

Trust Region Policy Optimization & the Natural Policy Gradient

Lucas Janson and Sham Kakade

**CS/Stat 184: Introduction to Reinforcement Learning
Fall 2022**

Today

- Announcements:

If you are an undergraduate student at Harvard and are possibly interested in pursuing research, formally or informally, with the [ML foundations group](#), please fill in the following form: <https://forms.gle/yCiTfbXn31x2RQtHA>

- Recap

- Today:

1. ~~Convergence of Fitted Policy Iteration~~
2. Trust Region Policy Optimization

of the natural gradient

Recap

[Policy Eval Subroutine]: TD Learning for “tabular” case

[Iterative Policy Eval Subroutine/TD]

input: policy π , sample size N , init w_0

1. Sample trajectories $\tau_1, \dots, \tau_N \sim \rho_\pi$ which gives us a dataset D
(each trajectory is of the form $\tau_i = \{s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1}, \dots\}$)
2. For $k = 0, \dots, K$:
 1. Sample a transition $(s_h, r_h, s_{h+1}) \in D$ and update:
$$V_{k+1}(s_h) = V_k(s_h) - \eta_k \left(V_k(s_h) - (r_h + V_k(s_{h+1})) \right)$$
3. Return the function V_K as an estimate of V^π

Policy eval method.

Q-Learning for “tabular” case

[Iterative Policy Eval Subroutine/TD]

input: **offline dataset** $\tau_1, \dots, \tau_N \sim \rho_{\pi_{\text{data}}}$

1. For $k = 0, \dots, K$:

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

$$Q_{k+1}(s_h, a_h) = Q_k(s_h, a_h) - \eta_k \left(Q_k(s_h, a_h) - \left(r_h + \max_{a'} Q_k(s_{h+1}, a') \right) \right)$$

2. Return the function Q_K as an estimate of Q^*

to learn/estimate Q^*

again the idea was to enforce
B. consistency.

Fitted Policy Iteration: (aka Approximate Policy Iteration API)

1. Initialize starting policy π_0 , samples size M
2. For $k = 0, \dots$:
 1. [Q-Evaluation Subroutine]
Using M sampled trajectories, $\tau_1, \dots, \tau_N \sim \rho_{\pi_k}$,
$$\widetilde{Q}_k(s, a) \approx Q_h^{\pi_k}(s, a)$$
 2. Policy Update
$$\pi_{k+1}(s) := \arg \max_a \widetilde{Q}^{\pi_k}(s, a)$$
3. Return \widetilde{Q}_K and π_K as an estimate of Q^\star and π^\star

Recap on Supervised Learning: regression

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

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

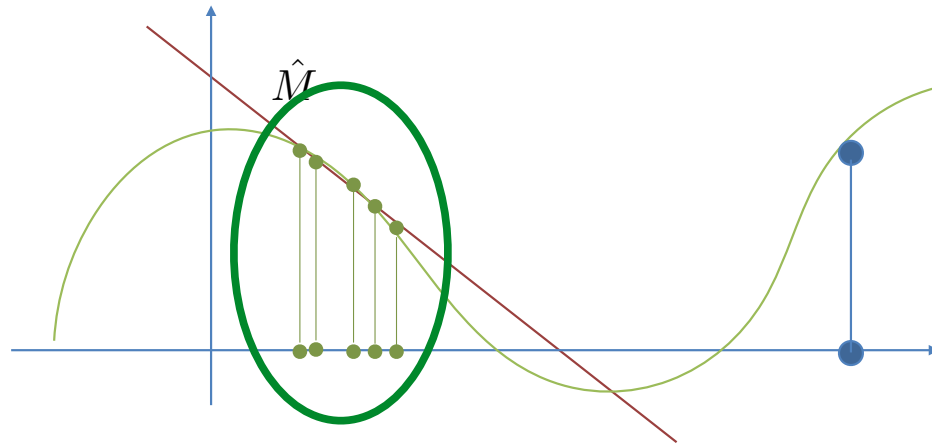
Supervised learning theory (e.g., VC theory) says that we can indeed **generalize**,
i.e., we can predict well **under the same distribution**:

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

$$\mathbb{E}_{x \sim \mathcal{D}} \left(\hat{f}(x) - f^\star(x) \right)^2 \leq \delta$$

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

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



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

Fitted Policy Improvement Guarantees

Fitted Policy Improvement Guarantees

- For all k , suppose that:

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

- 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)

Fitted Policy Improvement Guarantees

- For all k , suppose that:

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

- 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:

Fitted Policy Improvement Guarantees

- For all k , suppose that:

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

- δ : the average case supervised learning error (reasonable to expect this can be made small)
- δ_∞ : the worst case error (often unreasonable to expect to be small)

[Theorem:] We have that:

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

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

Fitted Policy Improvement Guarantees

- For all k , suppose that:

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

- 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**:

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

- For large enough K , final performance also governed by the **worst case error**:

$$Q^{\pi_K}(s) \geq Q^*(s) - 2H^2\delta_\infty$$

Fitted Policy Improvement Guarantees

- For all k , suppose that:

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

- 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**:

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

- For large enough K , final performance also governed by the **worst case error**:

$$Q^{\pi_K}(s) \geq Q^*(s) - 2H^2\delta_\infty$$

- (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

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

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

$$Q^{\pi_K}(s) \geq Q^*(s) - 2H^2 \cdot C_\infty \cdot \delta$$

Today:

~~Optimality in Markov Decision Processes~~

Outline:

1. Quick intro on KL-divergence
& the visitation measure
2. A Trust-Region Formulation for Policy Optimization
3. Algorithm: Natural Policy Gradient

KL-divergence: measures the distance between two distributions

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

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

KL-divergence: measures the distance between two distributions

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

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

always

Examples:

If $Q = P$, then $KL(P \parallel Q) = KL(Q \parallel P) = 0$

$$KL(P \parallel Q) \geq 0$$

KL-divergence: measures the distance between two distributions

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 $KL(P \mid Q) = KL(Q \mid P) = 0$

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

KL-divergence: measures the distance between two distributions

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 $KL(P \mid Q) = KL(Q \mid P) = 0$

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

Fact:

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

Outline:

1. Quick intro on KL-divergence
2. A Trust-Region Formulation for Policy Optimization (TRPO)
3. Algorithm: Natural Policy Gradient

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:

$$\begin{aligned} \max_{\theta} \mathbb{E}_{s_0, \dots, s_{H-1} \sim \rho_{\theta_t}} & \left[\sum_{h=0}^{H-1} \mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \\ \text{s.t., } & KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) \leq \delta \end{aligned}$$

A trust region formulation for policy update: $A^\pi(s, a)$
 $= Q^\pi(s, a) - V^\pi(s)$

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

$$\begin{aligned} \max_{\theta} \mathbb{E}_{s_0, \dots, s_{H-1} \sim \rho_{\theta_t}} & \left[\sum_{h=0}^{H-1} \mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \\ \text{s.t., } & KL(\rho_{\pi_{\theta_t}} | \rho_{\pi_{\theta}}) \leq \delta \end{aligned}$$

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)**

Some Helpful Notation: Visitation Measures

- Visitation probability at time h : $\mathbb{P}_h(s_h, a_h | \mu, \pi)$
(recall that we absorb h , into the state, i.e. $s \leftarrow (s, h)$)
- Average Visitation Measure:

$$d_\mu^\pi(s, a) = \frac{1}{H} \sum_{h=0}^{H-1} \mathbb{P}_h(s, a | \mu, \pi)$$

$$d_\mu^\pi(s) = \frac{1}{H} \sum_{h=0}^{H-1} \mathbb{P}_h(s | \mu, \pi)$$

$s_0 \sim \mu$, following π

- With this def, we have:

$$J(\theta) := E_{s_0 \sim \mu_0} [V^{\pi_\theta}(s_0)]$$

$$DJ(\theta)$$

$$= H E_{s, a \sim d_\mu^\pi} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\theta}(s, a) \right]$$

$$= E \left[\sum_{h=0}^{H-1} r(s_h, a_h) \mid \mu_0, \pi_{\theta} \right] = H E_{s, a \sim d_\mu^\pi} [r(s, a)]$$

$$= H \cdot E_{s \sim d_\mu^{\pi_\theta}} E_{a \sim \pi_\theta(s)} [r(s, a)]$$

Visitation Measures: the discounted case

- Visitation probability at time h : $\mathbb{P}_h(s_h, a_h | \mu, \pi)$
(recall that we absorb h , into the state, i.e. $s \leftarrow (s, h)$)

- Average Visitation Measure:

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

- With this def, we have:

$$\begin{aligned} J(\theta) &:= E_{s_0 \sim \mu_0} [V^{\pi_\theta}(s_0)] \\ &= E \left[\sum_{h=0}^{\infty} \gamma^h r_h \mid \mu_0, \pi_\theta \right] \\ &= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s \sim d_\mu^{\pi_\theta}} E_{a \sim \pi(s)} [r(s, a)] \end{aligned}$$

Equivalently,

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

Equivalently,

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

$$\max_{\theta} H \cdot \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right]$$

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

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)

Equivalently,

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

$$\max_{\theta} H \cdot \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right]$$

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

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)

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

Outline:

1. Quick intro on KL-divergence
2. A Trust-Region Formulation for Policy Optimization
3. Algorithm: Natural Policy Gradient

A trust region formulation for policy update:

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

$$\begin{aligned} \max_{\theta} \quad & H \cdot \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \\ \text{s.t.}, \quad & KL \left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}} \right) \leq \delta \end{aligned}$$

High-level strategy:

1. First-order Taylor expansion on the objective at θ_t
2. second-order Taylor expansion of the constraint at θ_t

Simplify Objective Function

$$H. \quad \max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right]$$

Simplify Objective Function

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right]$$

Since the objective is also non-linear,
let's do first order-taylor expansion on it:

Simplify Objective Function $\nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}} [f(s)]$

$$H \cdot \max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) f(s)]$$

Since the objective is also non-linear,
let's do first order-taylor expansion on it:

$$H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \approx H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta_t}(s)} A^{\pi_{\theta_t}}(s, a) \right] + \underbrace{\mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta_t}(s)} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \right] A^{\pi_{\theta_t}}(s, a) \right]}_{\nabla_{\theta} J(\pi_{\theta_t})} \cdot (\theta - \theta_t)$$

$$f_{\theta_t}(\theta) = f_{\theta_t}(\theta_t) + \left(\nabla f_{\theta_t}(\theta) \right)_{\theta=\theta_t} (\theta - \theta_t) + o(\quad)$$

because $\mathbb{E}_{a \sim \pi(s)} [A^{\pi}(s, a)] = 0$

$$J(\theta) = H \cdot \mathbb{E}_{s \sim d^{\theta}} \mathbb{E}_{a \sim \pi^{\theta}} [r(s, a)]$$

Simplify Objective Function

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right]$$

Since the objective is also non-linear,
let's do first order-taylor expansion on it:

$$H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \approx H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta_t}(s)} A^{\pi_{\theta_t}}(s, a) \right] + \underbrace{\mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta_t}(s)} \nabla_{\theta} \ln \pi_{\theta_t}(a | s) A^{\pi_{\theta_t}}(s, a) \right]}_{\nabla_{\theta} J(\pi_{\theta_t})} \cdot (\theta - \theta_t)$$

$$\Rightarrow \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t)$$

Simplify Constraint via second-order Taylor Expansion:

Simplify Constraint via second-order Taylor Expansion:

$$KL(\rho_{\theta_t} | \rho_{\theta}) := \ell(\theta)$$

Simplify Constraint via second-order Taylor Expansion:

$$KL(\rho_{\theta_t} | \rho_{\theta}) := \ell(\theta)$$

$$\ell(\theta) \approx \ell(\theta_t) + \nabla \ell(\theta_t)^\top (\theta - \theta_t) + \frac{1}{2} (\theta - \theta_t)^\top \nabla_{\theta}^2 \ell(\theta_t) (\theta - \theta_t)$$

Simplify Constraint via second-order Taylor Expansion:

$$KL(\rho_{\theta_t} | \rho_{\theta}) := \ell(\theta)$$

$$\ell(\theta) \approx \ell(\theta_t) + \nabla \ell(\theta_t)^\top (\theta - \theta_t) + \frac{1}{2} (\theta - \theta_t)^\top \nabla_{\theta}^2 \ell(\theta_t) (\theta - \theta_t)$$

$$\ell(\theta_t) = KL(\rho_{\theta_t} | \rho_{\theta_t}) = 0$$

Simplify Constraint via second-order Taylor Expansion:

$$KL(\rho_{\theta_t} | \rho_{\theta}) := \ell(\theta)$$

$$\ell(\theta) \approx \ell(\theta_t) + \nabla \ell(\theta_t)^\top (\theta - \theta_t) + \frac{1}{2} (\theta - \theta_t)^\top \nabla_{\theta}^2 \ell(\theta_t) (\theta - \theta_t)$$

$$\ell(\theta_t) = KL(\rho_{\theta_t} | \rho_{\theta_t}) = 0$$

We will show that $\nabla_{\theta} \ell(\theta_t) = 0$, and $\nabla^2 \ell(\theta_t)$ has a nice form!

The gradient of the KL-divergence is zero at θ_t

Change from trajectory distribution to state-action distribution:

The gradient of the KL-divergence is zero at θ_t

Change from trajectory distribution to state-action distribution:

$$KL(\rho_{\pi_{\theta_t}} | \rho_{\pi_{\theta}}) = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \ln \frac{\rho_{\pi_{\theta_t}}(\tau)}{\rho_{\pi_{\theta}}(\tau)} = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \sum_{h=0}^{H-1} \ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_{\theta}(a_h | s_h)}$$

$$\rho_{\pi}(\tau) = p_0(s_0) \pi(a_0 | s_0) P(s_1 | s_0, a_0) \dots$$

The gradient of the KL-divergence is zero at θ_t

Change from trajectory distribution to state-action distribution:

$$\begin{aligned} KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) &= \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \ln \frac{\rho_{\pi_{\theta_t}}(\tau)}{\rho_{\pi_{\theta}}(\tau)} = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \sum_{h=0}^{H-1} \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \\ &= H \mathbb{E}_{s_h, a_h \sim d_{\mu}^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \right] := \ell(\theta) \end{aligned}$$

The gradient of the KL-divergence is zero at θ_t

Change from trajectory distribution to state-action distribution:

$$KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \ln \frac{\rho_{\pi_{\theta_t}}(\tau)}{\rho_{\pi_{\theta}}(\tau)} = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \sum_{h=0}^{H-1} \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)}$$

$$= H \mathbb{E}_{s_h, a_h \sim d_{\mu}^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \right] := \ell(\theta)$$

Handwritten note: $H \cdot \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \log \frac{\pi_{\theta_t}(a \mid s)}{\pi_{\theta}(a \mid s)}$

$$\nabla_{\theta} \ell(\theta) \big|_{\theta=\theta_t} = H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \left(-\nabla_{\theta} \ln \pi_{\theta}(a_h \mid s_h) \big|_{\theta=\theta_t} \right)$$

The gradient of the KL-divergence is zero at θ_t

Change from trajectory distribution to state-action distribution:

$$\begin{aligned} KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) &= \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \ln \frac{\rho_{\pi_{\theta_t}}(\tau)}{\rho_{\pi_{\theta}}(\tau)} = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \sum_{h=0}^{H-1} \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \\ &= H \mathbb{E}_{s_h, a_h \sim d_{\mu}^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \right] := \ell(\theta) \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} \ell(\theta) \mid_{\theta=\theta_t} &= H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \left(-\nabla_{\theta} \ln \pi_{\theta}(a_h \mid s_h) \mid_{\theta=\theta_t} \right) \\ &= -H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \frac{\nabla_{\theta} \pi_{\theta_t}(a \mid s)}{\pi_{\theta_t}(a \mid s)} \end{aligned}$$

The gradient of the KL-divergence is zero at θ_t

Change from trajectory distribution to state-action distribution:

$$\begin{aligned} KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) &= \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \ln \frac{\rho_{\pi_{\theta_t}}(\tau)}{\rho_{\pi_{\theta}}(\tau)} = \mathbb{E}_{\tau \sim \rho_{\pi_{\theta_t}}} \sum_{h=0}^{H-1} \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \\ &= H \mathbb{E}_{s_h, a_h \sim d_{\mu}^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)} \right] := \ell(\theta) \end{aligned}$$

$$= \nabla \sum_a \pi(a|s)$$

$$= \nabla 1$$

$$= 0$$

$$\nabla_{\theta} \ell(\theta) \big|_{\theta=\theta_t} = H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \left(-\nabla_{\theta} \ln \pi_{\theta}(a_h \mid s_h) \big|_{\theta=\theta_t} \right)$$

$$= -H \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \frac{\nabla_{\theta} \pi_{\theta_t}(a \mid s)}{\pi_{\theta_t}(a \mid s)} = 0$$

Let's compute the Hessian of the KL-divergence at θ_t

$$H \cdot \mathbb{E}_{s,a \sim d_{\mu}^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_{\theta}(a_h | s_h)} \right] := \ell(\theta)$$

Let's compute the Hessian of the KL-divergence at θ_t

$$H \cdot \mathbb{E}_{s,a \sim d_\mu^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_\theta(a_h | s_h)} \right] := \ell(\theta)$$

$$\frac{1}{H} \nabla_\theta^2 \ell(\theta) |_{\theta=\theta_t} = \mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(- \nabla_\theta^2 \ln \pi_\theta(a | s) |_{\theta=\theta_t} \right)$$

Let's compute the Hessian of the KL-divergence at θ_t

$$H \cdot \mathbb{E}_{s,a \sim d_\mu^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_\theta(a_h | s_h)} \right] := \ell(\theta)$$

$$\frac{1}{H} \nabla_\theta^2 \ell(\theta) |_{\theta=\theta_t} = \mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(-\nabla_\theta^2 \ln \pi_\theta(a | s) |_{\theta=\theta_t} \right)$$

$$= - \mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(\frac{\nabla_\theta^2 \pi_{\theta_t}(a | s)}{\pi_{\theta_t}(a | s)} - \frac{\nabla_\theta \pi_{\theta_t}(a | s) \nabla_\theta \pi_{\theta_t}(a | s)^\top}{\pi_{\theta_t}^2(a | s)} \right)$$

$$\sum_a \nabla^2 \pi_\theta(a | s)$$

$$= \nabla^2 \sum_a \pi_\theta(a | s)$$

$$\sum_a \pi_\theta(a | s) = \nabla^2 \mathbf{1} = 0$$

$$\nabla \log \pi_\theta = \frac{\nabla \pi_\theta}{\pi_\theta}$$

$$\nabla^2 \log \pi_\theta = \frac{\nabla^2 \pi_\theta}{\pi_\theta} + \nabla \pi_\theta \nabla \left(\frac{1}{\pi_\theta} \right)$$

Let's compute the Hessian of the KL-divergence at θ_t

$$H \cdot \mathbb{E}_{s,a \sim d_\mu^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_{\theta}(a_h | s_h)} \right] := \ell(\theta)$$

$$\begin{aligned} \frac{1}{H} \nabla_\theta^2 \ell(\theta) |_{\theta=\theta_t} &= \mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(-\nabla_\theta^2 \ln \pi_\theta(a | s) |_{\theta=\theta_t} \right) \\ &= -\mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(\frac{\nabla_\theta^2 \pi_{\theta_t}(a | s)}{\pi_{\theta_t}(a | s)} - \frac{\nabla_\theta \pi_{\theta_t}(a | s) \nabla_\theta \pi_{\theta_t}(a | s)^\top}{\pi_{\theta_t}^2(a | s)} \right) \\ &= \mathbb{E}_{s,a \sim d_\mu^{\pi_{\theta_t}}} \left[\nabla_\theta \ln \pi_{\theta_t}(a | s) \left(\nabla_\theta \ln \pi_{\theta_t}(a | s) \right)^\top \right] \in \mathbb{R}^{dim_\theta \times dim_\theta} \end{aligned}$$

Let's compute the Hessian of the KL-divergence at θ_t

$$H \cdot \mathbb{E}_{s,a \sim d_\mu^{\pi_{\theta_t}}} \left[\ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_{\theta}(a_h | s_h)} \right] := \ell(\theta)$$

$$\begin{aligned} \frac{1}{H} \nabla_{\theta}^2 \ell(\theta) |_{\theta=\theta_t} &= \mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(-\nabla_{\theta}^2 \ln \pi_{\theta}(a | s) |_{\theta=\theta_t} \right) \\ &= -\mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s) \left(\frac{\nabla_{\theta}^2 \pi_{\theta_t}(a | s)}{\pi_{\theta_t}(a | s)} - \frac{\nabla_{\theta} \pi_{\theta_t}(a | s) \nabla_{\theta} \pi_{\theta_t}(a | s)^{\top}}{\pi_{\theta_t}^2(a | s)} \right) \end{aligned}$$

$$= \mathbb{E}_{s,a \sim d_\mu^{\pi_{\theta_t}}} \left[\underbrace{\nabla_{\theta} \ln \pi_{\theta_t}(a | s) \left(\nabla_{\theta} \ln \pi_{\theta_t}(a | s) \right)^{\top}} \right] \in \mathbb{R}^{dim_{\theta} \times dim_{\theta}}$$

It's called fisher Information Matrix!

PSD matrix

$$x^{\top} M x \geq 0$$

Family of dist.

$$\{ \pi(\cdot | s) \mid s \in S \}$$

Summary so far:

We did second-order Taylor expansion on the KL constraint, and we get:

$$\frac{1}{H}KL\left(\rho_{\pi_{\theta_t}}|\rho_{\pi_{\theta}}\right)\approx\frac{1}{2}(\theta-\theta_t)^{\top}F_{\theta_t}(\theta-\theta_t)$$

$$F_{\theta_t}:=\mathbb{E}_{s,a\sim d_{\mu}^{\pi_{\theta_t}}}\left[\nabla_{\theta}\ln\pi_{\theta_t}(a|s)\left(\nabla_{\theta}\ln\pi_{\theta_t}(a|s)\right)^{\top}\right]\in\mathbb{R}^{dim_{\theta}\times dim_{\theta}}$$

Summary so far:

We did second-order Taylor expansion on the KL constraint, and we get:

$$\frac{1}{H}KL\left(\rho_{\pi_{\theta_t}}|\rho_{\pi_{\theta}}\right)\approx\frac{1}{2}(\theta-\theta_t)^{\top}F_{\theta_t}(\theta-\theta_t)$$

$$F_{\theta_t}:=\mathbb{E}_{s,a\sim d_{\mu}^{\pi_{\theta_t}}}\left[\nabla_{\theta}\ln\pi_{\theta_t}(a|s)\left(\nabla_{\theta}\ln\pi_{\theta_t}(a|s)\right)^{\top}\right]\in\mathbb{R}^{dim_{\theta}\times dim_{\theta}}$$

This leads to the following simplified constrained optimization:

Summary so far:

We did second-order Taylor expansion on the KL constraint, and we get:

$$\frac{1}{H} KL \left(\rho_{\pi_{\theta_t}} | \rho_{\pi_{\theta}} \right) \approx \frac{1}{2} (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t)$$

$$F_{\theta_t} := \mathbb{E}_{s, a \sim d_{\mu}^{\pi_{\theta_t}}} \left[\nabla_{\theta} \ln \pi_{\theta_t}(a | s) \left(\nabla_{\theta} \ln \pi_{\theta_t}(a | s) \right)^\top \right] \in \mathbb{R}^{dim_{\theta} \times dim_{\theta}}$$

This leads to the following simplified constrained optimization: “local”

$$\begin{aligned} \max_{\theta} \quad & \nabla_{\theta} J(\pi_{\theta_t})^\top (\theta - \theta_t) \\ \text{s.t.} \quad & (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

← TRPO objective.

Outlines



1. Quick intro on KL-divergence



2. A Trust-Region Formulation for Policy Optimization

3. Algorithm: Natural Policy Gradient

Put everything together, we get:

(dropping the H factors) At iteration t , we update to θ_{t+1} via:

$$\begin{aligned} & \max_{\theta} \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ & \text{s.t. } (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

Put everything together, we get:

(dropping the H factors) At iteration t, we update to θ_{t+1} via:

$$\begin{aligned} \max_{\theta} \quad & \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ \text{s.t.} \quad & (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

Linear objective and quadratic convex constraint, we can solve it optimally!

Put everything together, we get:

(dropping the H factors) At iteration t , we update to θ_{t+1} via:

$$\begin{aligned} \max_{\theta} \quad & \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ \text{s.t.} \quad & (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

substitution \Rightarrow
 $\theta \leftarrow \theta_t$

Linear objective and quadratic convex constraint, we can solve it optimally!

Indeed this gives us:

$$\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$$

Put everything together, we get:

(dropping the H factors) At iteration t, we update to θ_{t+1} via:

$$\begin{aligned} \max_{\theta} \quad & \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ \text{s.t.} \quad & (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

F has units
of
 $\frac{1}{\theta^2}$

Linear objective and quadratic convex constraint, we can solve it optimally!

Indeed this gives us:

$$\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$$

$$\text{Where } \eta = \sqrt{\frac{\delta}{\nabla_{\theta} J(\pi_{\theta_t})^{\top} F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})}}$$

Algorithm: Natural Policy Gradient

Initialize θ_0

For $t = 0, \dots$

Algorithm: Natural Policy Gradient

Initialize θ_0

For $t = 0, \dots$

Estimate PG $\nabla_{\theta} J(\pi_{\theta_t})$

Algorithm: Natural Policy Gradient

Initialize θ_0

For $t = 0, \dots$

Estimate PG $\nabla_{\theta} J(\pi_{\theta_t})$

Estimate Fisher info-matrix $F_{\theta_t} := \mathbb{E}_{s, a \sim d_{\mu}^{\pi_{\theta_t}}} \nabla_{\theta} \ln \pi_{\theta_t}(a | s) (\nabla_{\theta} \ln \pi_{\theta_t}(a | s))^{\top}$

Algorithm: Natural Policy Gradient

Initialize θ_0

For $t = 0, \dots$

Estimate PG $\nabla_{\theta} J(\pi_{\theta_t})$

Estimate Fisher info-matrix $F_{\theta_t} := \mathbb{E}_{s,a \sim d_{\mu}^{\pi_{\theta_t}}} \nabla_{\theta} \ln \pi_{\theta_t}(a | s) (\nabla_{\theta} \ln \pi_{\theta_t}(a | s))^{\top}$

Natural Gradient Ascent: $\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$

Algorithm: Natural Policy Gradient

Initialize θ_0

For $t = 0, \dots$

Estimate PG $\nabla_{\theta} J(\pi_{\theta_t})$

Estimate Fisher info-matrix $F_{\theta_t} := \mathbb{E}_{s,a \sim d_{\mu}^{\pi_{\theta_t}}} \nabla_{\theta} \ln \pi_{\theta_t}(a | s) (\nabla_{\theta} \ln \pi_{\theta_t}(a | s))^{\top}$

Natural Gradient Ascent: $\theta_{t+1} = \theta_t + \eta \hat{F}_{\theta_t}^{-1} \hat{\nabla}_{\theta} J(\pi_{\theta_t})$

$$\text{Where } \eta = \sqrt{\frac{\delta}{\nabla_{\theta} J(\pi_{\theta_t})^{\top} F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})}}$$

Algorithm: Natural Policy Gradient

Initialize θ_0

For $t = 0, \dots$

Estimate PG $\nabla_{\theta} J(\pi_{\theta_t})$

Estimate Fisher info-matrix $F_{\theta_t} := \mathbb{E}_{s, a \sim d_{\mu}^{\pi_{\theta_t}}} \nabla_{\theta} \ln \pi_{\theta_t}(a | s) (\nabla_{\theta} \ln \pi_{\theta_t}(a | s))^{\top}$

Natural Gradient Ascent: $\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$

$$\text{Where } \eta = \sqrt{\frac{\delta}{\nabla_{\theta} J(\pi_{\theta_t})^{\top} F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})}}$$

(We will implement it in HW4 on Cartpole)

Example of Natural Gradient on 1-d problem:

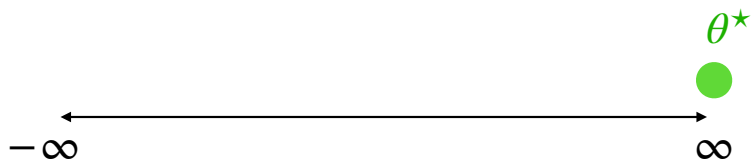
$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$

Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

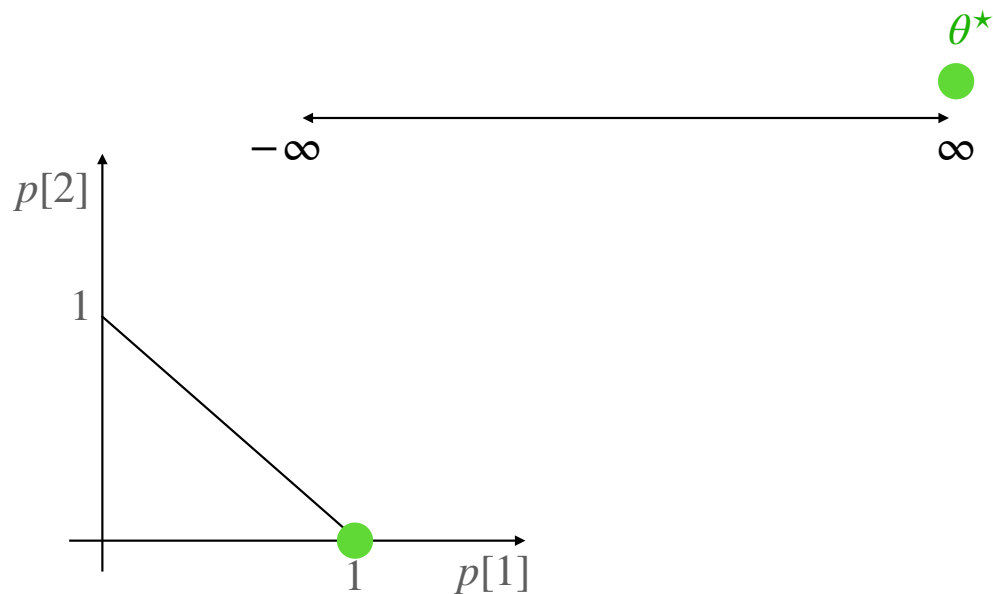
$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$



Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

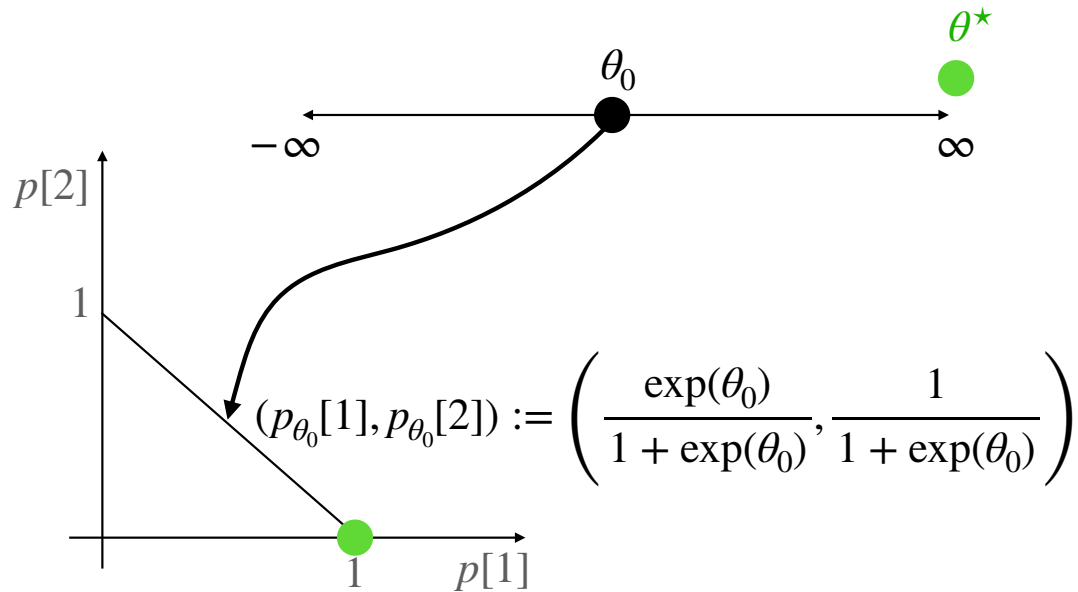
$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$



Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$

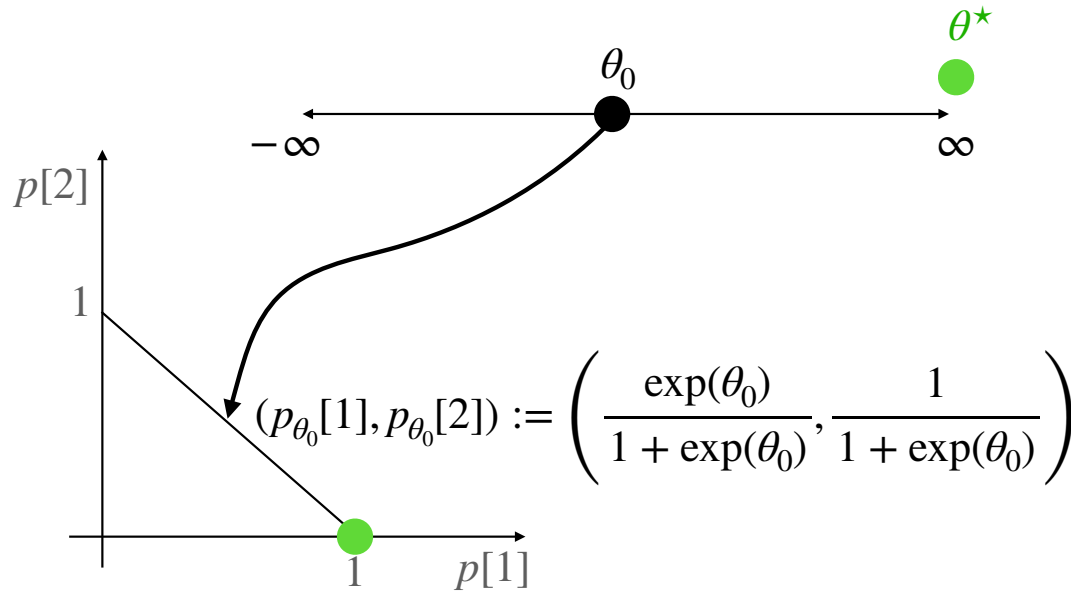


Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$\text{Fisher information scalar: } f_{\theta_0} = \frac{\exp(\theta_0)}{(1 + \exp(\theta_0))^2}$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$



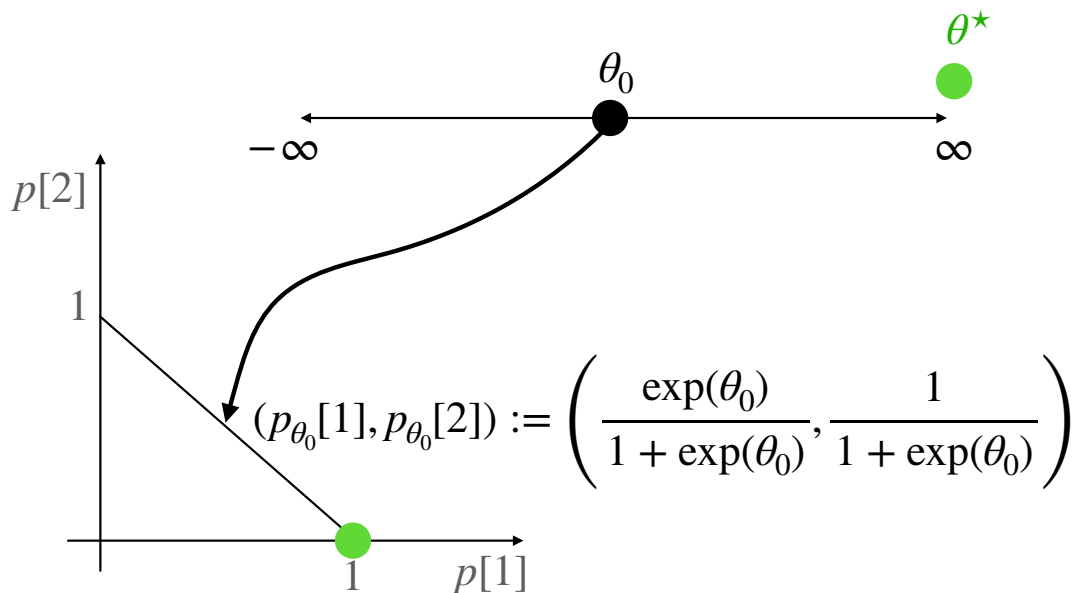
Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$

$$\text{Fisher information scalar: } f_{\theta_0} = \frac{\exp(\theta_0)}{(1 + \exp(\theta_0))^2}$$

Hence: $f_{\theta_0} \rightarrow 0^+$, as $\theta_0 \rightarrow \infty$



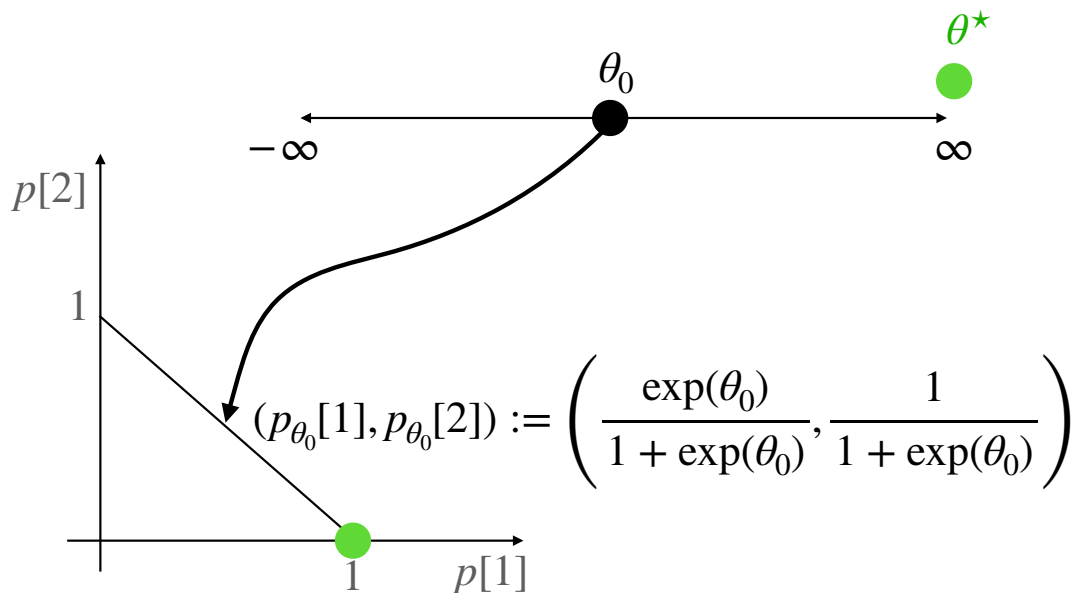
Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$

$$\text{Fisher information scalar: } f_{\theta_0} = \frac{\exp(\theta_0)}{(1 + \exp(\theta_0))^2}$$

Hence: $f_{\theta_0} \rightarrow 0^+$, as $\theta_0 \rightarrow \infty$

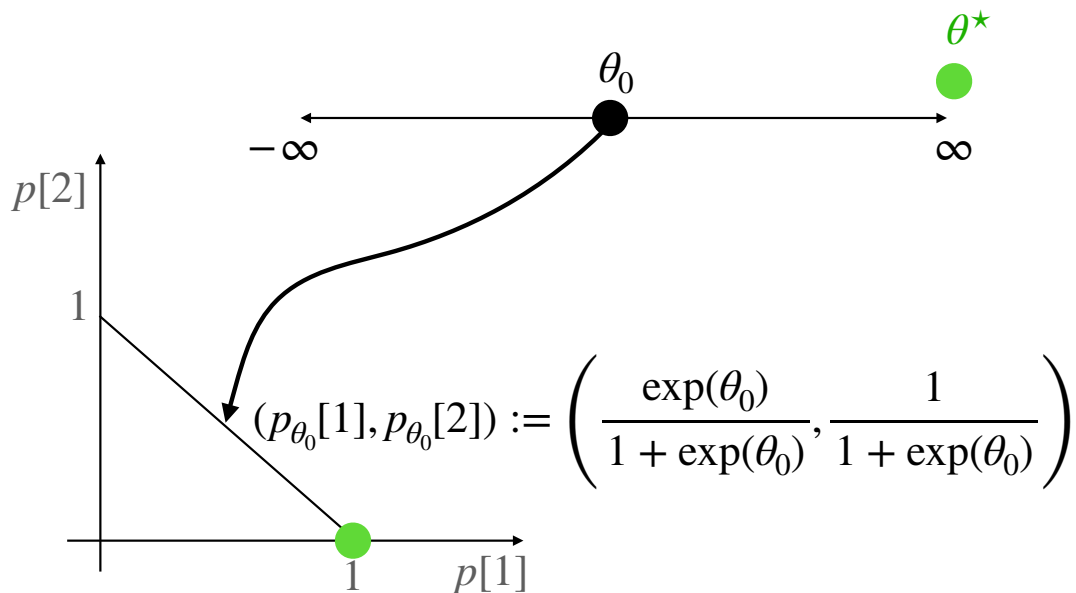


$$\text{NPG: } \theta_1 = \theta_0 + \eta \frac{g'(\theta_0)}{f_{\theta_0}}$$

Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$



Fisher information scalar: $f_{\theta_0} = \frac{\exp(\theta_0)}{(1 + \exp(\theta_0))^2}$

Hence: $f_{\theta_0} \rightarrow 0^+$, as $\theta_0 \rightarrow \infty$

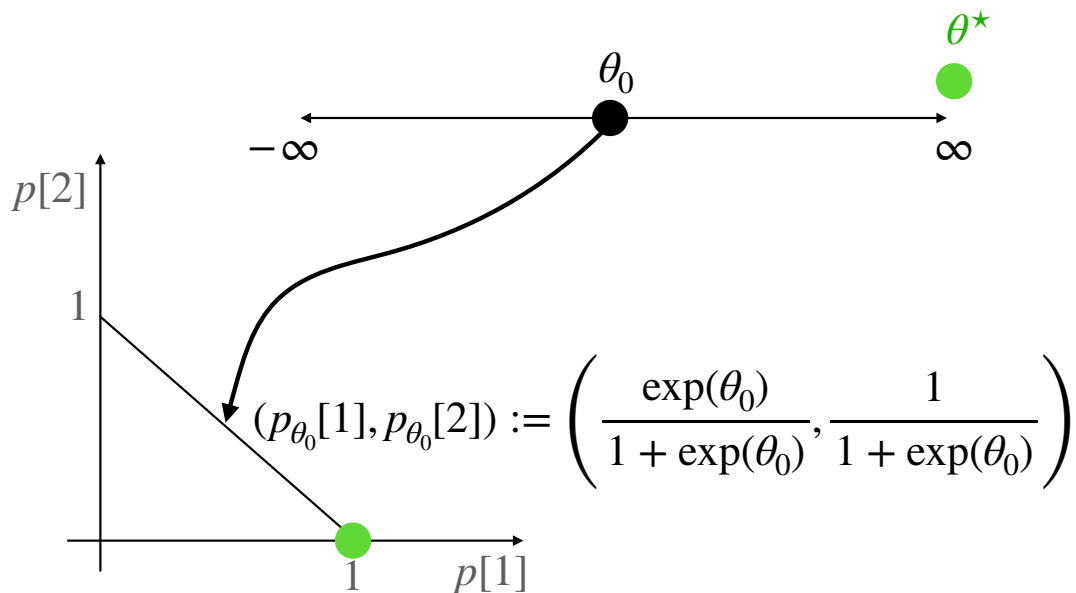
NPG: $\theta_1 = \theta_0 + \eta \frac{g'(\theta_0)}{f_{\theta_0}}$

GA: $\theta_1 = \theta_0 + \eta g'(\theta_0)$

Example of Natural Gradient on 1-d problem:

$$p_{\theta} = \left(\frac{\exp(\theta)}{1 + \exp(\theta)}, \frac{1}{1 + \exp(\theta)} \right)$$

$$g(\theta) = 100 \cdot p_{\theta}[1] + 1 \cdot p_{\theta}[2]$$



$$\text{Fisher information scalar: } f_{\theta_0} = \frac{\exp(\theta_0)}{(1 + \exp(\theta_0))^2}$$

Hence: $f_{\theta_0} \rightarrow 0^+$, as $\theta_0 \rightarrow \infty$

$$\text{NPG: } \theta_1 = \theta_0 + \eta \frac{g'(\theta_0)}{f_{\theta_0}}$$

$$\text{GA: } \theta_1 = \theta_0 + \eta g'(\theta_0)$$

i.e., Plain GA in θ will move to $\theta = \infty$ at a constant speed,
while Natural GA can **traverse faster and faster when θ gets bigger**
(subject to the same learning rate)

Summary for NPG:

Trust Region Policy Optimization and NPG

At iteration t:

$$\begin{aligned} \max_{\pi_{\theta}} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \\ \text{s.t., } KL \left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}} \right) \leq \delta \end{aligned}$$

Intuition: maximize local adv subject
to being incremental (in KL);

Summary for NPG:

Trust Region Policy Optimization and NPG

At iteration t:

$$\begin{aligned} \max_{\pi_{\theta}} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] & \longrightarrow \text{First-order Taylor expansion at } \theta_t \\ \text{s.t., } KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) \leq \delta & \longrightarrow \text{second-order Taylor expansion at } \theta_t \end{aligned}$$

Intuition: maximize local adv subject
to being incremental (in KL);

Summary for NPG:

Trust Region Policy Optimization and NPG

At iteration t:

$$\max_{\pi_{\theta}} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \longrightarrow \text{First-order Taylor expansion at } \theta_t$$

$$\text{s.t., } KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) \leq \delta \longrightarrow \text{second-order Taylor expansion at } \theta_t$$

Intuition: maximize local adv subject
to being incremental (in KL);

$$\begin{aligned} & \max_{\theta} \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ & \text{s.t. } (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

Summary for NPG:

Trust Region Policy Optimization and NPG

At iteration t:

$$\max_{\pi_{\theta}} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \longrightarrow \text{First-order Taylor expansion at } \theta_t$$

$$\text{s.t., } KL(\rho_{\pi_{\theta_t}} | \rho_{\pi_{\theta}}) \leq \delta \longrightarrow \text{second-order Taylor expansion at } \theta_t$$

Intuition: maximize local adv subject
to being incremental (in KL);

$$\begin{aligned} & \max_{\theta} \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ & \text{s.t. } (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

(Exercise: work out the $\arg \max_{\theta}$)

Summary for NPG:

Trust Region Policy Optimization and NPG

At iteration t:

$$\max_{\pi_{\theta}} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(s)} A^{\pi_{\theta_t}}(s, a) \right] \longrightarrow \text{First-order Taylor expansion at } \theta_t$$

$$\text{s.t., } KL\left(\rho_{\pi_{\theta_t}} \mid \rho_{\pi_{\theta}}\right) \leq \delta \longrightarrow \text{second-order Taylor expansion at } \theta_t$$

Intuition: maximize local adv subject
to being incremental (in KL);

$$\theta_{t+1} = \theta_t + \eta F_{\theta_t}^{-1} \nabla_{\theta} J(\pi_{\theta_t})$$

NPG

$$\begin{aligned} & \max_{\theta} \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta - \theta_t) \\ & \text{s.t. } (\theta - \theta_t)^{\top} F_{\theta_t} (\theta - \theta_t) \leq \delta \end{aligned}$$

(Exercise: work out the $\arg \max_{\theta}$)

An extension of NPG (even faster in practice):

Given an current policy π^t , we perform policy update to π^{t+1}

fifth attempt (new): **Proximal Policy Optimization (PPO)**

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} A^{\pi_{\theta_t}}(s, a) \right]$$

An extension of NPG (even faster in practice):

Given an current policy π^t , we perform policy update to π^{t+1}

fifth attempt (new): **Proximal Policy Optimization (PPO)**

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} A^{\pi_{\theta_t}}(s, a) \right] - \underbrace{\lambda \mathbb{E}_{s \sim d_{\mu}^{\pi_t}} \left[\text{KL} \left(\pi_{\theta_t}(a | s) \mid \pi_{\theta}(a | s) \right) \right]}_{\text{regularization}}$$

An extension of NPG (even faster in practice):

Given an current policy π^t , we perform policy update to π^{t+1}

fifth attempt (new): **Proximal Policy Optimization (PPO)**

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} A^{\pi_{\theta_t}}(s, a) \right] - \underbrace{\lambda \mathbb{E}_{s \sim d_{\mu}^{\pi_t}} \left[\text{KL} \left(\pi_{\theta_t}(a | s) \mid \pi_{\theta}(a | s) \right) \right]}_{\text{regularization}}$$

Use importance weighting & expand KL divergence:

An extension of NPG (even faster in practice):

Given an current policy π^t , we perform policy update to π^{t+1}

fifth attempt (new): **Proximal Policy Optimization (PPO)**

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} A^{\pi_{\theta_t}}(s, a) \right] - \underbrace{\lambda \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\text{KL} \left(\pi_{\theta_t}(a | s) \mid \pi_{\theta}(a | s) \right) \right]}_{\text{regularization}}$$

Use importance weighting & expand KL divergence:

$$\ell(\theta) := \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta_t}(\cdot | s)} \frac{\pi_{\theta}(a | s)}{\pi_{\theta_t}(a | s)} A^{\pi_{\theta_t}}(s, a) \right] - \lambda \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \mathbb{E}_{a \sim \pi_{\theta_t}(\cdot | s)} \left[-\ln \pi_{\theta}(a | s) \right]$$

An extension of NPG (even faster in practice):

Given an current policy π^t , we perform policy update to π^{t+1}

fifth attempt (new): **Proximal Policy Optimization (PPO)**

$$\max_{\theta} \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)} A^{\pi_{\theta_t}}(s, a) \right] - \underbrace{\lambda \mathbb{E}_{s \sim d_{\mu}^{\pi^t}} \left[\text{KL} \left(\pi_{\theta_t}(a | s) \mid \pi_{\theta}(a | s) \right) \right]}_{\text{regularization}}$$

Use importance weighting & expand KL divergence:

$$\ell(\theta) := \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \left[\mathbb{E}_{a \sim \pi_{\theta_t}(\cdot | s)} \frac{\pi_{\theta}(a | s)}{\pi_{\theta_t}(a | s)} A^{\pi_{\theta_t}}(s, a) \right] - \lambda \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \mathbb{E}_{a \sim \pi_{\theta_t}(\cdot | s)} \left[-\ln \pi_{\theta}(a | s) \right]$$

PPO: Perform a few steps of mini-batch SGA on $\ell(\theta)$ to approximate $\arg \max_{\theta} \ell(\theta)$

Next a few lectures:

**Imitation Learning
(Learning from Demonstrations)**

Can we learn a good policy purely from expert demonstrations?

Summary:

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

1-minute feedback form: <https://bit.ly/3RHtlxy>

