

Planning in Markov Decision Processes

Announcements

HW0: due this Thursday 11:59pm
Gradescope (please self-enroll)

Announcements

HW0: due this Thursday 11:59pm
Gradescope (please self-enroll)

Waiting List

Recap: Infinite Horizon MDPs

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

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

Recap: Infinite Horizon MDPs

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

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

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

Recap: Infinite Horizon MDPs

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

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

Stationary Policy $\pi : S \mapsto \Delta(A)$ $\xrightarrow{\pi^*}$ $\xrightarrow{A^S}$ deterministic

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

Recap: Infinite Horizon MDPs

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

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

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

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

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

Recap: Bellman Optimality

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

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

Theorem 1: Bellman Optimality

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

Recap: Bellman Optimality

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

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

Theorem 1: Bellman Optimality

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

Bell consistent

Theorem 2:

For any $Q : S \times A \rightarrow \mathbb{R}$, if $\boxed{Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a'} Q(s', a') \right]}$ for all s, a , then $Q(s, a) = Q^*(s, a), \forall s, a$

Main Question for Today:

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

P. r. are known

$$\begin{array}{l} \forall s, a, s', P(s'|s, a) \quad \checkmark \\ \forall s, a, r(s, a) \quad \checkmark \end{array} \quad \Rightarrow \pi^*$$

Main Question for Today:

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

Two Approaches:

- 1. Value Iteration
 - 2. Policy Iteration
- $\hat{Q} \approx Q^*$

$$f: S \times A \rightarrow \mathbb{R}$$

Define Bellman Operator $\mathcal{T}: \mathcal{T}f: S \times A \rightarrow \mathbb{R}$

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

$$\boxed{\mathcal{T}f: S \times A \mapsto \mathbb{R}},$$

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

Example: $f = Q^*$

$$\mathcal{T}Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} Q^*(s', a') = Q^*(s, a)$$
$$\mathcal{T}Q^* = Q^*$$

Value Iteration Algorithm:

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

2. Iterate until convergence: $\widehat{Q}^{t+1} = \mathcal{T}\widehat{Q}^t \quad t \geq 0, \dots, \infty$

$$\sum_{i=0}^{\infty} \gamma^i = \frac{1}{1-\gamma}$$

$$Q^* \leftarrow \mathcal{T}Q^*$$

Intuition:

Via Bellman optimality theorem:

$$Q^* = \mathcal{T}Q^* \xleftarrow{\text{Fixed point}}$$

Intuition:

Via Bellman optimality theorem:

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

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

Intuition:

Via Bellman optimality theorem:

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

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

$$\ell: \mathbb{R} \rightarrow \mathbb{R}$$

Consider the simple problem: finding fixed point solution $x^* = \underset{\Delta}{\ell}(x^*)$

Intuition:

Via Bellman optimality theorem:

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

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

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

$$x_0, \underbrace{x_{t+1}}_{\Delta} = \ell(x_t), t = 0, \dots,$$

Intuition:

Via Bellman optimality theorem:

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

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

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

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots, \quad x_0 \xrightarrow{} x^*$$

$$\underline{|x_t - x^*|} =$$

Intuition:

Via Bellman optimality theorem:

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

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

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

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

$$\underline{\underline{= \ell(x_{t+1})}}$$

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

Intuition:

Via Bellman optimality theorem:

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

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

$\ell: \mathbb{R} \rightarrow \mathbb{R}$

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

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

Lipchiz constant

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

Intuition:

Via Bellman optimality theorem:

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

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

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

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

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

$\cdots L^t |x_0 - x^*|$
 $L < 1$ $t \rightarrow \infty$
 $L^t \rightarrow 0$

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

Convergence of Value Iteration:

Bellman-operator

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

$$\|\mathcal{T}_\Delta Q - \mathcal{T}_\Delta Q'\|_\infty \leq \gamma \|\underline{Q} - \underline{Q}'\|_\infty \quad \|f\|_\infty \\ = \max_{s,a} |f(s,a)|$$

Proof:

Convergence of Value Iteration:

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

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

$\forall s, a \in S \times A$

by def of Bellman operator

Proof:

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

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:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \underbrace{\max_{a'} Q(s', a')}_{\text{red}} - \underbrace{\max_{a'} Q'(s', a')}_{\text{red}} \right| \end{aligned}$$

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:

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

Assume $f \geq 0$

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:

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

Convergence of Value Iteration:

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

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

$\forall s, a \in S \times A$

Proof:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \\ &\leq \gamma \max_{s'} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| = \gamma \|Q - Q'\|_\infty \end{aligned}$$

$\Rightarrow \|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$
 $\gamma < 1$

Convergence of Value Iteration:

Lemma [Convergence]: Given \widehat{Q}^0 , we have:

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

Proof:

Convergence of Value Iteration:

Lemma [Convergence]: Given \widehat{Q}^0 , we have:

$$\|\widehat{Q}^t - Q^*\|_\infty \leq \gamma^t \|\widehat{Q}^0 - Q^*\|_\infty$$

Proof: $\widehat{Q}^* = \mathcal{T}\widehat{Q}^*$

$$\|\widehat{Q}^{t+1} - Q^*\|_\infty = \|\mathcal{T}\widehat{Q}^t - \mathcal{T}Q^*\|_\infty \leq \gamma \|\widehat{Q}^t - Q^*\|_\infty \quad \leftarrow \text{Recursion}$$

Convergence of Value Iteration:

Lemma [Convergence]: Given \widehat{Q}^0 , we have:

$$\|\widehat{Q}^t - Q^*\|_\infty \leq \gamma^t \|\widehat{Q}^0 - Q^*\|_\infty$$

Proof:

$$\|\widehat{Q}^{t+1} - Q^*\|_\infty = \underbrace{\|\mathcal{T}\widehat{Q}^t - \mathcal{T}Q^*\|_\infty}_{\mathcal{T}\text{-operator}} \leq \gamma \|\widehat{Q}^t - Q^*\|_\infty$$

$$\dots \leq \underbrace{\gamma^{t+1} \|\widehat{Q}^0 - Q^*\|_\infty}_{\text{Recursion}}$$

$$\gamma < 1, \quad t \rightarrow \infty \quad \gamma^t \rightarrow 0$$

$$\widehat{Q}_{(t,\infty)}^t \rightarrow Q^*(s,a) + \gamma a$$

① Bell-operator
② Contraction

③ Recursion

Final Quality of the Policy:

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

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

$$\begin{aligned}\hat{Q} &\rightarrow Q^\star \\ \pi^\star &= \arg \max_a \hat{Q}^\star(s, a) \\ \hat{\pi} &= \arg \max_a \hat{Q}(s, a)\end{aligned}$$

Proof:

Final Quality of the Policy:

$$\hat{\pi}^t : \hat{\pi}^t(s) = \arg \max_a \widehat{Q}^t(s, a)$$

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

$s \in S$

Bell-Eqn

Proof:

$$\underbrace{V^{\hat{\pi}^t}(s) - V^\star(s)}_{\text{Bell-Eqn}} = \underbrace{Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s))}_{}$$

Final Quality of the Policy:

$$\hat{\pi}^t : \hat{\pi}^t(s) = \arg \max_a \widehat{Q}^t(s, a)$$

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

Proof:

$$\begin{aligned} V^{\hat{\pi}^t}(s) - V^\star(s) &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\ &\stackrel{\Delta}{=} Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - \underbrace{Q^\star(s, \hat{\pi}^t(s))}_{+/-} + \underbrace{Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s))}_{\text{red wavy line}} \end{aligned}$$

Final Quality of the Policy:

$$\hat{\pi}^t : \hat{\pi}^t(s) = \arg \max_a \widehat{Q}^t(s, a)$$

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

Proof:

$$\begin{aligned} \underbrace{V^{\hat{\pi}^t}(s) - V^\star(s)}_{\text{Bellman Error}} &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\ &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \hat{\pi}^t(s)) + Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(\underbrace{V^{\hat{\pi}^t}(s') - V^\star(s')}_{\Delta} \right) + \underbrace{Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s))}_{\Delta} \end{aligned}$$

Final Quality of the Policy:

$$\hat{\pi}^t : \hat{\pi}^t(s) = \arg \max_a \widehat{Q}^t(s, a)$$

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

Proof:

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

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

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

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

≤ 0

Final Quality of the Policy:

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

output of VI

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

Proof:



$$\begin{aligned}
 V^{\hat{\pi}^t}(s) - V^\star(s) &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\
 &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \hat{\pi}^t(s)) + Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\
 &= \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(V^{\hat{\pi}^t}(s') - V^\star(s') \right) + Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\
 &\geq \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(V^{\hat{\pi}^t}(s') - V^\star(s') \right) + \boxed{Q^\star(s, \hat{\pi}^t(s)) - \hat{Q}^t(s, \hat{\pi}^t(s))} + \boxed{\hat{Q}^t(s, \pi^\star(s)) - Q^\star(s, \pi^\star(s))} \\
 &\geq \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(V^{\hat{\pi}^t}(s') - V^\star(s') \right) - 2\gamma^t \|\hat{Q}^0 - Q^\star\|_\infty
 \end{aligned}$$

Recursion

Final Quality of the Policy:

$$\hat{\pi}^t : \hat{\pi}^t(s) = \arg \max_a \widehat{Q}^t(s, a)$$

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

Proof:

$$\begin{aligned} V^{\hat{\pi}^t}(s) - V^\star(s) &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\ &= Q^{\hat{\pi}^t}(s, \hat{\pi}^t(s)) - Q^\star(s, \hat{\pi}^t(s)) + Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(V^{\hat{\pi}^t}(s') - V^\star(s') \right) + Q^\star(s, \hat{\pi}^t(s)) - Q^\star(s, \pi^\star(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(V^{\hat{\pi}^t}(s') - V^\star(s') \right) + Q^\star(s, \hat{\pi}^t(s)) - \widehat{Q}^t(s, \hat{\pi}^t(s)) + \widehat{Q}^t(s, \pi^\star(s)) - Q^\star(s, \pi^\star(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}^t(s))} \left(V^{\hat{\pi}^t}(s') - V^\star(s') \right) - 2\gamma^t \|\widehat{Q}^0 - Q^\star\|_\infty \dots \text{Recursion} \end{aligned}$$

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto \underline{\pi}(A)$

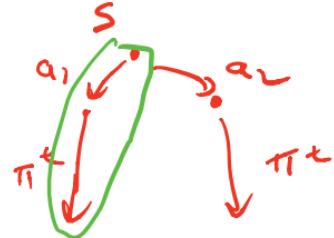
(focus on
Deterministic policies)

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto \Delta(A)$
2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto \Delta(A)$
2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$
3. Policy Improvement $\underline{\pi^{t+1}(s)} = \arg \max_a \underline{Q^{\pi^t}(s, a)}, \forall s$



Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto \Delta(A)$
 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: $\pi^0 : S \mapsto \Delta(A)$
 2. Policy Evaluation: $\underline{Q^{\pi^t}(s, a)}, \forall s, a$
 3. Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$
- 

Analysis of Policy Iteration

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

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

Analysis of Policy Iteration

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

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

$\forall s, a \in S \times A$

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

\uparrow
Bell-Eqⁿ \uparrow
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 $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned} Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &\stackrel{4}{=} \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \end{aligned}$$

Analysis of Policy Iteration

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

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

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

Analysis of Policy Iteration

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

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

$$\begin{aligned} \underline{Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a)} &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[\cancel{Q^{\pi^{t+1}}(s', \pi^{t+1}(s'))} - \cancel{Q^{\pi^t}(s', \pi^{t+1}(s'))} \right] \cancel{+} \dots \cancel{+} \cancel{\gamma^\infty} = 0 \end{aligned}$$

Analysis of Policy Iteration

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

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

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

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

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

$$\geq \max_a Q^{\pi^t}(s, a) \geq Q^{\pi^t}(s, \pi^t(s)) = V^{\pi^t}(s)$$

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

Analysis of Policy Iteration

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

Theorem: Convergence $\|\underline{V}^{\pi^{t+1}} - \underline{V}^{\star}\|_{\infty} \leq \gamma \|\underline{V}^{\pi^t} - \underline{V}^{\star}\|_{\infty}$ --- $\gamma^t \|V^{\pi^0} - V^{\star}\|_{\infty}$
 $t \rightarrow \infty$
 $\delta^t \rightarrow 0$

Analysis of Policy Iteration

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

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

$\forall s \in S$

Bell-opt

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

Analysis of Policy Iteration

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

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

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

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

Analysis of Policy Iteration

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

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

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

Analysis of Policy Iteration

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

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

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

Analysis of Policy Iteration

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

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

$$\begin{aligned} V^\star(s) - V^{\pi^{t+1}}(s) &= \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^\star(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^\star(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \\ &= \max_a (r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^\star(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s')) \\ &\leq \max_a \left(\cancel{r(s, a)} + \underbrace{\gamma \mathbb{E}_{s' \sim P(s, a)} V^\star(s')}_{\text{red}} - \left(\cancel{r(s, a)} + \underbrace{\gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s')}_{\text{red}} \right) \right) \\ &\leq \underbrace{\gamma \|V^\star - V^{\pi^t}\|_\infty}_{\text{red}} \quad \dots \quad \gamma^t \|V^\star - V^{\pi^0}\|_\infty \end{aligned}$$

$$\mathcal{T} \hat{Q} \rightarrow \hat{Q}$$

Value Iteration vs Policy Iteration

$$\hat{Q}^\infty = Q^*$$

$\pi^0: S \rightarrow A$
 π^{t+1} is better than π^t
 A^S many policies

Which one is faster? ✓

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

A key Lemma: Performance Difference Lemma

forall s.

Q: What's the performance difference of two policies π, π' , i.e., $\underline{V^\pi(s)} - \underline{V^{\pi'}(s)}$

A key Lemma: Performance Difference Lemma

$s_0 \in S$

PDL Statement:

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

\uparrow
from π

$\approx A^{\pi'}(s, a)$

$s_0 \xrightarrow{a}$

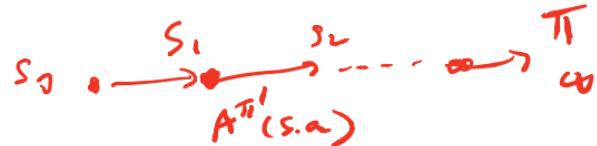
$\pi' \rightarrow \pi'$

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

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

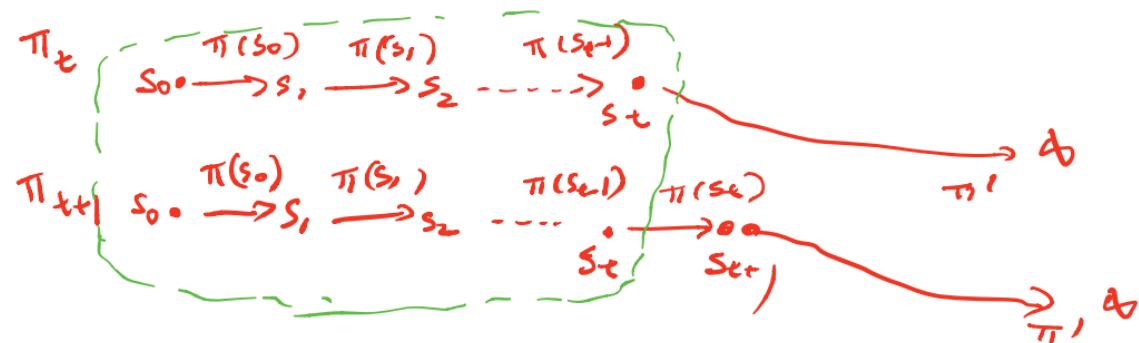
Proof of PDL

$$\underbrace{V^\pi(s_0) - V^{\pi'}(s_0)}_{\text{Red}} = \frac{1}{1-\gamma} \mathbb{E}_{\underbrace{s,a \sim d_{s_0}^\pi}_{\text{Red}}} [A^{\pi'}(s, a)]$$



Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

$$\pi_0 = \underline{\{\pi', \dots, \pi'\}}, \pi_1 = \underline{\{\pi, \pi', \dots, \pi'\}}, \pi_2 = \underline{\{\pi, \pi, \pi', \dots, \pi'\}}, \dots \underline{\pi_\infty = \{\pi, \dots, \pi\}}$$



Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

$$\underbrace{\pi_0 = \{\pi', \dots, \pi'\}, \pi_1 = \{\pi, \pi', \dots, \pi'\}, \pi_2 = \{\pi, \pi, \pi', \dots, \pi'\}, \dots}_{\pi_\infty = \{\pi, \dots, \pi\}}$$

$$\underbrace{V^\pi(s_0)}_{= V^{\pi_\infty}(s_0)}, \quad \underbrace{V^{\pi'}(s_0)}_{= V^{\pi_0}(s_0)}$$

Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

$$\pi_0 = \{\pi', \dots, \pi'\}, \pi_1 = \{\pi, \pi', \dots, \pi'\}, \pi_2 = \{\pi, \pi, \pi', \dots, \pi'\}, \dots \pi_\infty = \{\pi, \dots, \pi\}$$

$$V^\pi(s_0) = V^{\pi_\infty}(s_0), \quad V^{\pi'}(s_0) = V^{\pi_0}(s_0)$$

✓ ✓

$$V^{\pi_{n+1}}(s_0) - V^{\pi_n}(s_0)$$

Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

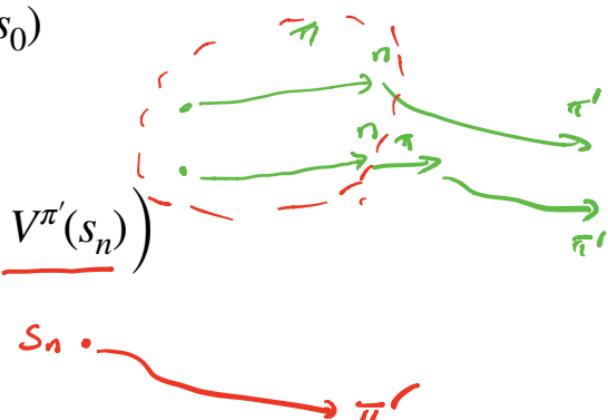
$$\pi_0 = \{\pi', \dots, \pi'\}, \pi_1 = \{\pi, \pi', \dots, \pi'\}, \pi_2 = \{\pi, \pi, \pi', \dots, \pi'\}, \dots \pi_\infty = \{\pi, \dots, \pi\}$$

$$V^\pi(s_0) = V^{\pi_\infty}(s_0), \quad V^{\pi'}(s_0) = V^{\pi_0}(s_0)$$

$$V^{\pi_{n+1}}(s_0) - V^{\pi_n}(s_0)$$

$$= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(r(s_n, \pi(s_n)) + \gamma \mathbb{E}_{s_{n+1} \sim P_{s_n, \pi(s_n)}} V^{\pi'}(s_{n+1}) - V^{\pi'}(s_n) \right)$$

$$s_n \xrightarrow{\pi} \pi'$$



Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

$$\pi_0 = \{\pi', \dots, \pi'\}, \pi_1 = \{\pi, \pi', \dots, \pi'\}, \pi_2 = \{\pi, \pi, \pi', \dots, \pi'\}, \dots \pi_\infty = \{\pi, \dots, \pi\}$$

$$V^\pi(s_0) = V^{\pi_\infty}(s_0), \quad V^{\pi'}(s_0) = V^{\pi_0}(s_0)$$

$$\begin{aligned}
 & V^{\pi_{n+1}}(s_0) - V^{\pi_n}(s_0) \\
 &= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(\underbrace{r(s_n, \pi(s_n)) + \gamma \mathbb{E}_{s_{n+1} \sim P_{s_n, \pi(s_n)}} V^{\pi'}(s_{n+1}) - V^{\pi'}(s_n)}_{\text{red underline}} \right) \\
 &= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(\underbrace{Q^{\pi'}(s_n, \pi(s_n)) - V^{\pi'}(s_n)}_{\text{red underline}} \right)
 \end{aligned}$$

Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

$$\pi_0 = \{\pi', \dots, \pi'\}, \pi_1 = \{\pi, \pi', \dots, \pi'\}, \pi_2 = \{\pi, \pi, \pi', \dots, \pi'\}, \dots \pi_\infty = \{\pi, \dots, \pi\}$$

$$V^\pi(s_0) = V^{\pi_\infty}(s_0), \quad V^{\pi'}(s_0) = V^{\pi_0}(s_0)$$

$$\begin{aligned}
 & V^{\pi_{n+1}}(s_0) - V^{\pi_n}(s_0) \\
 &= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(r(s_n, \pi(s_n)) + \gamma \mathbb{E}_{s_{n+1} \sim P_{s_n, \pi(s_n)}} V^{\pi'}(s_{n+1}) - V^{\pi'}(s_n) \right) \\
 &= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(Q^{\pi'}(s_n, \pi(s_n)) - V^{\pi'}(s_n) \right) \\
 &= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(\underset{\Delta}{\cancel{A}}^{\pi'}(s_n, \pi(s_n)) \right)
 \end{aligned}$$

Proof of PDL $V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} [A^{\pi'}(s, a)]$

$$\pi_0 = \{\pi', \dots, \pi'\}, \pi_1 = \{\pi, \pi', \dots, \pi'\}, \pi_2 = \{\pi, \pi, \pi', \dots, \pi'\}, \dots \pi_\infty = \{\pi, \dots, \pi\}$$

$$V^\pi(s_0) = V^{\pi_\infty}(s_0), \quad V^{\pi'}(s_0) = V^{\pi_0}(s_0)$$

$$V^{\pi_{n+1}}(s_0) - V^{\pi_n}(s_0)$$

$$= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(r(s_n, \pi(s_n)) + \gamma \mathbb{E}_{s_{n+1} \sim P_{s_n, \pi(s_n)}} V^{\pi'}(s_{n+1}) - V^{\pi'}(s_n) \right)$$

$$= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} (Q^{\pi'}(s_n, \pi(s_n)) - V^{\pi'}(s_n))$$

$$= \gamma^n \mathbb{E}_{s_n \sim \mathbb{P}_n(\cdot; s_0, \pi)} \left(A^{\pi'}(s_n, \pi(s_n)) \right)$$

Telescoping: $V^{\pi'} - V^\pi = \sum_{n=0}^{\infty} V^{\pi_n} - V^{\pi_{n+1}}$

$$\sqrt{V^{T_{10}}(s_0)} \rightarrow \sqrt{V^{T_{100}}(s_0)}$$

$$\begin{aligned} & \checkmark \pi_0 = \checkmark \pi_1 \swarrow \\ & = \pm \checkmark \pi_1 = -\checkmark \pi_2 \\ 1 & \quad \pm \checkmark \pi_2 = \checkmark \pi_3 \\ & = \checkmark \pi_0 = \checkmark \pi_{10} \quad i \end{aligned}$$