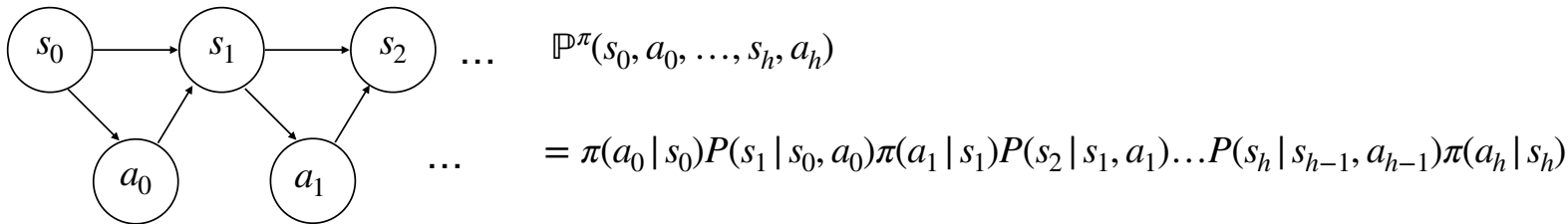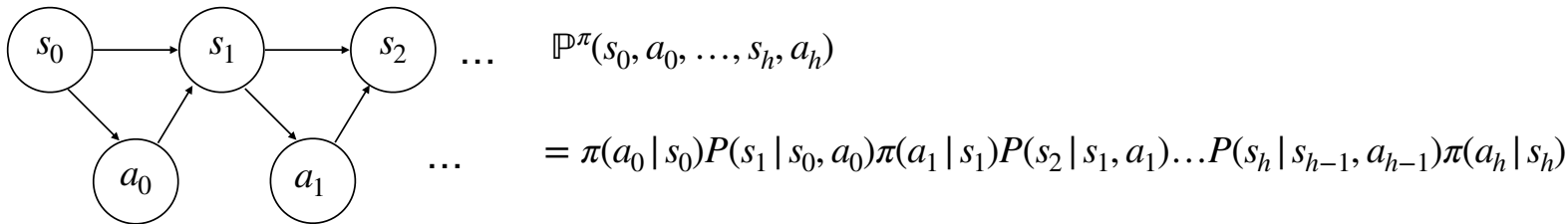Q: what is the probability of $\pi$ generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \ldots, s_h, a_h\}$?



$\mathbb{P}^\pi(s_0, a_0, \ldots, s_h, a_h)$

$= \pi(a_0 \,|\, s_0) P(s_1 \,|\, s_0, a_0) \pi(a_1 \,|\, s_1) P(s_2 \,|\, s_1, a_1) \ldots P(s_h \,|\, s_{h-1}, a_{h-1}) \pi(a_h \,|\, s_h)$

Q: what's the probability of $\pi$ visiting state ($s$,a) at time step h?

# Trajectory distribution and state-action distribution

Q: what is the probability of $\pi$ generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \ldots, s_h, a_h\}$?



$$\mathbb{P}^\pi(s_0, a_0, \ldots, s_h, a_h)$$

$$= \pi(a_0 \mid s_0)P(s_1 \mid s_0, a_0)\pi(a_1 \mid s_1)P(s_2 \mid s_1, a_1)\ldots P(s_h \mid s_{h-1}, a_{h-1})\pi(a_h \mid s_h)$$

Q: what's the probability of $\pi$ visiting state $(s,a)$ at time step h?

$$\mathbb{P}_h(s, a; s_0, \pi) = \sum_{a_0, s_1, a_1, \ldots, s_{h-1}, a_{h-1}} \mathbb{P}^\pi(s_0, a_0, \ldots, s_{h-1}, a_{h-1}s_h = s, a_h = a)$$

# State action occupancy measure

$\mathbb{P}_h(s, a; s_0, \pi)$: probability of $\pi$ visiting $(s, a)$ at time step $h \in \mathbb{N}$, starting at $s_0$

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h(s, a; s_0, \pi)$$

$\sum_{sa} \mathbb{P}_h(sa) = 1$

normalizing

$\sum_{sa} d_{s_0}^{\pi}(s, a) = 1$

$h$     $h-1$     $h-2$

$\gamma^h$     $\gamma^{h+1}$     $\gamma^{h+2}$

# State action occupancy measure

$\mathbb{P}_h(s, a; s_0, \pi)$: probability of $\pi$ visiting $(s, a)$ at time step $h \in \mathbb{N}$, starting at $s_0$

$$d_{s_0}^{\pi}(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h(s, a; s_0, \pi)$$

$$V^{\pi}(s_0) = \frac{1}{1 - \gamma} \sum_{s,a} d_{s_0}^{\pi}(s, a) r(s, a)$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^{\pi}} \left[ r(s,a) \right]$$

# Summary for today

**Two planning algorithms (no learning so far):**

**VI**: fixed point iteration $Q^{t+1} = \mathcal{T}Q^t$

Key property: it's a contraction map

# Summary for today

**Two planning algorithms (no learning so far):**

**VI**: fixed point iteration $Q^{t+1} = \mathcal{T} Q^t$

Key property: it's a contraction map

**PI**: $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a)$

Key property: monotonic improvement $V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s), \forall s$