"Convergence" &

Trust Region Policy Optimization

Lucas Janson and Sham Kakade

CS/Stat 184: Introduction to Reinforcement Learning Fall 2022

1

Today

- Next class: embedded ethics (from Jenna Donohue, a postdoc in the Philosophy dept)
 - Course Plan: consider different ethical implications of different possible utility functions for a (fictional) RL algorithm that was setting dynamic prices for rides.
 - Please come to the next class. (There will be a discussion.)
 - Please do the assigned reading (John Rawls) in advance.
- Recap++
- Today:
 - 1. Convergence of Fitted Policy Iteration
 - 2. Trust Region Policy Optimization

Recap + Examples

Is there an iterative version of Policy Evaluation?

(that is faster, but approximate?)



[Iterative Policy Eval Subroutine/TD] input: policy π , sample size N, init w_0



[Iterative Policy Eval Subroutine/TD] input: policy π , sample size N, init w_0

1. Sample trajectories $\tau_1, \ldots \tau_N \sim \rho_{\pi}$ which gives us a dataset *D* (each trajectory is of the form $\tau_i = \{s_0, a_0, r_0, \ldots s_{H-1}, a_{H-1}, r_{H-1}, \}$)

2. For
$$k = 0, ..., K$$
:



[Iterative Policy Eval Subroutine/TD] input: policy π , sample size N, init w_0

1. Sample trajectories $\tau_1, \ldots \tau_N \sim \rho_{\pi}$ which gives us a dataset *D* (each trajectory is of the form $\tau_i = \{s_0, a_0, r_0, \ldots s_{H-1}, a_{H-1}, r_{H-1}, \}$)

2. For
$$k = 0, ..., K$$
:

1. Sample a transition $(s_h, r_h, s_{h+1}) \in D$ and update:

 $V_{k+1}(s_h) = V_k(s_h) - \eta_k \Big(V_k(s_h) - (r_h + V_k(s_{h+1})) \Big)$

3. Return the function V_K as an estimate of V^π

Fit $V^{\pi}(s)$ using the iterative policy evaluation alg.

Fit $V^{\pi}(s)$ using the iterative policy evaluation alg.

[Iterative Policy Eval Subroutine/TD]

input: policy π , sample size N, init w_0

- 1. Sample trajectories $\tau_1, \ldots \tau_N \sim \rho_{\pi}$ which gives us a dataset D
- 2. For k = 0, ..., K:
 - 1. Construct an *empirical loss function*:

$$L_{k}(w) = \frac{1}{NH} \sum_{i=1}^{N} \sum_{(s_{h}, r_{h}, s_{h+1}) \in \tau_{i}} \left(f_{w}(s_{h}) - \left(r_{h} + f_{w_{k}}(s_{h+1}) \right) \right)^{2}$$

2. Update with either: full minimization:

$$w_{k+1} \approx \arg\min_{w} L_k(w)$$

Fit $V^{\pi}(s)$ using the iterative policy evaluation alg.



Fit $V^{\pi}(s)$ using the iterative policy evaluation alg.

[Iterative Policy Eval Subroutine/TD]

input: policy π , sample size N, init w_0

- 1. Sample trajectories $\tau_1, \ldots \tau_N \sim \rho_{\pi}$ which gives us a dataset D
- 2. For k = 0, ..., K:
 - 1. Construct an *empirical loss function*:

$$L_{k}(w) = \frac{1}{NH} \sum_{i=1}^{N} \sum_{(s_{h}, r_{h}, s_{h+1}) \in \tau_{i}} \left(f_{w}(s_{h}) - \left(r_{h} + f_{w_{k}}(s_{h+1}) \right) \right)^{2}$$

2. Update with either: full minimization:

 $w_{k+1} \approx \arg \min_{w} L_k(w)$ TD learning: (one step of SGD) $w_{k+1} = w_k - \eta_k \,\overline{\nabla} \, L_k(w_k)$

3. Return the function $f_{_{\!W_K}}$ as an estimate of V^π

Fitted Dynamic Programming Methods for learning Q^* and π^*

Policy Iteration (PI)

- Initialization: choose a policy $\pi^0: S \mapsto A$
- For k = 0, 1, ...
 - 1. Policy Evaluation: compute $Q^{\pi^k}(s, a)$
 - 2. Policy Improvement: set

$$\pi^{k+1}(s) := \arg\max_{a} Q^{\pi^{k}}(s, a)$$

Fitted Policy Iteration: (aka Approximate Policy Iteration API)

1. Initialize staring policy π_0 , samples size M 2. For k = 0, ... : 1. [Q-Evaluation Subroutine] Using N sampled trajectories, $\tau_1, \ldots, \tau_N \sim \rho_{\pi_k}$, try to learn \widetilde{a} $\widetilde{Q}_k(s,a) \approx Q_h^{\pi_k}(s,a) \simeq \int_{\mathcal{W}} (5,q,h)$ 2. Policy Update $\pi_{k+1}(s) := \arg\max[Q^{\pi_k}(s, a)]$ 3. Return Q_{K} and π_{K} as an estimate of Q^{\star} and π^{\star}

SE(SG)

Alternative Version: Bellman Operator \mathcal{T} on Q (HW2 Q2 is the Q-version of the Bellman Equations)

- Given a function $Q: S \times A \mapsto \mathbb{R}$, define $\mathcal{T}Q: S \times A \mapsto \mathbb{R}$ as $(\mathcal{T}Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a' \in A} Q(s', a')$
- (Bellman equations for Q) Q is equal to Q^* if and only if $\mathcal{T}Q = Q$.

Q-Value Iteration Algorithm:

The Offline Learning Setting:

We don't know the MDP and our data collection is under some fixed distribution.

The Finite Horizon, Offline Learning Setting:

- We have N trajectories $au_1, \ldots au_N \sim
 ho_{\pi_{data}}$
- π_{data} is often referred to as our data collection policy.

t dala 7 ' · · r. ~ tt. (s, q, r, s') $p = s \sim P(sq)$ r(s,q)

Q-Learning for "tabular" case



Fitted Q-Iteration

input: offline dataset $au_1, \ldots au_N \sim
ho_{\pi_{data}}$, init w_0 1. For k = 0, 1, ..., K: 1. Construct an *empirical loss function*: $L_{k}(w) = \frac{1}{NH} \sum_{i=1}^{N} \sum_{(s_{h}, a_{h}, r_{h}, s_{h+1}) \in \tau_{i}}^{N} \left(f_{w}(s_{h}, a_{h}) - \left(r_{h} + \max_{a} f_{w_{k}}(s_{h+1}, a) \right) \right)^{2}$ 2. Update with either: full minimization: $w_{k+1} \approx \arg\min L_k(w)$ Q-learning: (one step of SGD) $w_{k+1} = w_k - \eta_k \nabla L_k(w_k)$ 2. Return the function $f_{\widetilde{w}_{arkappa}}$ as an estimate of Q^{\star}

Today: "Convergence" & Trust Region Policy Optimization

Outline:

- 1. Convergence of Fitted Policy Iteration
 - 1. "Tabular" case
 - 2. Fitted case
- 2. Trust Region Policy Optimization
 - 1. Quick intro on KL-divergence
 - 2. TRPO formulation

(the easiest case to think about fitted Policy Iteration)

(the easiest case to think about fitted Policy Iteration)

1. For k = 0, ... :

(the easiest case to think about fitted Policy Iteration)

- 1. For k = 0, ...:
 - 1. Q-Evaluation:

For each (s, a, h), suppose that:

(the easiest case to think about fitted Policy Iteration)

- 1. For k = 0, ...:
 - 1. Q-Evaluation:

For each (s, a, h), suppose that:

1. we are able to draw M trajectories as follows:

start at $(s_h = s, a_h = a)$, run π_k , and end trajectory at time H

(the easiest case to think about fitted Policy Iteration)

1. For k = 0, ...: 1. Q-Evaluation: For each (s, a, h), suppose that: 1. we are able to draw M trajectories as follows: start at $(s_h = s, a_h = a)$, run π_k , and end trajectory at time H2. Set $\widetilde{Q}_k(s, a)$ as the empirical average of the cumulative reward on these trajectories. (i.e. $\widetilde{Q}_k(s, a)$ is unbiased sample of $Q_h^{\pi_k}(s, a)$, with N samples)

(the easiest case to think about fitted Policy Iteration)

- 1. For k = 0, ...:
 - 1. Q-Evaluation:

For each (s, a, h), suppose that:

1. we are able to draw M trajectories as follows:

start at $(s_h = s, a_h = a)$, run π_k , and end trajectory at time H

- 2. Set $\widetilde{Q}_k(s, a)$ as the empirical average of the cumulative reward on these trajectories. (i.e. $\widetilde{Q}_k(s, a)$ is unbiased sample of $Q_h^{\pi_k}(s, a)$, with *N* samples)
- 2. Policy Update

 $\pi_{k+1}(s) := \arg \max_{a} \widetilde{Q}^{\pi_k}(s, a)$

(the easiest case to think about fitted Policy Iteration)

- 1. For k = 0, ...:
 - 1. Q-Evaluation:

For each (s, a, h), suppose that:

1. we are able to draw M trajectories as follows:

start at $(s_h = s, a_h = a)$, run π_k , and end trajectory at time H

- 2. Set $\widetilde{Q}_k(s, a)$ as the empirical average of the cumulative reward on these trajectories. (i.e. $\widetilde{Q}_k(s, a)$ is unbiased sample of $Q_h^{\pi_k}(s, a)$, with *N* samples)
- 2. Policy Update

$$\pi_{k+1}(s) := \arg\max \widetilde{Q}^{\pi_k}(s,a)$$

2. Return \widetilde{Q}_K and π_K as an estimate of Q^{\star} and π^{\star}

(the easiest case to think about fitted Policy Iteration)



[Theorem] Using polynomial many total samples and polynomial computation time (in $|S|, |A|, H, 1/\epsilon$), we have that $\|\widetilde{Q}_K - Q^{\star}\|_{\infty} \leq \epsilon$ and $\|Q^{\pi_K} - Q^{\star}\|_{\infty} \leq \epsilon$.

Outline:

- 1. Convergence of Fitted Policy Iteration
 - 1. "Tabular" case
 - 2. Fitted case
- 2. Trust Region Policy Optimization
 - 1. Quick intro on KL-divergence
 - 2. TRPO formulation

First: let's summarize a few things about Supervised Learning

Recap on Supervised Learning: Classification

Recap on Supervised Learning: Classification

Given i.i.d examples at training:


Recap on Supervised Learning: Classification

Given i.i.d examples at training:





Recap on Supervised Learning: Classification

Given i.i.d examples at training:





Using function approximator, we are able to predict on cats/dogs that we **never see before** (i.e., we **generalize**)







Using function approximation, we are able to predict on the value of some house not from the training data

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

Recap on Supervised Learning: regression $\int_{1}^{\infty} (x_{i}) = Q^{\pi_{k}}(x_{i})$

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$ We want to approximate f^* using finite training samples;

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

We want to approximate f^{\star} using finite training samples;

Let us introduce an abstract function class $\mathcal{F} = \{f : \mathcal{X} \mapsto \mathbb{R}\}$, and do least squares:

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

We want to approximate f^{\star} using finite training samples;

Let us introduce an abstract function class $\mathcal{F} = \{f : \mathcal{X} \mapsto \mathbb{R}\}$, and do least squares:

$$\hat{f} = \arg\min_{f\in\mathscr{F}}\sum_{i=1}^{N} (f(x_i) - y_i)^2$$

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

We want to approximate f^{\star} using finite training samples;

Let us introduce an abstract function class $\mathscr{F} = \{f : \mathscr{X} \mapsto \mathbb{R}\}$, and do least squares:

Empirical Risk Minimizer (ERM)
$$\hat{f} = \arg \min_{f \in \mathscr{F}} \sum_{i=1}^{N} (f(x_i) - y_i)^2$$

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

We want to approximate f^{\star} using finite training samples;

Let us introduce an abstract function class $\mathscr{F} = \{f : \mathscr{X} \mapsto \mathbb{R}\}$, and do least squares:

Empirical Risk Minimizer (ERM)
$$\hat{f} = \arg \min_{f \in \mathscr{F}} \sum_{i=1}^{N} (f(x_i) - y_i)^2$$

Q: quality of ERM \hat{f} ?

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^{\star}(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

$$\hat{f} = \arg\min_{f\in\mathscr{F}}\sum_{i=1}^{N} (f(x_i) - y_i)^2$$

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

$$\hat{f} = \arg\min_{f \in \mathscr{F}} \sum_{i=1}^{N} (f(x_i) - y_i)^2$$

Supervised learning theory (e.g., VC theory) says that we can indeed **generalize**, i.e., we can predict well **under the same distribution:**

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

$$\hat{f} = \arg\min_{f \in \mathscr{F}} \sum_{i=1}^{N} (f(x_i) - y_i)^2$$

Supervised learning theory (e.g., VC theory) says that we can indeed **generalize**, i.e., we can predict well **under the same distribution:**

Assume $f^{\star} \in \mathcal{F}$ (this is called realizability), we can expect:

We have a data distribution \mathcal{D} , $x_i \sim \mathcal{D}$, $y_i = f^*(x_i) + \epsilon_i$, where noise $\mathbb{E}[\epsilon_i] = 0$, $|\epsilon_i| \le c$

$$\hat{f} = \arg\min_{f\in\mathscr{F}}\sum_{i=1}^{N} (f(x_i) - y_i)^2$$

Supervised learning theory (e.g., VC theory) says that we can indeed **generalize**, i.e., we can predict well **under the same distribution:**

Assume $f^* \in \mathcal{F}$ (this is called realizability), we can expect: $\mathbb{E}_{x \sim \mathcal{D}} \left(\hat{f}(x) - f^*(x) \right)^2 \leq \delta \quad \text{for } f^*(x) \quad f^*($

Supervise Learning can fail if there is train-test distribution mismatch

However, for some $\mathscr{D}' \neq \mathscr{D}$, $\mathbb{E}_{x \sim \mathscr{D}'} (f(x) - f^{\star}(x))^2$ might be arbitrarily large



Supervise Learning can fail if there is train-test distribution mismatch

However, for some $\mathscr{D}' \neq \mathscr{D}$, $\mathbb{E}_{x \sim \mathscr{D}'} (f(x) - f^{\star}(x))^2$ might be arbitrarily large



Deeper neural nets and larger datasets are typically not enough to address "distribution shift"

Back to RL

• For all k, suppose that:

$$E_{\tau \sim \rho_{\pi_k}} \left[\sum_{h=1}^n \left(\widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right)^2 \right] \le \delta, \text{ and } \max_{s, a} \left| \widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right| \le \delta_{\infty}$$

• For all k, suppose that:

• δ : the average case supervised learning error (reasonable to expect this can be made small) δ_{∞} : the worse case error (often unreasonable to expect to be small)

• For all k, suppose that:

$$E_{\tau \sim \rho_{\pi_k}} \left[\sum_{h=1}^{n} \left(\widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right)^2 \right] \le \delta, \text{ and } \max_{s, a} \left| \widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right| \le \delta_{\infty}$$

• δ : the average case supervised learning error (reasonable to expect this can be made small) δ_{∞} : the worse case error (often unreasonable to expect to be small)

[Theorem:] We have that:

• For all k, suppose that:

$$E_{\tau \sim \rho_{\pi_k}} \left[\sum_{h=1}^n \left(\widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right)^2 \right] \le \delta, \text{ and } \max_{s, a} \left| \widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right| \le \delta_{\infty}$$

• δ : the average case supervised learning error (reasonable to expect this can be made small) PXE δ_{∞} : the worse case error (often unreasonable to expect to be small)

[Theorem:] We have that:

• One step performance degradation is bounded by the worst case error: $Q^{k+1}(s) \ge Q^k(s) - 2H\delta_{\infty}$ (and equality possible in some examples).

Fitted Policy Improvement Guarantees • For all k, suppose that: $\leq \delta$, and max • δ : the average case supervised learning error (reasonable to expect this can be made small) δ_{∞} : the worse case error (often unreasonable to expect to be small) [Theorem:] We have that: • One step performance degradation is bounded by the worst case error: $Q^{k+1}(s) \ge Q^k(s) - 2H\delta_{\infty}$ (and equality possible in some examples). • For large enough K, final performance also governed by the worst case error: $Q^{\pi_K}(s) \ge Q^{\star}(s)$

• For all k, suppose that:

$$E_{\tau \sim \rho_{\pi_k}} \left[\sum_{h=1}^n \left(\widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right)^2 \right] \le \delta, \text{ and } \max_{s, a} \left| \widetilde{Q}_k(s_h, a_h) - Q_h^{\pi_k}(s_h, a_h) \right| \le \delta_{\infty}$$

• δ : the average case supervised learning error (reasonable to expect this can be made small) δ_{∞} : the worse case error (often unreasonable to expect to be small)

[Theorem:] We have that:

• One step performance degradation is bounded by the worst case error:

 $-Q^{k+1}(s) \ge Q^k(s) - 2H\delta_{\infty}$ (and equality possible in some examples).

• For large enough *K*, final performance also governed by the worst case error: $Q^{\pi_{K}}(s) \ge Q^{\star}(s) - 2H^{2}\delta_{\infty}$

• (Intuition) If it somehow turns out that, for all iterations k, the density under the next policy, uniformly does not differ from that of previous policy, i.e. that along the path of the algo.

26

$$\max_{s,a,h} \left(\frac{\Pr(s_h = s, a_h = a \mid \pi_{k+1})}{\Pr(s_h = s, a_h = a \mid \pi_k)} \right) \le C_{\infty}$$

then we can bound our sub-optimality by the average case error:

 $Q^{\pi_{K}}(s) \geq Q^{\star}(s) - 2H^{2} \cdot C_{m} \cdot \delta$

Outline:

- 1. Convergence of Fitted Policy Iteration
 - 1. "Tabular" case
 - 2. Fitted case
- 2. Trust Region Policy Optimization
 - 1. Quick intro on KL-divergence
 - 2. TRPO formulation

Given two distributions P & Q, where $P \in \Delta(X), Q \in \Delta(X)$, KL Divergence is defined as:

$$KL(P \mid Q) = \mathbb{E}_{x \sim P} \left[\ln \frac{P(x)}{Q(x)} \right]$$

Given two distributions P & Q, where $P \in \Delta(X), Q \in \Delta(X)$, KL Divergence is defined as:

$$KL(P \mid Q) = \mathbb{E}_{x \sim P} \left[\ln \frac{P(x)}{Q(x)} \right]$$

Examples:

If Q = P, then KL(P | Q) = KL(Q | P) = 0

Given two distributions P & Q, where $P \in \Delta(X), Q \in \Delta(X)$, KL Divergence is defined as:

$$KL(P \mid Q) = \mathbb{E}_{x \sim P} \left[\ln \frac{P(x)}{Q(x)} \right]$$

Examples:

If
$$Q = P$$
, then $KL(P | Q) = KL(Q | P) = 0$
If $P = \mathcal{N}(\mu_1, \sigma^2 I), Q = \mathcal{N}(\mu_2, \sigma^2 I)$, then $KL(P | Q) = ||\mu_1 - \mu_2||_2^2 / \sigma^2$

Given two distributions P & Q, where $P \in \Delta(X), Q \in \Delta(X)$, KL Divergence is defined as:

$$KL(P \mid Q) = \mathbb{E}_{x \sim P} \left[\ln \frac{P(x)}{Q(x)} \right]$$

Examples:

If
$$Q = P$$
, then $KL(P | Q) = KL(Q | P) = 0$

If $P = \mathcal{N}(\mu_1, \sigma^2 I), Q = \mathcal{N}(\mu_2, \sigma^2 I)$, then $KL(P | Q) = \|\mu_1 - \mu_2\|_2^2 / \sigma^2$

Fact:

 $KL(P \mid Q) \ge 0$, and being 0 if and only if P = Q

Outline:

- 1. Convergence of Fitted Policy Iteration
- 2. Trust Region Policy Optimization
 - 1. Quick intro on KL-divergence
 - 2. TRPO formulation

An "idealized" trust region formulation for policy update: (back to direct policy optimization)

At iteration t, with π_{θ_t} at hand, we compute θ_{t+1} as follows:

 $\max_{\pi_{\theta}} J(\theta) - J(\theta_t)$

s.t., $KL\left(\rho_{\pi_{\theta_t}} | \rho_{\pi_{\theta}}\right) \leq \delta$

An "idealized" trust region formulation for policy update: (back to direct policy optimization)

At iteration t, with π_{θ_t} at hand, we compute θ_{t+1} as follows:

 $\max_{\pi_{\theta}} J(\theta) - J(\theta_t)$

s.t.,
$$KL\left(\rho_{\pi_{\theta_{t}}}|\rho_{\pi_{\theta}}\right) \leq \delta$$

We want to maximize performance improvement starting at π_{θ_t} , but we want the new policy to be close to π_{θ_t} (in the KL sense)

A trust region formulation for policy update:

At iteration t, with π_{θ_t} at hand, we compute θ_{t+1} as follows:

A trust region formulation for policy update:

At iteration t, with π_{θ_t} at hand, we compute θ_{t+1} as follows:

$$\max_{\pi_{\theta}} \mathbb{E}_{s_{0},...s_{H-1} \sim \rho_{\theta_{t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} \left[A^{\pi_{\theta_{t}}}(s,a) \right] \right]$$

s.t., $KL\left(\rho_{\pi_{\theta_{t}}} | \rho_{\pi_{\theta}} \right) \leq \delta$

A trust region formulation for policy update:

At iteration t, with π_{θ_t} at hand, we compute θ_{t+1} as follows:

$$\max_{\pi_{\theta}} \mathbb{E}_{s_{0},...s_{H-1} \sim \rho_{\theta_{t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} \left[A^{\pi_{\theta_{t}}}(s,a) \right] \right]$$

s.t., $KL\left(\rho_{\pi_{\theta_{t}}} | \rho_{\pi_{\theta}} \right) \leq \delta$

We want to maximize local advantage against π_{θ_t} , but we want the new policy to be close to π_{θ_t} (in the KL sense)
A trust region formulation for policy update:

At iteration t, with π_{θ_t} at hand, we compute θ_{t+1} as follows:

$$\max_{\pi_{\theta}} \mathbb{E}_{s_{0},...s_{H-1} \sim \rho_{\theta_{t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} \left[A^{\pi_{\theta_{t}}}(s,a) \right] \right]$$

s.t., $KL\left(\rho_{\pi_{\theta_{t}}} | \rho_{\pi_{\theta}} \right) \leq \delta$

We want to maximize local advantage against π_{θ_i} , but we want the new policy to be close to π_{θ_i} (in the KL sense)

How we can actually do the optimization here? After all, we don't even know the analytical form of trajectory likelihood...

Summary:

- 1. Convergence of Fitted Policy Iteration
- 2. Trust Region Policy Optimization
 - 1. Quick intro on KL-divergence
 - 2. TRPO formulation



