

HW₂ due ~~to~~ Friday

Approximate Policy Iteration & Conservative Policy Iteration

Recap

Recall Policy Iteration (PI):

Recap

Recall Policy Iteration (PI):

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

Recap

Recall Policy Iteration (PI):

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

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

Recap

Recall Policy Iteration (PI):

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

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

i.e., be greedy with respect to π at every state s ,

Recap

Recall Policy Iteration (PI):

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

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

i.e., be greedy with respect to π at every state s ,

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

Recap

Recall Policy Iteration (PI):

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

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

i.e., be greedy with respect to π at every state s ,

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

However, for large scale, unknown MDP
there is no way we will be able to know $A^\pi(s, a)$ at all s, a ,
so how can we do policy update?

Recap

Recall Policy Iteration (PI):

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

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

Recap

Recall Policy Iteration (PI):

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

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

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

Recap

Recall Policy Iteration (PI):

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

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

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

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

Recap

Recall Policy Iteration (PI):

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

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

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

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

$$= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi'}} \left[\underbrace{\max_{a \in A} A^\pi(s, a)} \right] \geq 0 \quad \checkmark$$

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

$\max_a A^\pi(s, a) \geq 0$

Recap

Recall Policy Iteration (PI):

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

Monotonic improvement of PI: $Q^{\pi'}(s, a) \geq Q^\pi(s, a), \forall s, a$

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

$$\begin{aligned} V^{\pi'}(s_0) - V^\pi(s_0) &= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi'}} [A^\pi(s, \pi'(s))] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi'}} \left[\max_{a \in A} A^\pi(s, a) \right] \geq 0 \end{aligned}$$

However, for large scale, unknown MDP
there is no way we will be able to know $A^\pi(s, a)$ at all s, a ,
so how can we do policy update?

Setting and Notation

Discounted infinite horizon MDP:

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

↓

Setting and Notation

Discounted infinite horizon MDP:

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

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

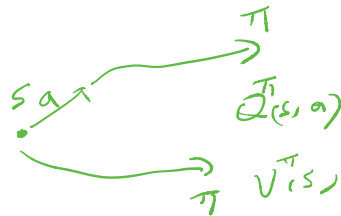
Setting and Notation

Discounted infinite horizon MDP:

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

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

$$\text{Unbiased estimate of } A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$$



Setting and Notation

Discounted infinite horizon MDP:

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

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

$$\text{Unbiased estimate of } A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$$

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

$$\Pi = \{\pi : S \mapsto \Delta(A)\}$$

Δ

Setting and Notation

Discounted infinite horizon MDP:

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

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

$$\text{Unbiased estimate of } A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$$

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

$$\Pi = \{\pi : S \mapsto \Delta(A)\}$$

Note that the optimal policy π^* may not be in Π

$$\pi^* \notin \Pi$$

Attempt One: Approximate Policy Iteration (API)

Attempt One: Approximate Policy Iteration (API)

Given the current policy π^t , let's act greedily wrt π under $d_{\mu}^{\pi^t}$

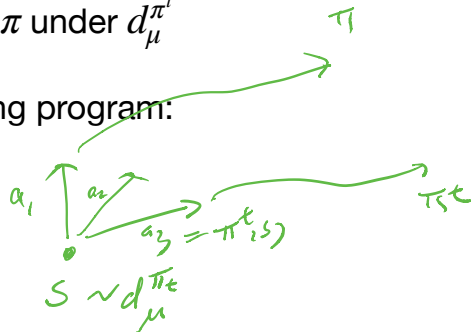
Attempt One: Approximate Policy Iteration (API)

Given the current policy π^t , let's act greedily wrt π under $d_\mu^{\pi^t}$

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

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

A



Attempt One: Approximate Policy Iteration (API)

Given the current policy π^t , let's act greedily wrt π under $d_\mu^{\pi^t}$

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

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

Attempt One: Approximate Policy Iteration (API)

Given the current policy π^t , let's act greedily wrt π under $d_\mu^{\pi^t}$

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

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

But we can only sample from $d_\mu^{\pi^t}$, and we can only get an approximation of $A^{\pi^t}(s, a)$

Attempt One: Approximate Policy Iteration (API)

Given the current policy π^t , let's act greedily wrt π under $d_\mu^{\pi^t}$

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

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

(Note: In the original image, a green arrow points from the expectation symbol to the distribution $d_\mu^{\pi^t}$, and a green circle highlights the distribution and the action $\pi(s)$ in the value function.)

But we can only sample from $d_\mu^{\pi^t}$, and we can only get an approximation of $A^{\pi^t}(s, a)$

(Note: In the original image, a green circle highlights $d_\mu^{\pi^t}$ and a green oval highlights $A^{\pi^t}(s, a)$.)

We can hope for an Approximate Greedy Policy Selector a reduction to Regression

Implementing Approximate Greedy Policy Selector via Regression

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

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

(Note: A green checkmark is drawn under the expression $(\approx A^{\pi'})$)

Implementing Approximate Greedy Policy Selector via Regression

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

$$\mathcal{F} = \{f: S \times A \mapsto \mathbb{R}\} \quad (\approx A^{\pi'})$$
$$\Pi = \{\pi(s) = \arg \max_a \underbrace{f(s, a)}\}$$

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

Implementing Approximate Greedy Policy Selector via Regression

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

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

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

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

$\Delta \quad \dot{\Delta}$ \uparrow unbiased from MC rollout

Implementing Approximate Greedy Policy Selector via Regression

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

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

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

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

Regression oracle:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_i \left(\underset{\Delta \Delta}{f(s_i, a_i)} - \underset{\Delta}{\widetilde{A}_i} \right)^2$$

$\rightarrow A^{\pi'}(s_i, a_i)$

Implementing Approximate Greedy Policy Selector via Regression

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

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

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

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

Regression oracle:

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

Act greedily wrt the estimator \hat{f} (as we hope $\hat{f} \approx A^{\pi'}$).

Implementing Approximate Greedy Policy Selector via Regression

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

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

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

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

Regression oracle:

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

Act greedily wrt the estimator \hat{f} (as we hope $\hat{f} \approx A^{\pi'}$):

$$\hat{\pi}(s) = \arg \max_a \hat{f}(s, a), \forall s$$

Implementing Approximate Greedy Policy Selector via Regression

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

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

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

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

Regression oracle:

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

Act greedily wrt the estimator \hat{f} (as we hope $\hat{f} \approx A^{\pi'}$):

$$\hat{\pi}(s) = \arg \max_a \hat{f}(s, a), \forall s$$

Do **finite sample analysis for Regression** first, and then transfer the guarantee to greedy policy selection

Analyzing Approximation error via Regression

Greedy Policy Selector

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

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

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

$$\hat{\pi}(s) = \arg \max_a \hat{f}(s, a), \forall s$$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \hat{\pi}(s))] = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [\hat{f}(s, \hat{\pi}(s)) + A^{\pi^t}(s, \hat{\pi}(s)) - \hat{f}(s, \hat{\pi}(s))]$$

$$\geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [\hat{f}(s, \tilde{\pi}(s)) + A^{\pi^t}(s, \hat{\pi}(s)) - \hat{f}(s, \hat{\pi}(s))]$$

$$\geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \tilde{\pi}(s)) + \hat{f}(s, \tilde{\pi}(s)) - A^{\pi^t}(s, \tilde{\pi}(s)) + A^{\pi^t}(s, \hat{\pi}(s)) - \hat{f}(s, \hat{\pi}(s))]$$

$$\geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \tilde{\pi}(s))] + \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [\hat{f}(s, \tilde{\pi}(s)) - A^{\pi^t}(s, \tilde{\pi}(s)) + A^{\pi^t}(s, \hat{\pi}(s)) - \hat{f}(s, \hat{\pi}(s))]$$

Summary So Far:

By reduction to Supervised Learning (i.e., classification using Π or Regression using \mathcal{F}), with high probability, we get:

Summary So Far:

By reduction to Supervised Learning (i.e., classification using Π or Regression using \mathcal{F}), with high probability, we get:

$$\mathbb{E}_{s \sim d_{\mu}^t} \left[A^{\pi^t}(s, \hat{\pi}(s)) \right] \geq \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^t} \left[A^{\pi^t}(s, \pi(s)) \right] - \underbrace{\frac{A}{1-\gamma} \sqrt{\frac{\ln(|\mathcal{F}|/\delta)}{N}}}_{\text{statistical error: } \epsilon}$$

N # of Training Samples

Summary So Far:

By reduction to Supervised Learning (i.e., classification using Π or Regression using \mathcal{F}), with high probability, we get:

$$\mathbb{E}_{s \sim d_{\mu}^t} \left[A^{\pi^t}(s, \hat{\pi}(s)) \right] \geq \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^t} \left[A^{\pi^t}(s, \pi(s)) \right] - \underbrace{\frac{A}{1-\gamma} \sqrt{\frac{\ln(|\mathcal{F}|/\delta)}{N}}}_{\text{statistical error: } \epsilon} \rightarrow \infty$$

In the rest of the lecture, as we will focus on convergence rather than sample complexity, we ignore the statistical error (goes to zero as N increases),


Summary So Far:

By reduction to Supervised Learning (i.e., classification using Π or Regression using \mathcal{F}), with high probability, we get:

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \hat{\pi}(s)) \right] \geq \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right] - \underbrace{\frac{A}{1-\gamma} \sqrt{\frac{\ln(|\mathcal{F}|/\delta)}{N}}}_{\text{statistical error: } \epsilon}$$

In the rest of the lecture, as we will focus on convergence rather than sample complexity, we ignore the statistical error (goes to zero as N increases),

i.e., we assume we can do the exact greedy policy selector: $\arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$



Algorithm: Approximate Policy Iteration (API)

Iterate: π_t

$$\text{API: } \pi^{t+1} \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s, a \sim d_{\mu}^{\pi}} [A^{\pi_t}(s, \pi(s))]$$

A



Reduce Regression

$N \rightarrow \infty$

Algorithm: Approximate Policy Iteration (API)

Iterate:

$$\text{API: } \pi^{t+1} \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s, a \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))]$$

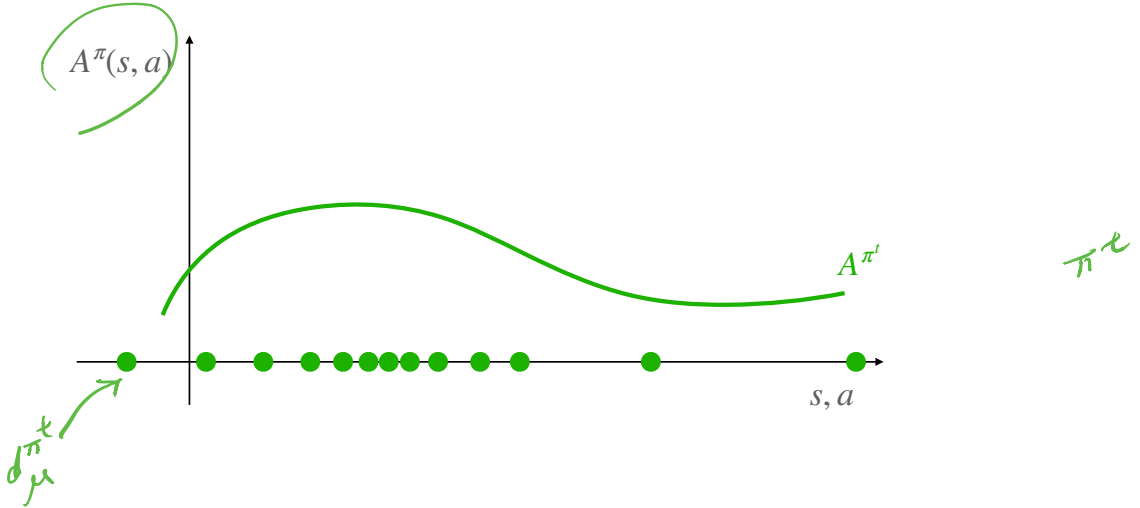
Question:

- (1) Does API has monotonic improvement?
- (2) Does it convergence?

The Oscillation of API from Abrupt Distribution Change

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

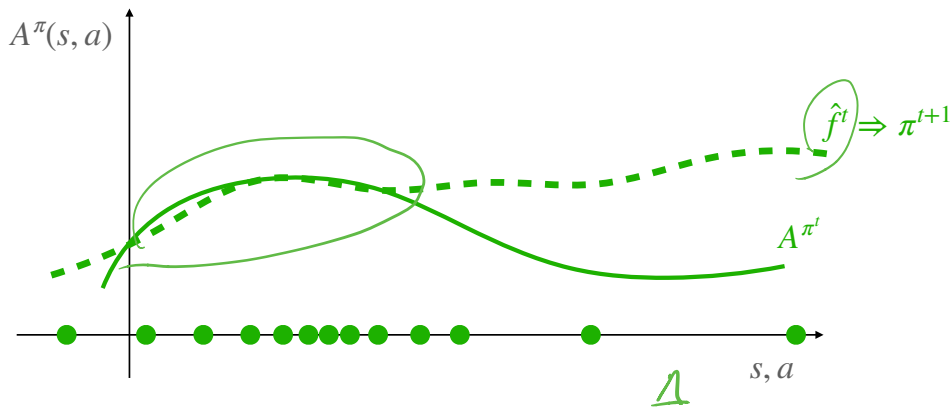
Concrete example in Chapter 3



The Oscillation of API from Abrupt Distribution Change

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

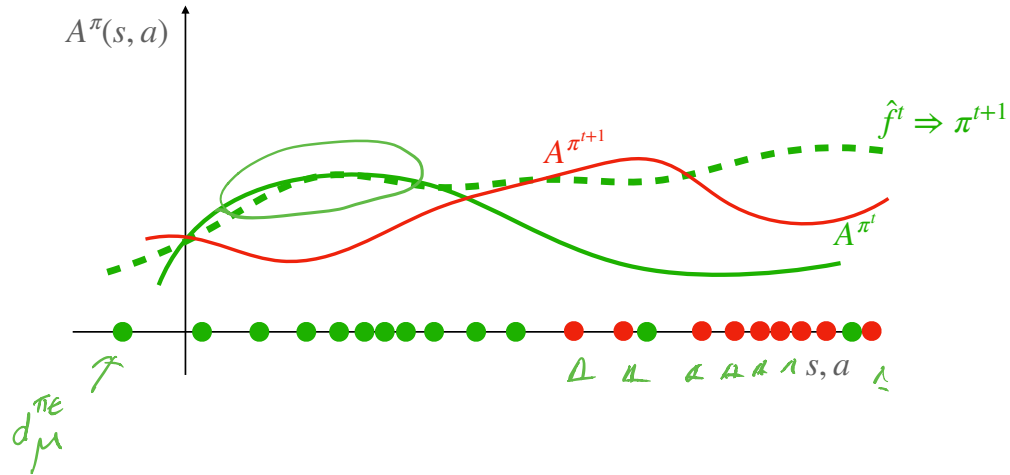
Concrete
example in
Chapter 3



The Oscillation of API from Abrupt Distribution Change

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

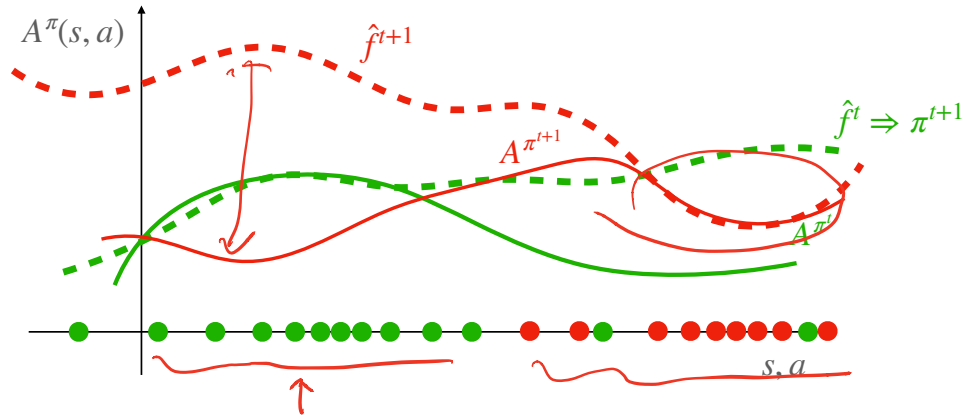
Concrete example in Chapter 3



The Oscillation of API from Abrupt Distribution Change

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

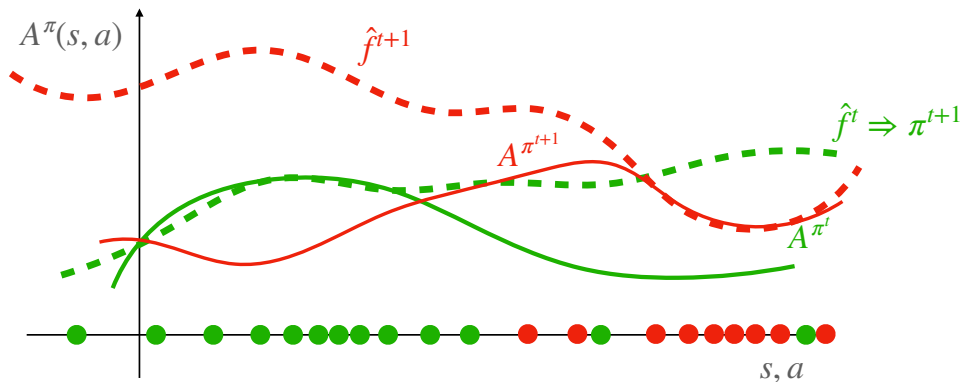
Concrete
example in
Chapter 3



The Oscillation of API from Abrupt Distribution Change

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

Concrete
example in
Chapter 3



Oscillation between two updates:
No monotonic improvement

Key Issue: Abrupt Policy Change, i.e., $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ could be widely different

Key Issue: Abrupt Policy Change, i.e., $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ could be widely different

Our estimator \hat{f}^t is only good under $d_{\mu}^{\pi^t}$, i.e.,

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}, a \sim U(A)} (\hat{f}^t(s, a) - A^{\pi^t}(s, a))^2 \text{ small,}$$

but $\mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}, a \sim U(A)} (\hat{f}^t(s, a) - A^{\pi^t}(s, a))^2$ might be arbitrarily big

$$d_{\mu}^{\pi^*} \neq d_{\mu}^{\pi^{t+1}}$$

Key Issue: Abrupt Policy Change, i.e., $d_{\mu}^{\pi^{t+1}}$ and $d_{\mu}^{\pi^t}$ could be widely different

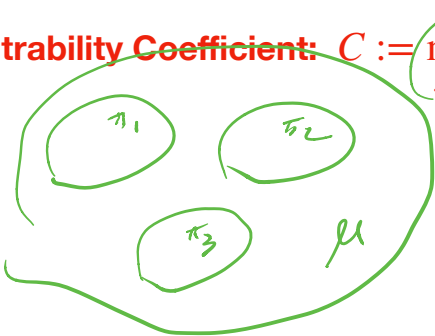
Our estimator \hat{f}^t is only good under $d_{\mu}^{\pi^t}$, i.e.,

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}, a \sim U(A)} (\hat{f}^t(s, a) - A^{\pi^t}(s, a))^2 \text{ small,}$$

but $\mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}, a \sim U(A)} (f^t(s, a) - A^{\pi^t}(s, a))^2$ might be arbitrarily big

To make API to make improvement, we need a much stronger coverage assumption:

A strong Concentrability Coefficient: $C := \max_{\pi \in \Pi} \sup_s \frac{d_{\mu}^{\pi}(s)}{\mu(s)} < \infty$



Key Issue: Abrupt Policy Change, i.e., $d_\mu^{\pi^{t+1}}$ and $d_\mu^{\pi^t}$ could be widely different

Our estimator \hat{f}^t is only good under $d_\mu^{\pi^t}$, i.e.,

$$\mathbb{E}_{s \sim d_\mu^{\pi^t}, a \sim U(A)} (\hat{f}^t(s, a) - A^{\pi^t}(s, a))^2 \text{ small,}$$

but $\mathbb{E}_{s \sim d_\mu^{\pi^{t+1}}, a \sim U(A)} (f^t(s, a) - A^{\pi^t}(s, a))^2$ might be arbitrarily big

To make API to make improvement, we need a much stronger coverage assumption:

A strong Concentrability Coefficient: $C := \max_{\pi \in \Pi} \sup_s \frac{d_\mu^\pi(s)}{\mu(s)} < \infty$

If $C < \infty$, i.e., μ covers all d_μ^π , then we can expect error

$\mathbb{E}_{s \sim d_\mu^{\pi^{t+1}}, a \sim U(A)} (\hat{f}^t(s, a) - A^{\pi^t}(s, a))^2$ is reasonably under control;

$$\textcircled{\text{O}} \quad \underline{d_\mu^\pi} \approx \underline{d_\mu^{\pi^{t+1}}}$$

$$\textcircled{d_\mu^{\pi^{t+1}}} \textcircled{\mu}$$

Conservative Policy Iteration – An Incremental Policy Optimization Approach

(And the benefit of being incremental)

NPG or Trust Region

$$KL(P_{\pi^{t+1}} || P_{\pi^t}) \leq \delta$$

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

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

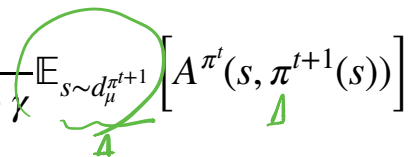
Δ

Δ

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

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

Recall Performance Difference Lemma:

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


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

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

Recall Performance Difference Lemma:

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

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

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

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

Recall Performance Difference Lemma:

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

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

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

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

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

Recall Performance Difference Lemma:

$$V^{\pi^{t+1}} - V^{\pi^t} = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right] \geq 0$$

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

$$\text{s.t.}, \mathbb{E}_{s \sim d^{\pi^t}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right] \approx \mathbb{E}_{s \sim d^{\pi^{t+1}}} \left[A^{\pi^t}(s, \pi^{t+1}(s)) \right] \geq 0$$

This we know how to optimize: the Greedy Policy Selector

→ Regression, $s \sim d^{\pi^t}$, $a \sim U(A)$, $A(s, a)$

CPI Algorithm:

CPI Algorithm:

1. Greedy Policy Selector:

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

CPI Algorithm:

1. Greedy Policy Selector:

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

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

↪ local optimal

CPI Algorithm:

1. Greedy Policy Selector:

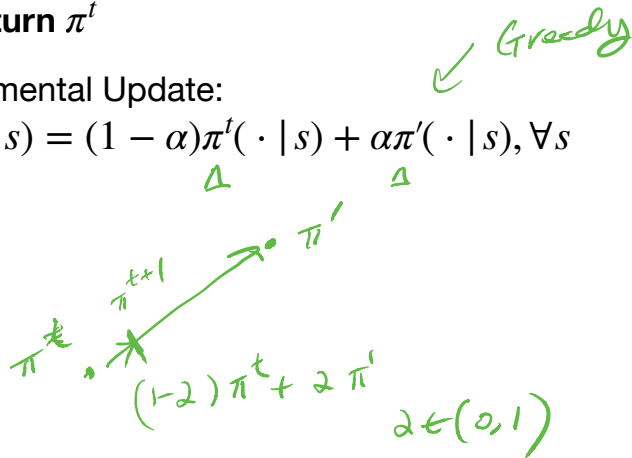
$$\pi' \in \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi'}} [A^{\pi'}(s, \pi(s))]$$

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

Return π^t

3. Incremental Update:

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



$$\pi^{t+1}(c_i | s) = \begin{cases} \pi^t(c_i | s) & \text{wp } 1-\alpha \\ \pi'(c_i | s) & \text{wp } \alpha \end{cases}$$

$\alpha \rightarrow 0$

CPI Algorithm:

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'

3. Incremental Update:

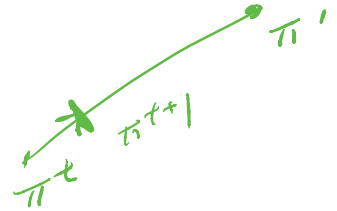
$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

Q: Why this is incremental? In what sense?

Q: Can we get monotonic policy improvement?

The incremental Nature of CPI:

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



The incremental Nature of CPI:

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

Key observation 1: $\forall s$

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

$$\begin{aligned} & \pi^{t+1}(c|s) - \pi^t(c|s) \\ &= \underbrace{\alpha}_{\text{small}} \left(\pi^t(c|s) - \pi^t(c|s) \right) \end{aligned}$$

The incremental Nature of CPI:

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

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:

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

The incremental Nature of CPI:

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

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:

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

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

Policy Improvement before Termination:

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

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

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:

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

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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:

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

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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:

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

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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

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

Return π^t

3. Incremental Update:

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

$$\pi^{t+1} \equiv \begin{cases} \pi^t & 1-\alpha \\ \pi' & \alpha \end{cases}$$

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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:

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

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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:

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

$$\begin{aligned} & \left| \mathbb{E}_P(f(x)) - \mathbb{E}_Q(f(x)) \right| \\ &= \sup_x |f(x)| \|P - Q\|_{TV} \end{aligned}$$

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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:

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

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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

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:

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

Policy 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 $V^{\pi^{t+1}} - V^{\pi^t}$?

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

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:

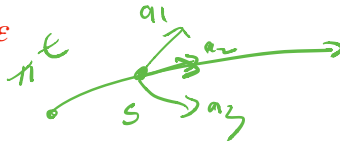
$$\pi^{t+1}(\cdot | s) = (1 - \alpha)\pi^t(\cdot | s) + \alpha\pi'(\cdot | s), \forall s$$

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

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

1. No more positive advantage by one-step deviation from π^t 's own states

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



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

Return π^t

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

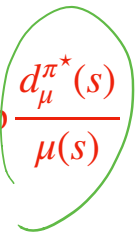
Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

Δ



Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

Recall Π is restricted, denote $\epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \underbrace{\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))]}_{\Delta}$

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

arg max _{π} $A^{\pi^t}(s, a)$
 $\notin \Pi$
 $\pi \nearrow \in \Pi \searrow$

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

Recall Π is restricted, denote $\epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^t(s)) \right] + \epsilon_{\Pi} \geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] = \mathbb{E}_{s \sim d^*} \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \max_a A^{\pi^t}(s, a)$$

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

Recall Π is restricted, denote $\epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^t(s)) \right] + \epsilon_{\Pi} \geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] = \mathbb{E}_{s \sim d^*} \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \max_a A^{\pi^t}(s, a)$$

$$\geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} \max_a A^{\pi^t}(s, a)$$

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

Recall Π is restricted, denote $\epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^t(s)) \right] + \epsilon_{\Pi} \geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] = \mathbb{E}_{s \sim d^*} \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \max_a A^{\pi^t}(s, a)$$

$$\geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} \max_a A^{\pi^t}(s, a) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} A^{\pi^t}(s, \pi^*(s))$$

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

Recall Π is restricted, denote $\epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi(s))]$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} [A^{\pi^t}(s, \pi^t(s))] + \epsilon_{\Pi} \geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] = \mathbb{E}_{s \sim d^*} \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \max_a A^{\pi^t}(s, a)$$

← PPL +

$$\geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} \max_a A^{\pi^t}(s, a) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} A^{\pi^t}(s, \pi^*(s)) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) (V^* - V^{\pi^t})(1 - \gamma)$$

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

$$\text{Recall } \Pi \text{ is restricted, denote } \epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^t(s)) \right] + \epsilon_{\Pi} \geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] = \mathbb{E}_{s \sim d^*} \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \max_a A^{\pi^t}(s, a)$$

$$\geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} \max_a A^{\pi^t}(s, a) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} A^{\pi^t}(s, \pi^*(s)) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) (V^* - V^{\pi^t})(1 - \gamma)$$

$$V^* - V^{\pi^t} \leq \sup_s \left(\frac{d^*(s)}{d_{\mu}^{\pi^t}(s)} \right) \frac{\varepsilon + \epsilon_{\Pi}}{1 - \gamma}$$

Upon Termination we get a locally optimal solution (or globally optimal if μ is nice)

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

Return π^t

1. No more positive advantage by one-step deviation from π^t 's own states

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

2. Indeed, we can say more if μ covers $d_{\mu}^{\pi^*}$, i.e., $C^* := \sup_s \frac{d_{\mu}^{\pi^*}(s)}{\mu(s)} < \infty$

$$\text{Recall } \Pi \text{ is restricted, denote } \epsilon_{\Pi} = \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] - \max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi(s)) \right]$$

$$\mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[A^{\pi^t}(s, \pi^t(s)) \right] + \epsilon_{\Pi} \geq \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\max_a A^{\pi^t}(s, a) \right] = \mathbb{E}_{s \sim d^*} \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \max_a A^{\pi^t}(s, a)$$

$$\geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} \max_a A^{\pi^t}(s, a) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) \mathbb{E}_{s \sim d^*} A^{\pi^t}(s, \pi^*(s)) \geq \inf_s \left(\frac{d_{\mu}^{\pi^t}(s)}{d^*(s)} \right) (V^* - V^{\pi^t})(1 - \gamma)$$

$$V^* - V^{\pi^t} \leq \sup_s \left(\frac{d^*(s)}{d_{\mu}^{\pi^t}(s)} \right) \frac{\varepsilon + \epsilon_{\Pi}}{1 - \gamma} \leq C^* \frac{\varepsilon + \epsilon_{\Pi}}{(1 - \gamma)^2}$$

Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

↑ ↑

Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

We start with a set of classifiers Π ; but we cannot guarantee $\pi^* \in \Pi$;




Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

We start with a set of classifiers Π ; but we cannot guarantee $\pi^* \in \Pi$;
What we can hope is that we can find **the best classifier in the class Π**



Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

We start with a set of classifiers Π ; but we cannot guarantee $\pi^* \in \Pi$;

What we can hope is that we can find **the best classifier in the class Π**

$$\mathbb{E}_{s \sim \nu, y = \pi^*(s)} \ell(\hat{\pi}(s), y) \leq \underbrace{\min_{\pi \in \Pi} \mathbb{E}_{s \sim \nu, y = \pi^*(s)} \ell(\pi(s), y)}_{\epsilon_{\Pi}} + \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{N}}}_{\text{statistical error: } \epsilon}$$

N: # of samples

Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

We start with a set of classifiers Π ; but we cannot guarantee $\pi^* \in \Pi$;
What we can hope is that we can find **the best classifier in the class Π**

$$\mathbb{E}_{s \sim \nu, y = \pi^*(s)} \mathcal{L}(\hat{\pi}(s), y) \leq \underbrace{\min_{\pi \in \Pi} \mathbb{E}_{s \sim \nu, y = \pi^*(s)} \mathcal{L}(\pi(s), y)}_{\epsilon_{\Pi}} + \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{N}}}_{\text{statistical error: } \epsilon}$$

In RL:

Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

We start with a set of classifiers Π ; but we cannot guarantee $\pi^* \in \Pi$;

What we can hope is that we can find **the best classifier in the class Π**

$$\mathbb{E}_{s \sim \nu, y = \pi^*(s)} \ell(\hat{\pi}(s), y) \leq \underbrace{\min_{\pi \in \Pi} \mathbb{E}_{s \sim \nu, y = \pi^*(s)} \ell(\pi(s), y)}_{\epsilon_{\Pi}} + \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{N}}}_{\text{statistical error: } \epsilon}$$

$$\text{In RL: } V^* - V^{\hat{\pi}} \leq \sup_s \left(\frac{d^*(s)}{\mu(s)} \right) \frac{\epsilon_{\Pi} + \epsilon}{(1 - \gamma)^2}$$

Connection to Agnostic Guarantees in Supervised Learning

Multi-class Classification (A many classes):

$$s \sim \nu, y \sim \pi^*(s), y \in [A]$$

We start with a set of classifiers Π ; but we cannot guarantee $\pi^* \in \Pi$;
What we can hope is that we can find **the best classifier in the class Π**

$$\mathbb{E}_{s \sim \nu, y = \pi^*(s)} \ell(\hat{\pi}(s), y) \leq \underbrace{\min_{\pi \in \Pi} \mathbb{E}_{s \sim \nu, y = \pi^*(s)} \ell(\pi(s), y)}_{\epsilon_{\Pi}} + \underbrace{\sqrt{\frac{\ln(|\Pi|/\delta)}{N}}}_{\text{statistical error: } \epsilon}$$

$$\text{In RL: } V^* - V^{\hat{\pi}} \leq \sup_s \left(\frac{d^*(s)}{\mu(s)} \right) \frac{\epsilon_{\Pi} + \epsilon}{(1 - \gamma)^2}$$

1. Multi-step prediction (not i.i.d), 2. We don't get to see samples from d^*

Compare the two Concentrability Coefficients from CPI and API:

$$\pi^{t+1} \leftarrow \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\pi}} \left[\sum_{\pi'} \mu(s, \pi') \right]$$

$$\text{API: } \max_{\pi \in \Pi} \sup_s \frac{d^{\pi}(s)}{\mu(s)} < \infty$$

Wide enough to cover **all policies**,
 i.e., making sure \hat{f} is accurate at all places
 where **any** policy would go

Just need to **cover the best in Π** ,
 steady improvement via incremental update

$$\text{CPI: } \sup_s \frac{d^*(s)}{\mu(s)} < \infty$$

Compare the two Concentrability Coefficients from CPI and API:

$$\text{API: } \max_{\pi \in \Pi} \sup_s \frac{d^\pi(s)}{\mu(s)} < \infty$$

Wide enough to cover **all policies**,
i.e., making sure \hat{A} is accurate at all places
where **any** policy would go

Just need to **cover the best in Π** ,
steady improvement via incremental update

$$\text{CPI: } \sup_s \frac{d^*(s)}{\mu(s)} < \infty$$

expert's distribution

1. Prior knowledge of how the optimal trajectories look like
2. Expert demonstrations (Imitation + RL)

$\approx \pi^$*

Summary of Policy Gradient Learning

1. PG formulation: $\mathbb{E}_{s,a \sim d^{\pi_\theta}} \nabla \ln \pi_\theta(a | s) A^{\pi_\theta}(s, a)$

Summary of Policy Gradient Learning

1. PG formulation: $\mathbb{E}_{s,a \sim d^{\pi_\theta}} \nabla \ln \pi_\theta(a | s) A^{\pi_\theta}(s, a)$

2. For tabular MDP, gradient ascent on KL-regularized objective converges to global optimality:

$$V^{\pi_\theta} + \lambda \sum_s \sum_a \ln \pi_\theta(a | s), \text{ where } \pi_\theta(a | s) \propto \exp(\theta_{s,a})$$

A

$$\ln \pi_\theta(a | s) \rightarrow -\infty \quad \pi_\theta(a | s) \rightarrow 0$$

Summary of Policy Gradient Learning

1. PG formulation: $\mathbb{E}_{s,a \sim d^{\pi_\theta}} \nabla \ln \pi_\theta(a | s) A^{\pi_\theta}(s, a)$

2. For tabular MDP, gradient ascent on KL-regularized objective converges to global optimality:

$$V^{\pi_\theta} + \lambda \sum_s \sum_a \ln \pi_\theta(a | s), \text{ where } \pi_\theta(a | s) \propto \exp(\theta_{s,a})$$

$KL(P_r^{\pi^k} || P_r^{\pi^{k+1}}) \leq \delta$

3. Natural Policy Gradient (trust region optimization) and its convergence (tabular, linear, & neural)

$$\hat{w} \in \arg \min_w \mathbb{E}_{s,a \sim d_v^{\pi_\theta}} \left[(w^\top \nabla_\theta \ln \pi_\theta(a | s) - A^{\pi_\theta}(a | s))^2 \right], \quad \theta' = \theta + \eta \hat{w}$$

$$\sigma_{\min} \begin{pmatrix} E \phi(s,a)^\top \\ \text{sample } \phi(s,a) \end{pmatrix} > c \epsilon R^d$$

Summary of Policy Gradient Learning

1. PG formulation: $\mathbb{E}_{s,a \sim d^{\pi_\theta}} \nabla \ln \pi_\theta(a | s) A^{\pi_\theta}(s, a)$

2. For tabular MDP, gradient ascent on KL-regularized objective converges to global optimality:

$$V^{\pi_\theta} + \lambda \sum_s \sum_a \ln \pi_\theta(a | s), \text{ where } \pi_\theta(a | s) \propto \exp(\theta_{s,a})$$

3. Natural Policy Gradient (trust region optimization) and its convergence (tabular, linear, & neural)

$$\widehat{w} \in \arg \min_w \mathbb{E}_{s,a \sim d_v^{\pi_\theta}} \left[\left(w^\top \nabla_\theta \ln \pi_\theta(a | s) - A^{\pi_\theta}(a | s) \right)^2 \right], \quad \theta' = \theta + \eta \widehat{w}$$

4. The incremental nature of NPG/CPI/PPO and its advantage comparing to naive API

CPI (TRPO): $V^{\pi^{t+1}} > V^{\pi^t} > V^{\pi^{t-1}}$, thanks to $\|d_\mu^{\pi^{t+1}} - d_\mu^{\pi^t}\|_1$ is small (i.e., incremental)

Δ

and API could oscillate and never converges

Next week on Control Theory:

$$x_{t+1} = Ax_t + Bu_t$$

Basics of Optimal Control on Linear Quadratic Regulators
(no learning, just planning/control)

$$\max_{\pi} E_{S \sim d_{\mu}^{\pi t}} [A^{\pi t}(s, \pi(s))]]$$

$$\text{s.t. } \forall s \in S$$

$$\| \pi(\cdot | s) - \pi^t(\cdot | s) \|_{TV} \leq \delta$$

$$\textcircled{1} \text{ CPI : } \| \pi^{t+1}(\cdot | s) - \pi^t(\cdot | s) \|_{TV} \leq \alpha$$