

Conservative Policy Iteration

Recap

Recall Policy Iteration (PI) with known (P, r)

Recap

Recall Policy Iteration (PI) with known (P, r)

Assume MDP is known, we compute $A^{\pi_{old}}(s, a)$ exactly for all s, a , PI updates policy as:

Recap

Recall Policy Iteration (PI) with known (P, r)

Assume MDP is known, we compute $A^{\pi_{old}}(s, a)$ exactly for all s, a , PI updates policy as:

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

Recap

Recall Policy Iteration (PI) with known (P, r)

Assume MDP is known, we compute $A^{\pi_{old}}(s, a)$ exactly for all s, a , PI updates policy as:

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

i.e., pick an action that has the largest advantage against π_{old} at every state s ,

Recap

Recall Policy Iteration (PI) with known (P, r)

Assume MDP is known, we compute $A^{\pi_{old}}(s, a)$ exactly for all s, a , PI updates policy as:

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

i.e., pick an action that has the largest advantage against π_{old} at every state s ,

Maximize advantage is great, as it gives monotonic improvement:

$$Q^{\pi_{new}}(s, a) \geq Q^{\pi_{old}}(s, a), \forall s, a$$

Recap

Recall Policy Iteration (PI):

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

Recap

Recall Policy Iteration (PI):

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

Performance Difference Lemma (PDL): for all $s_0 \in S$

Recap

Recall Policy Iteration (PI):

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

Performance Difference Lemma (PDL): for all $s_0 \in S$

$$V^{\pi_{new}}(s_0) - V^{\pi_{old}}(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d_{s_0}^{\pi_{new}}} [A^{\pi_{old}}(s, a)]$$

Recap

Recall Policy Iteration (PI):

$$\pi_{new}(s) = \arg \max_a A^{\pi_{old}}(s, a)$$

Performance Difference Lemma (PDL): for all $s_0 \in S$

$$V^{\pi_{new}}(s_0) - V^{\pi_{old}}(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d_{s_0}^{\pi_{new}}} [A^{\pi_{old}}(s, a)]$$

The advantage against π_{old} averaged over π_{new} 's own distribution

Today:
Conservative Policy Iteration

Q: How to enforce incremental policy update and ensure monotonic improvement

Outline

1. Greedy Policy Selection (via reduction to regression) and recap of API
2. Conservative Policy Iteration
3. Monotonic Improvement of CPI

Setting and Notation

Discounted infinite horizon MDP:

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

initial $s_0 \sim \mu$

Setting and Notation

Discounted infinite horizon MDP:

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

State visitation: $d_{\mu}^{\pi}(s) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s; \mu)$

Setting and Notation

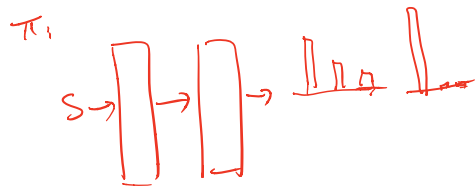
Discounted infinite horizon MDP:

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

$$\text{State visitation: } d_{\mu}^{\pi}(s) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s; \mu)$$

As we will consider large scale unknown MDP here, we start with a (restricted) function class Π :

$$\Pi = \{\pi : S \mapsto A\}$$



Setting and Notation

Discounted infinite horizon MDP:

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

$$\text{State visitation: } d_{\mu}^{\pi}(s) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s; \mu)$$

As we will consider large scale unknown MDP here, we start with a (restricted) function class Π :

$$\Pi = \{\pi : S \mapsto A\}$$

(We can think about each policy as a classifier)

Note that the optimal policy π^{\star} may not be in Π

Setting and Notation

Discounted infinite horizon MDP:

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

$$\text{State visitation: } d_{\mu}^{\pi}(s) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^{\pi}(s; \mu)$$

As we will consider large scale unknown MDP here, we start with a (restricted) function class Π :

$$\Pi = \{\pi : S \mapsto A\}$$

From now on, think about
deterministic policy as a special
stochastic policy

(We can think about each policy as a classifier)
Note that the optimal policy π^{\star} may not be in Π

Recall Approximate Policy Iteration (API)

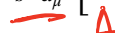
Recall Approximate Policy Iteration (API)

Given the current policy π^t , let's find a new policy that has large local adv over π^t under $d_{\mu}^{\pi^t}$

Recall Approximate Policy Iteration (API)

Given the current policy π^t , let's find a new policy that has large local adv over π^t under $d_\mu^{\pi^t}$

i.e., let's aim to (approximately) solve the following program:

$$\arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$


Recall Approximate Policy Iteration (API)

Given the current policy π^t , let's find a new policy that has large local adv over π^t under $d_\mu^{\pi^t}$

i.e., let's aim to (approximately) solve the following program:

$$\arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

Greedy Policy Selector

Recall Approximate Policy Iteration (API)

Given the current policy π^t , let's find a new policy that has large local adv over π^t under $d_\mu^{\pi^t}$

i.e., let's aim to (approximately) solve the following program:

$$\arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right] \quad \text{Greedy Policy Selector}$$

How to implement such greedy policy selector?
We talked about a regression process..

Implementing Approximate Greedy Policy Selector via Regression

We can do a **reduction to Regression** via Advantage function approximation

$$\mathcal{F} = \{f: S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi'}) \rightsquigarrow A^{\pi'}(s) = Q^{\pi'}(s, a) - V^{\pi'}(s)$$

Implementing Approximate Greedy Policy Selector via Regression

We can do a **reduction to Regression** via Advantage function approximation

$$\mathcal{F} = \{f: S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi'}) \quad \leftarrow \text{Given}$$

$$\Pi = \{\pi(s) = \arg \max_a f(s, a) : f \in \mathcal{F}\} \quad \leftarrow \text{implicit}$$

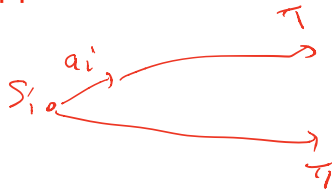
Implementing Approximate Greedy Policy Selector via Regression

We can do a **reduction to Regression** via Advantage function approximation

$$\mathcal{F} = \{f: S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi'})$$

$$\Pi = \{\pi(s) = \arg \max_a f(s, a) : f \in \mathcal{F}\}$$

$$\{s_i, a_i, y_i\}, s_i \sim d_{\mu}^{\pi'}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi'}(s_i, a_i)$$



Implementing Approximate Greedy Policy Selector via Regression

We can do a **reduction to Regression** via Advantage function approximation

$$\mathcal{F} = \{f: S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi^t})$$

$$\Pi = \{\pi(s) = \arg \max_a f(s, a) : f \in \mathcal{F}\}$$

$$\{s_i, a_i, y_i\}, s_i \sim d_\mu^{\pi^t}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi^t}(s_i, a_i)$$

Regression oracle:

$$\hat{A}^t = \arg \min_{f \in \mathcal{F}} \sum_i (f(s_i, a_i) - y_i)^2$$

Handwritten red annotations:
A red arrow points from the handwritten $\approx A^{\pi^t}$ to the \mathcal{F} in the regression equation.
A red arrow points from the circled $\mathbb{E}[y_i] = A^{\pi^t}(s_i, a_i)$ to the y_i in the regression equation.

Implementing Approximate Greedy Policy Selector via Regression

We can do a **reduction to Regression** via Advantage function approximation

$$\mathcal{F} = \{f : S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi^t})$$

$$\Pi = \{\pi(s) = \arg \max_a f(s, a) : f \in \mathcal{F}\}$$

$$\{s_i, a_i, y_i\}, s_i \sim d_{\mu}^{\pi^t}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi^t}(s_i, a_i)$$

Regression oracle:

$$\widehat{A}^t = \arg \min_{f \in \mathcal{F}} \sum_i (f(s_i, a_i) - y_i)^2$$

Act greedily wrt the estimator \widehat{A}^t (as we hope $\widehat{A}^t \approx A^{\pi^t}$):

Implementing Approximate Greedy Policy Selector via Regression

We can do a **reduction to Regression** via Advantage function approximation

$$\mathcal{F} = \{f: S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi'})$$

$$\Pi = \{\pi(s) = \arg \max_a f(s, a) : f \in \mathcal{F}\}$$

$$\{s_i, a_i, y_i\}, s_i \sim d_{\mu}^{\pi'}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi'}(s_i, a_i)$$

$$\begin{aligned} & \underset{a}{\operatorname{argmax}} A^{\pi}(s, a) \\ & = \underset{a}{\operatorname{argmax}} Q^{\pi}(s, a) \end{aligned}$$

Regression oracle:

$$\widehat{A}^t = \arg \min_{f \in \mathcal{F}} \sum_i (f(s_i, a_i) - y_i)^2$$

Act greedily wrt the estimator \widehat{A}^t (as we hope $\widehat{A}^t \approx A^{\pi'}$):

$$\widehat{\pi}(s) = \arg \max_a \widehat{A}^t(s, a), \forall s$$

Successful Regression ensures approximate greedy operator

$$\{s_i, a_i, y_i\}, s_i \sim d_\mu^{\pi^t}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi^t}(s_i, a_i)$$

Regression oracle:

$$\widehat{A}^t = \arg \min_{f \in \mathcal{F}} \sum_i (f(s_i, a_i) - y_i)^2$$

Successful Regression ensures approximate greedy operator

$$\{s_i, a_i, y_i\}, s_i \sim d_\mu^{\pi^t}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi^t}(s_i, a_i)$$

Regression oracle:

$$\widehat{A}^t = \arg \min_{f \in \mathcal{F}} \sum_i (f(s_i, a_i) - y_i)^2$$

Assume this regression is successful, i.e.,

$$\mathbb{E}_{s \sim d_\mu^{\pi^t}, a \sim U(A)} \left(\widehat{A}^t(s, a) - A^{\pi^t}(s, a) \right)^2 \leq \delta$$

$$\sqrt{\frac{1}{N}} \text{ or } \frac{1}{N}$$

Successful Regression ensures approximate greedy operator

$$\{s_i, a_i, y_i\}, s_i \sim d_\mu^{\pi^t}, a_i \sim U(A), \mathbb{E}[y_i] = A^{\pi^t}(s_i, a_i)$$

Regression oracle:

$$\widehat{A}^t = \arg \min_{f \in \mathcal{F}} \sum_i (f(s_i, a_i) - y_i)^2$$

Assume this regression is successful, i.e.,

$$\mathbb{E}_{s \sim d_\mu^{\pi^t}, a \sim U(A)} \left(\widehat{A}^t(s, a) - A^{\pi^t}(s, a) \right)^2 \leq \delta \leftarrow$$

Then, $\widehat{\pi}(s) = \arg \max_a \widehat{A}^t(s, a)$ is an approximate greedy policy:

$$\underbrace{\mathbb{E}_{s \sim d_\mu^{\pi^t}} A^{\pi^t}(s, \widehat{\pi}(s))}_{\Delta} \geq \underbrace{\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_\mu^{\pi^t}} A^{\pi^t}(s, \pi(s))}_{\Delta} - \underbrace{O(\sqrt{\delta})}_{\delta \approx \frac{1}{\sqrt{n}} \text{ (or } \# \text{ of)}} \rightarrow 0$$

$$\begin{aligned} & \mathbb{E}_{s \sim d_\mu^{\pi^t}} A^{\pi^t}(s, \widehat{\pi}(s)) - \mathbb{E}_{s \sim d_\mu^{\pi^t}} A^{\pi^t}(s, \widetilde{\pi}(s)) \\ &= \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[A^{\pi^t}(s, \widehat{\pi}(s)) - A^{\pi^t}(s, \widetilde{\pi}(s)) \right. \\ & \quad \left. + \widehat{A}^t(s, \widehat{\pi}(s)) - A^{\pi^t}(s, \widehat{\pi}(s)) \right. \\ & \quad \left. - A^{\pi^t}(s, \widetilde{\pi}(s)) \right] \\ & \geq \widehat{A}^t(s, \widehat{\pi}(s)) - A^{\pi^t}(s, \widetilde{\pi}(s)) \end{aligned}$$

Summary So Far:

By reduction to Supervised Learning (i.e., via Regression to approximate A^{π^t} under $d_{\mu}^{\pi^t}$), we can expect to find an approximate greedy optimizer $\hat{\pi}$, s.t.,

Summary So Far:

By reduction to Supervised Learning (i.e., via Regression to approximate A^{π^t} under $d_{\mu}^{\pi^t}$), we can expect to find an approximate greedy optimizer $\hat{\pi}$, s.t.,

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \hat{\pi}(s)) \right] \approx \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

PI:

$$\frac{A^{\pi^k}(s,a), \forall s,a}{\operatorname{argmax}_a A^{\pi^k}(s,a)}$$

Summary So Far:

API:

For $l=0, \dots, T$

$$\pi^{l+1} = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^l}} \left[A^{\pi^l}(s, \pi(s)) \right]$$

By reduction to Supervised Learning (i.e., via Regression to approximate A^{π^l} under $d_{\mu}^{\pi^l}$), we can expect to find an approximate greedy optimizer $\hat{\pi}$, s.t.,

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^l}} \left[A^{\pi^l}(s, \hat{\pi}(s)) \right] \approx \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^l}} \left[A^{\pi^l}(s, \pi(s)) \right]$$

from iteration l .

Throughout this lecture,
we will simply assume we can achieve $\operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^l}} \left[A^{\pi^l}(s, \pi(s)) \right]$

(think about

$N \rightarrow \infty$

Greedy policy selector

Outline



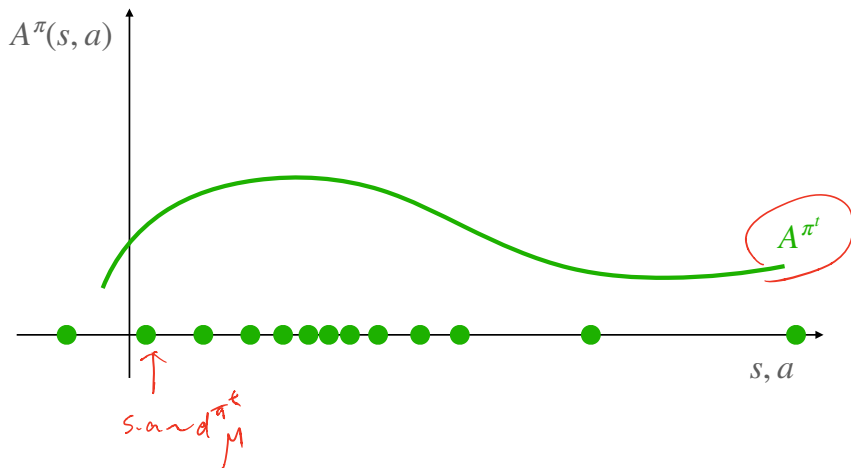
1. Greedy Policy Selection

2. Conservative Policy Iteration

3. Monotonic Improvement of CPI

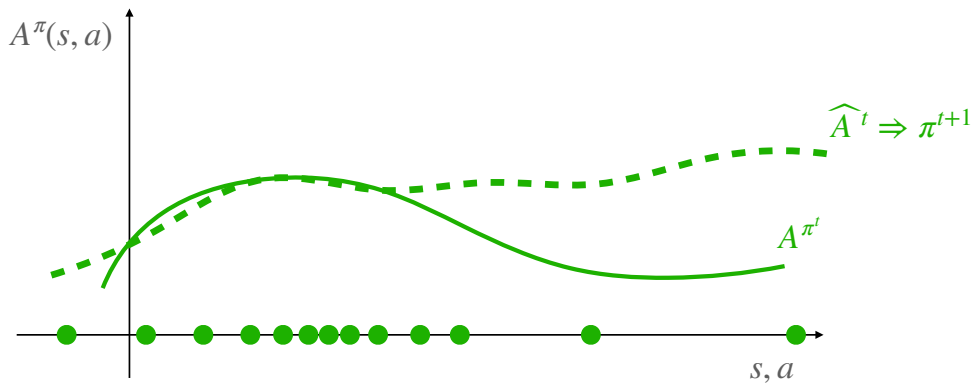
The Failure case of API: Abrupt distribution change

API cannot guarantee to succeed (let's think about advantage function approximation setting)



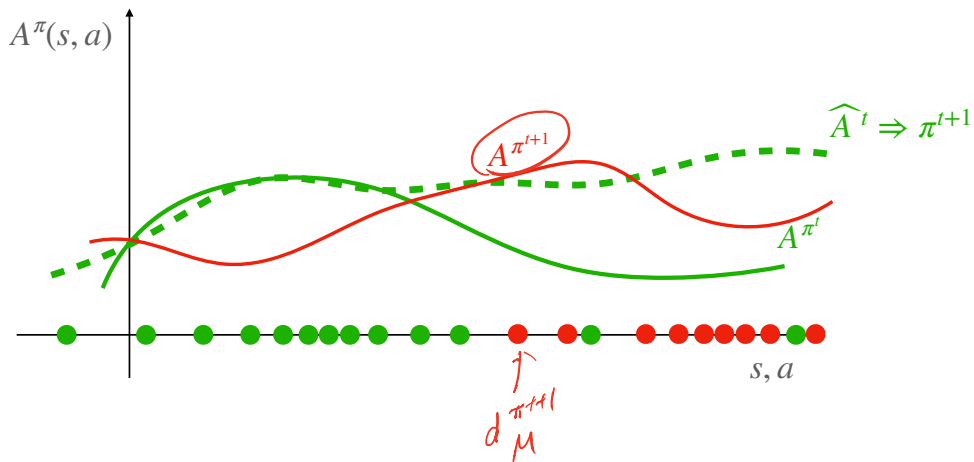
The Failure case of API: Abrupt distribution change

API cannot guarantee to succeed (let's think about advantage function approximation setting)



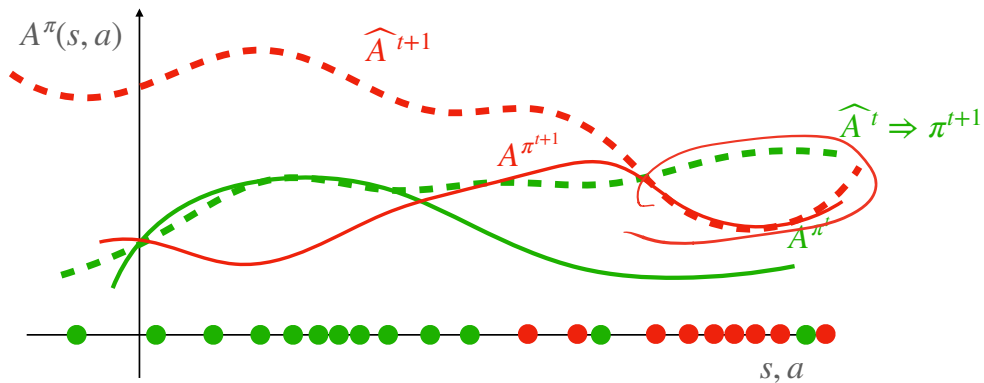
The Failure case of API: Abrupt distribution change

API cannot guarantee to succeed (let's think about advantage function approximation setting)



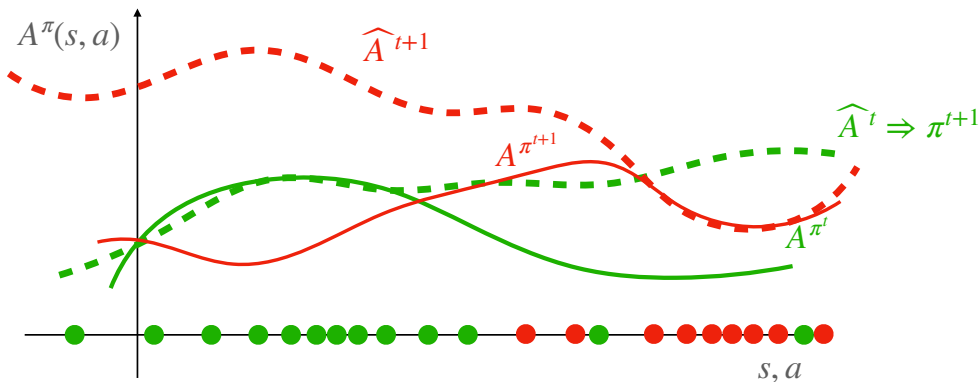
The Failure case of API: Abrupt distribution change

API cannot guarantee to succeed (let's think about advantage function approximation setting)



The Failure case of API: Abrupt distribution change

API cannot guarantee to succeed (let's think about advantage function approximation setting)



**Oscillation between two updates:
No monotonic improvement**

Key Idea of CPI: Incremental Update—No Abrupt Distribution Change

Let's design policy update rule such that $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ are not that different!

Key Idea of CPI: Incremental Update – No Abrupt Distribution Change

Let's design policy update rule such that $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ are not that different!

Recall Performance Difference Lemma:

$$V^{\pi^{t+1}} - V^{\pi^t} = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\underbrace{A^{\pi^t}(s, \pi^{t+1}(s))}_{\Delta} \right]$$

Key Idea of CPI: Incremental Update – No Abrupt Distribution Change

Let's design policy update rule such that $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ are not that different!

Recall Performance Difference Lemma:

$$V^{\pi^{t+1}} - V^{\pi^t} = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right]$$

Somehow: $d_{\mu}^{\pi^t} \approx d_{\mu}^{\pi^{t+1}}$

Key Idea of CPI: Incremental Update – No Abrupt Distribution Change

Let's design policy update rule such that $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ are not that different!

Recall Performance Difference Lemma:

$$V^{\pi^{t+1}} - V^{\pi^t} = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right]$$

$$d^{\pi^t} \approx d^{\pi^{t+1}}$$

$$\text{s.t.}, \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right] \approx \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right]$$

Key Idea of CPI: Incremental Update – No Abrupt Distribution Change

Let's design policy update rule such that $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ are not that different!

Recall Performance Difference Lemma:

$$V^{\pi^{t+1}} - V^{\pi^t} = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right]$$

$$d^{\pi^t} \approx d^{\pi^{t+1}} \leftarrow$$

$$\text{s.t.}, \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right] \approx \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right]$$

This we know how to optimize: the Greedy Policy Selector

(Recall Regression)

The CPI Algorithm

Initialize π^0

For $t = 0 \dots$

The CPI Algorithm

Initialize π^0

For $t = 0 \dots$

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi'}(s, \pi(s)) \right]$$

*local-Adv
against π^t*

$$\text{API: } \pi^{t+1} = \pi'$$

The CPI Algorithm

Initialize π^0

For $t = 0 \dots$

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \varepsilon$

Return π^t

The CPI Algorithm

Initialize π^0

For $t = 0 \dots$

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))]$$

2. If ~~$\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \epsilon$~~

~~Return π^t~~

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

with prob $1-\alpha$, $a \sim \pi^t(\cdot | s)$

w/ prob α , $a \sim \pi'(\cdot | s)$

$\alpha = 1 \rightarrow \text{API}$

$\alpha = 0 \rightarrow \pi^{t+1} = \pi^t$

$\pi^t \rightarrow \pi^{t+1}$
 $\pi' \leftarrow \text{Greedy policy selector}$

The CPI Algorithm

Initialize π^0

For $t = 0 \dots$

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))]$$

$$2. \text{ If } \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \varepsilon$$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\underline{\pi'}(\cdot | s), \forall s$$

$$\pi_{\theta} : s \rightarrow [a] \rightarrow [v] \rightarrow \dots$$

Q: Why this is incremental? In what sense?

Q: Can we get monotonic policy improvement?

Today: Policy Optimization



1. Greedy Policy Selection



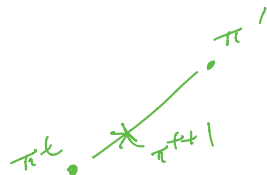
2. Conservative Policy Iteration

3. Monotonic Improvement of CPI

Q1: Why this is incremental? In what sense?

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

← Greedy policy selector



Q1: Why this is incremental? In what sense?

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Key observation 1:

$$\rightarrow \sum_a \left| \pi^{t+1}(a|s) - \pi^t(a|s) \right|$$

For any state s , we have $\|\pi^{t+1}(\cdot | s) - \pi^t(\cdot | s)\|_1 \leq 2\alpha$

$$\begin{aligned} & \left\| \pi^{t+1}(\cdot | s) - \pi^t(\cdot | s) \right\|_1 \\ &= \left\| -\alpha\pi^t(\cdot | s) + \alpha\pi'(\cdot | s) \right\|_1 \\ &\leq \alpha \left\| \pi^t(\cdot | s) \right\|_1 + \alpha \left\| \pi'(\cdot | s) \right\|_1 \\ &= \alpha + \alpha = 2\alpha \end{aligned}$$

Q1: Why this is incremental? In what sense?

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Key observation 1:

For any state s , we have $\|\pi^{t+1}(\cdot | s) - \pi^t(\cdot | s)\|_1 \leq 2\alpha$

Key observation 2 (Lemma 12.1 in AJKS)

For any two policies π and π' , if $\|\pi(\cdot | s) - \pi'(\cdot | s)\|_1 \leq \delta, \forall s$, then $\|d_{\mu}^{\pi}(\cdot) - d_{\mu}^{\pi'}(\cdot)\|_1 \leq \frac{\gamma\delta}{1-\gamma}$

(Handwritten notes: green arrows point from δ to d_{μ}^{π} and $d_{\mu}^{\pi'}$; a green underline is under $\|\pi(\cdot | s) - \pi'(\cdot | s)\|_1$; the fraction $\frac{\gamma\delta}{1-\gamma}$ is circled in green; the text "State-Distribution" is written below the fraction.)

Q1: Why this is incremental? In what sense?

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Key observation 1:

For any state s , we have $\|\pi^{t+1}(\cdot | s) - \pi^t(\cdot | s)\|_1 \leq 2\alpha$

Key observation 2 (Lemma 12.1 in AJKS)

For any two policies π and π' , if $\|\pi(\cdot | s) - \pi'(\cdot | s)\|_1 \leq \delta, \forall s$, then $\|d_\mu^\pi(\cdot) - d_\mu^{\pi'}(\cdot)\|_1 \leq \frac{\gamma\delta}{1-\gamma}$

CPI ensures incremental update, i.e., $\|d_\mu^{\pi^{t+1}}(\cdot) - d_\mu^{\pi^t}(\cdot)\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg local advantage:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \varepsilon$$

Δ

Max local Adv against π^t

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \varepsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \varepsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \varepsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \varepsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$(1 - \gamma) \left(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \right) = \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot | s)} \underbrace{A^{\pi^t}(s, a)} \right]$$

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi}(s, \pi(s))] \leq \varepsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \varepsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1 - \gamma) \left(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \right) &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot | s)} A^{\pi^t}(s, a) \right] \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \end{aligned}$$

← def of π^{t+1}

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi}(s, \pi(s))] \leq \varepsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \epsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1 - \gamma) \left(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \right) &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot|s)} A^{\pi^t}(s, a) \right] \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \end{aligned}$$

Why?

plug in def of π^{t+1}

$$\underbrace{\sum_a (1 - \alpha) \pi^t(a|s) A^{\pi^t}(s, a)}_{=0} + \sum_a \alpha \pi^t(a|s) A^{\pi^t}(s, a)$$

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

$$\mathbb{E}_{a \sim \pi^t(\cdot|s)} A^{\pi^t}(s, a) = V^{\pi}(s) - V^{\pi}(s) = 0$$

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

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

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \epsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1 - \gamma) \left(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \right) &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot | s)} A^{\pi^t}(s, a) \right] \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \quad \text{Why?} \\ &= \underbrace{\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right]}_{\mathbb{A}} + \underbrace{\left(\mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \right)}_{\leq 2 \cdot (\text{Local-Avg})} \end{aligned}$$

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

$$\begin{aligned} &\leq \frac{1}{1 - \gamma} \\ &\leq 2 \left(\max_{s, a} |A^{\pi^t}(s, a)| \right) \|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \\ &= \left| \mathbb{E}_{x \sim p} f(x) - \mathbb{E}_{x \sim q} f(x) \right| \\ &\leq \max_x |f(x)| \|p - q\|_1 \end{aligned}$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \epsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1 - \gamma) \left(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \right) &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot | s)} A^{\pi^t}(s, a) \right] \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \quad \text{Why?} \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] + \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \\ &\geq \underbrace{\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right]}_{2 \cdot [\text{local-Adv}]} - \underbrace{\frac{\alpha}{1 - \gamma} \|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1}_{\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_2} \in \frac{\delta \cdot 2\alpha}{1 - \gamma} \end{aligned}$$

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:


$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \varepsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1 - \gamma) \left(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \right) &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot | s)} A^{\pi^t}(s, a) \right] \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \quad \text{Why?} \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] + \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \\ &\geq \underbrace{\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right]}_{= \mathbb{A}} - \frac{\alpha}{1 - \gamma} \|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \\ &\geq \underbrace{\alpha \mathbb{A}}_{\substack{\uparrow \\ \text{max-local-Adv}}} - \frac{2\gamma\alpha^2}{(1 - \gamma)^2} \end{aligned}$$


Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \varepsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

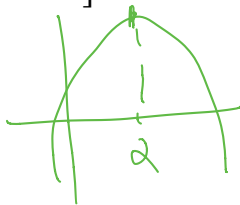
Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_\mu^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \varepsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1 - \gamma) \left(V_\mu^{\pi^{t+1}} - V_\mu^{\pi^t} \right) &= \mathbb{E}_{s \sim d_\mu^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot | s)} A^{\pi^t}(s, a) \right] \\ &= \mathbb{E}_{s \sim d_\mu^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \quad \text{Why?} \\ &= \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] + \mathbb{E}_{s \sim d_\mu^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \\ &\geq \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \frac{\alpha}{1 - \gamma} \|d_\mu^{\pi^{t+1}} - d_\mu^{\pi^t}\|_1 \\ &\geq \alpha \mathbb{A} - \frac{2\gamma\alpha^2}{(1 - \gamma)^2} \end{aligned}$$

$$\text{(Set } \alpha = \frac{(1 - \gamma)^2 \mathbb{A}}{4\gamma} \text{)}$$



Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_\mu^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_\mu^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \varepsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Monotonic Improvement before Termination:

Before terminate, we have non-trivial avg **local advantage**:

$$\mathbb{A} := \mathbb{E}_{d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi'(s)) \right] \geq \epsilon$$

Can we translate local advantage \mathbb{A} to $V^{\pi^{t+1}} - V^{\pi^t}$? (Yes, by PDL)

$$\begin{aligned} (1-\gamma)(V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t}) &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\mathbb{E}_{a \sim \pi^{t+1}(\cdot|s)} A^{\pi^t}(s, a) \right] \quad \leftarrow \text{PDL} \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \quad \leftarrow \text{Def of } \pi^{t+1} \\ &= \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] + \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] \\ &\geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\alpha A^{\pi^t}(s, \pi'(s)) \right] - \frac{\alpha}{1-\gamma} \|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \quad \leftarrow \text{Incremental-update} \\ &\geq \alpha \mathbb{A} - \frac{2\gamma\alpha^2}{(1-\gamma)^2} = \frac{\mathbb{A}^2(1-\gamma)}{8\gamma} > 0 \quad \left(\text{Set } \alpha = \frac{(1-\gamma)^2 \mathbb{A}}{4\gamma} \right) \end{aligned}$$

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1-\alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

$$\begin{aligned} V^{\pi^{t+1}} - V^{\pi^t} &\geq \frac{1}{1-\gamma} \cdot \frac{\mathbb{A}^2(1-\gamma)}{8\gamma} \\ &\approx \frac{\mathbb{A}^2}{8} \quad \mathbb{A} \geq \epsilon \\ &\geq \frac{\epsilon^2}{8} \end{aligned}$$

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$



Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$

2. Before terminate, monotonic improvement (Thm 12.2 in AJKS):

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi'}} [A^{\pi'}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi'}} [A^{\pi'}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$

2. Before terminate, monotonic improvement (Thm 12.2 in AJKS):

(By setting step size α properly...)

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi'}} [A^{\pi'}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi'}} [A^{\pi'}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$

2. Before terminate, monotonic improvement (Thm 12.2 in AJKS):

$$V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \geq \frac{\epsilon^2}{8\gamma}$$

(By setting step size α properly...)

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$

2. Before terminate, monotonic improvement (Thm 12.2 in AJKS):

$$V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \geq \frac{\epsilon^2}{8\gamma}$$

(By setting step size α properly...)

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

$$V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t}$$

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$

2. Before terminate, monotonic improvement (Thm 12.2 in AJKS):

$$V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \geq \frac{\epsilon^2}{8\gamma}$$

(By setting step size α properly...)

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))]$$

2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

$$\begin{aligned} & V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \\ & \geq \underbrace{\alpha \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi^t(s))]}_{\text{max-local-adv}} - \frac{\alpha}{1-\gamma} \underbrace{\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1} \end{aligned}$$

Summary of CPI so far:

1. Incremental update (Lemma 12.1 in AJKS)

$$\|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \leq \frac{2\gamma\alpha}{1-\gamma}$$

2. Before terminate, monotonic improvement (Thm 12.2 in AJKS):

$$V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \geq \frac{\epsilon^2}{8\gamma}$$

(By setting step size α properly...)

Recall CPI:

1. Greedy Policy Selector:

$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))]$$


2. If $\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi(s))] \leq \epsilon$

Return π^t

3. Incremental Update:

$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Local adv versus distribution change:

$$\begin{aligned} & V_{\mu}^{\pi^{t+1}} - V_{\mu}^{\pi^t} \\ & \geq \alpha \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi'}(s, \pi^t(s))] - \frac{\alpha}{1-\gamma} \|d_{\mu}^{\pi^{t+1}} - d_{\mu}^{\pi^t}\|_1 \\ & \geq \alpha\epsilon - \frac{\gamma\alpha^2}{(1-\gamma)^2} \end{aligned}$$


Finite Horizon H steps

$$\|P_n^{\pi^t} - P_n^{\pi^{t+1}}\|_1 \leq \delta \cdot h$$

Summary for today:

$$\underbrace{(0 + \delta + 2\delta + 3\delta \dots (H-1)\delta)}_{\delta \cdot H/2} / H \approx H \cdot \delta$$

1. Algorithm: Conservative Policy Iteration:
Find the local greedy policy, and move towards it a little bit

2. Small change in policies results small change in state distributions

2. Unlike API, incremental policy update ensures monotonic improvement

$$\|P_h^{\pi^t} - P_h^{\pi^{t+1}}\|_1$$

$$\approx h \cdot \delta$$

$$(1-\gamma) [0 + \delta + \delta^2 + \delta^3 \dots]$$

$$V^{\pi^{\text{new}}}(s_0) - V^{\pi^{\text{old}}}(s_0) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi^{\text{new}}}_{s_0}} \left[A^{\pi^{\text{old}}}(s, \pi^{\text{new}}(s)) \right]$$

$$\mathbb{E}_{s \sim p_0} \left[V^{\pi^{\text{new}}}(s) - V^{\pi^{\text{old}}}(s) \right] = \frac{1}{1-\gamma} \left(\mathbb{E}_{s \sim d^{\pi^{\text{new}}}_{s_0}} \right) \left[\dots \right]$$

↓ π^{new}
d μ

$$s' = \underset{\Delta}{f}(s, a) + \Sigma \rightarrow P(s' | s, a) \leftarrow \text{CPI (PG)}$$

$$\rightarrow \mathbb{Q}^{\pi}(s, a), \forall s, a ??$$

$$\min_{\pi} \sum_{h=0}^{H-1} C(s_h, a_h)$$

$$\text{s.t. } s_{h+1} = \underset{\Delta}{f}(s_h, a_h), a_h = \pi(s_h)$$

$$\pi^{t+1} = \left(\frac{1}{2}\right)\pi^t + 2\pi' = \pi^t - 2\pi^t + 2\pi'$$

$$\left\| \frac{\pi^{t+1} - \pi^t}{1} \right\|_1 = \left\| -2\pi^t + 2\pi' \right\|_1$$

$$\leq 2\|\pi^t\|_1 + 2\|\pi'\|_1$$

$$\leq 2 + 2$$

$$\|\pi^t(\cdot|s)\|_1 = 1$$

Finite Horizon H

$$\sqrt{\|P_n^{\pi^{t+1}}(\cdot) - P_n^{\pi^t}(\cdot)\|} \leq h \cdot \delta$$



$$\left(\overset{\delta}{\downarrow} 0 + \overset{\delta}{\downarrow} \delta + \overset{\delta}{\downarrow} 2\delta + \dots + \overset{\delta}{\downarrow} (H-1)\delta \right) / H$$

$$\approx \frac{\delta H(H-1)}{2} / H \approx \delta H$$

$$\pi^t(\cdot|s) \in \Delta(A)$$

