# Trust Region
# Policy Optimization

# Announcements

Thanks for providing midterm feedback!

1. HW2 will be out this Friday

2. I will have an additional office hour every Monday morning
(11am - noon)

# Recap Policy Gradient

$$J(\pi_\theta) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 \sim \mu, a \sim \pi_\theta\right]$$

# Recap Policy Gradient

$$J(\pi_\theta) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,|\, s_0 \sim \mu, a \sim \pi_\theta\right]$$

The most commonly used formulation:

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

# Recap Policy Gradient

$$J(\pi_\theta) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,|\, s_0 \sim \mu, a \sim \pi_\theta\right]$$

The most commonly used formulation:

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

Algorithm: Stochastic Gradient Ascent

# Recap on Conservative Policy Iteration

For t = 0 …

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. Incremental Update:
$$\pi^{t+1}( \cdot \mid s) = (1 - \alpha)\pi^t( \cdot \mid s) + \alpha\pi'( \cdot \mid s), \forall s$$

# Recap on Conservative Policy Iteration

For t = 0 …

    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. Incremental Update:

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

Q: Why this is incremental? In what sense?

Q: Can we get monotonic policy improvement?

# Recap of CPI:

Incremental update (Lemma 12.1 in AJKS)

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

# Pros and Cons of CPI:

Pros:
**This is fundamental!**
The idea of incremental update and the theorem behind it are still being used today…

Cons:
**Practical Issue (e.g., memory issue)**
e.g., what if my policies are all extremely large neural networks…

# Today's Question

Can we develop some practical version of CPI?

# Outlines

1. Quick intro on KL-divergence

2. A Trust-Region Formulation for Policy Optimization

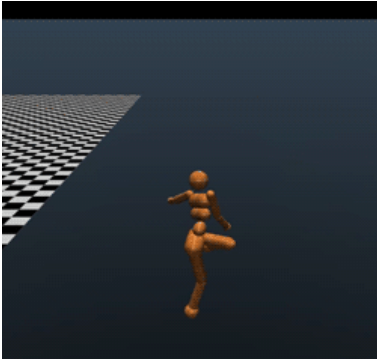3. Algorithm: Natural Policy Gradient

# Interesting videos from the today's algorithm

**Train a robot to "run" forward as fast as possible:**
**State**: joint angles, center of mass, velocity, etc
**Action**: torques on joints
**Reward**: distance of moving forward between two steps

# Interesting videos from the today's algorithm

**Train a robot to "run" forward as fast as possible:**
**State**: joint angles, center of mass, velocity, etc
**Action**: torques on joints
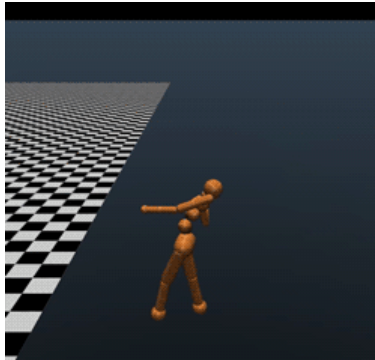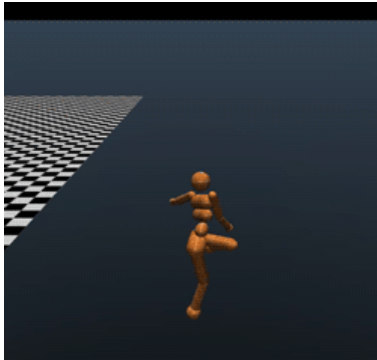**Reward**: distance of moving forward between two steps

# Interesting videos from the today's algorithm

**Train a robot to "run" forward as fast as possible:**
**State**: joint angles, center of mass, velocity, etc
**Action**: torques on joints
**Reward**: distance of moving forward between two steps

# Interesting videos from the today's algorithm

**Train a robot to "run" forward as fast as possible:**
**State**: joint angles, center of mass, velocity, etc
**Action**: torques on joints
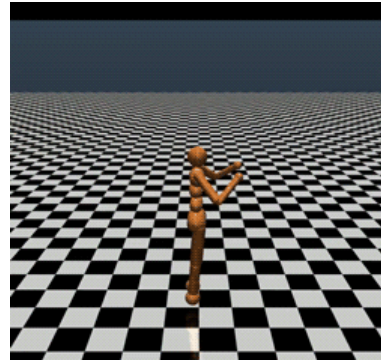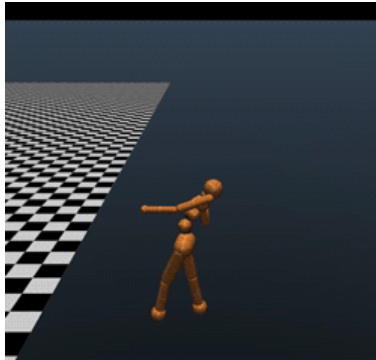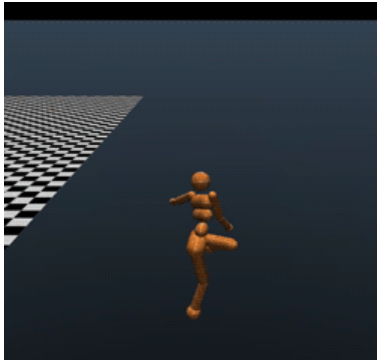**Reward**: distance of moving forward between two steps

# Interesting videos from the today's algorithm

**Train a robot to "run" forward as fast as possible:**
**State**: joint angles, center of mass, velocity, etc
**Action**: torques on joints
**Reward**: distance of moving forward between two steps



(BTW, This reveals an issue on reward design—we will study it in Learning from Demonstrations)

# 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:

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

# 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:

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

$\forall x$
$P(x) = Q(x)$

**Examples:**

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

# 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:

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

$$KL(P\,|\,Q)$$
$$\neq KL(Q\,|\,P)$$

**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 $KL(P \,|\, Q) = \|\mu_1 - \mu_2\|_2^2 / \sigma^2$  ✓

# 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:

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

**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:**

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

$P(x)$.

Entropy, $\mathbb{E}_{x \sim P} - \log_2 P(x)$

$\mathbb{E}_{x \sim P}\left[ f\left(\frac{P(x)}{Q(x)}\right)\right]$

$f:$

$= \mathbb{E}_{x \sim P}[\ln P(x)]$
$- \mathbb{E}_{x \sim P} \ln Q(x)$

$\int_x P(x) = 1$

$\int_x Q(x) = 1$

$\frac{1}{2}\|P - Q\|_1 \leq \sqrt{KL(P\|Q)}$

# Outlines

✅ 1. Quick intro on KL-divergence

2. A Trust-Region Formulation for Policy Optimization

3. Algorithm: Natural Policy Gradient

# Policy Parameterization

Recall that we consider parameterized policy $\pi_\theta(\,\cdot\,|\,s) \in \Delta(A), \forall s$

### 1. Softmax linear Policy (We will try this in HW2)

Feature vector $\phi(s, a) \in \mathbb{R}^d$, and parameter $\theta \in \mathbb{R}^d$

$$\pi_\theta(a\,|\,s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$$

### 2. Neural Policy:

Neural network
$f_\theta : S \times A \mapsto \mathbb{R}$

$$\pi_\theta(a\,|\,s) = \frac{\exp(f_\theta(s, a))}{\sum_{a'} \exp(f_\theta(s, a'))}$$

# A trust region formulation for policy update:

At iteration t, with $\pi_{\theta_t}$ at hand, we compute $\theta_{t+1}$ as follows:

# A trust region formulation for policy update:

At iteration t, with $\pi_{\theta_t}$ at hand, we compute $\theta_{t+1}$ as follows:

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

$$\text{s.t., } KL\left( \rho_{\pi_{\theta_t}} \,|\, \rho_{\pi_\theta} \right) \leq \delta \quad \leftarrow \text{hyper-prameter}$$

$\delta$

$P_{\pi_{\theta_t}}$

Truse Region

Trejy-Distnbutiby

$\rho(\tau) = \mu(s_0) \, \Pi(a_0 | s_0) \, P(s_1 | s_0, a_0) \; - \; - \; -$

**A trust region formulation for policy update:**

At iteration t, with $\pi_{\theta_t}$ at hand, we compute $\theta_{t+1}$ as follows:

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

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

We want to maximize local advantage against $\pi_{\theta_t}$, but we want the new policy to be close to $\pi_{\theta_t}$ (in the KL sense)

# A trust region formulation for policy update:

At iteration t, with $\pi_{\theta_t}$ at hand, we compute $\theta_{t+1}$ as follows:

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

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

We want to maximize local advantage against $\pi_{\theta_t}$, but we want
the new policy to be close to $\pi_{\theta_t}$ (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…

# A trust region formulation for policy update:

At iteration t, with $\pi_{\theta_t}$ at hand, we compute $\theta_{t+1}$ as follows:

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

$$\text{s.t., } KL \left( \rho_{\pi_{\theta_t}} | \rho_{\pi_\theta} \right) \leq \delta \quad \leftarrow \text{second-order}$$
$$\text{Taylor - Exp at } \theta_t$$

High-level strategy:
1. First-order Taylor expansion on the objective at $\theta_t$
2. second-order Taylor expansion of the constraint at $\theta_t$

## Simplify Objective Function

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

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

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

# 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-talyor expansion on it:

# 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-talyor expansion on it:

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

$= 0$

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

Inner product

# Simplify Objective Function

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

$$\underbrace{\phantom{\nabla_\theta J(\theta_t)}}_{\mathcal{A}}$$

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

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

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

**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:

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

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

First

Hessian

# 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!

$$KL(\rho_{\theta_t} | \rho_\theta) \approx \frac{1}{2}(\theta - \theta_t)^\top \nabla_\theta^2 \ell(\theta_t)(\theta - \theta_t)$$

# The gradient of the KL-divergence is zero at $\theta_t$

Change from trajectory distribution to state-action distribution:

# The gradient of the KL-divergence is zero at $\theta_t$

Change from trajectory distribution to state-action distribution:

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

*KL-DIV*

$$\frac{\rho_{\pi_{\theta_t}}(\tau)}{\rho_{\pi_\theta}(\tau)} = \frac{\mu_0(s_0)\,\pi_{\theta_t}(a_0 \mid s_0)\, P_1(s_1 \mid s_0, a_0) \cdots}{\mu_0(s_0)\,\pi_\theta(a_0 \mid s)\, P_1(s_1 \mid s_0, a_0) \cdots}$$

$$\ln\left[\frac{\rho_{\pi_{\theta_t}}(t)}{\rho_{\pi_\theta}(t)}\right] = \sum_{h=0}^{\infty} \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_{\theta}(a_h \mid s_h)}$$

# The gradient of the KL-divergence is zero at $\theta_t$

Change from trajectory distribution to state-action distribution:

$\leftarrow$ finite Horizon setting

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

$$= \frac{1}{1-\gamma} \cdot H \cdot \mathbb{E}_{s_h, a_h \sim d_\mu^{\pi_{\theta_t}}} \left[ \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_\theta(a_h \mid s_h)} \right] := \ell(\theta)$$

# The gradient of the KL-divergence is zero at $\theta_t$

Change from trajectory distribution to state-action distribution:

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

$$= \frac{\cancel{1}^{\,H}}{\cancel{1-\gamma}} \mathbb{E}_{s_h, a_h \sim d_\mu^{\pi_{\theta_t}}} \left[ \ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_\theta(a_h | s_h)} \right] := \ell(\theta)$$

$$\ln\left(\tfrac{\pi_{\theta_t}}{\pi_\theta}\right) = \ln \pi_{\theta_t} - \ln \pi_\theta$$

$$\nabla_\theta \ell(\theta)|_{\theta=\theta_t} = \mathbb{E}_{s \sim d_\mu^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a | s)\left(-\nabla_\theta \ln \pi_\theta(a_h | s_h)|_{\theta=\theta_t}\right) = 0$$

# The gradient of the KL-divergence is zero at $\theta_t$

Change from trajectory distribution to state-action distribution:

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

$$= \frac{1}{1-\gamma} \mathbb{E}_{s_h, a_h \sim d_\mu^{\pi_{\theta_t}}} \left[ \ln \frac{\pi_{\theta_t}(a_h \mid s_h)}{\pi_\theta(a_h \mid s_h)} \right] := \ell(\theta)$$

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

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

$$= -\mathbb{E}_s \nabla_\theta \cdot 1 = 0$$

# The gradient of the KL-divergence is zero at $\theta_t$

Change from trajectory distribution to state-action distribution:

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

$$= \frac{1}{1-\gamma} \mathbb{E}_{s_h, a_h \sim d_\mu^{\pi_{\theta_t}}} \left[ \ln \frac{\pi_{\theta_t}(a_h | s_h)}{\pi_\theta(a_h | s_h)} \right] := \ell(\theta)$$

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

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

# Let's compute the Hessian of the KL-divergence at $\theta_t$

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

# Let's compute the Hessian of the KL-divergence at $\theta_t$

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

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

$\ln \pi_{\theta_t} - \ln \pi_\theta$

$$\left( \frac{f}{g} \right)' = \frac{f'}{g} - \frac{f \cdot g'}{g^2}$$

$$\nabla_\theta \ln \pi_\theta(a \mid s) = \frac{\nabla_\theta \pi_\theta(a \mid s)}{\pi_\theta(a \mid s)}$$

$$\nabla_\theta \left[ \nabla_\theta \ln \pi_\theta(a \mid s) \right] = \nabla_\theta \left[ \frac{\nabla_\theta \pi_\theta(a \mid s)}{\pi_\theta(a \mid s)} \right]$$

$$= \frac{\nabla_\theta^2 \pi_\theta(a \mid s)}{\pi_\theta(a \mid s)} - \frac{\nabla_\theta \pi_\theta(a \mid s) \cdot \nabla_\theta \pi(a \mid s)^\top}{\pi_\theta^2(a \mid s)}$$

# Let's compute the Hessian of the KL-divergence at $\theta_t$

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

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

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

$$= 0$$

$$\sum_a \pi_{\theta_t}(a\,|\,s)\,\frac{\nabla_\theta^2\pi_{\theta_t}(a\,|\,s)}{\pi_{\theta_t}(a\,|\,s)} = \nabla_\theta^2\left[\sum_a \pi_{\theta_t}(a\,|\,s)\right] = \nabla_\theta^2\,1 = 0$$

# Let's compute the Hessian of the KL-divergence at $\theta_t$

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

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

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

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

# Let's compute the Hessian of the KL-divergence at $\theta_t$

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

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

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

$$\frac{\nabla_\theta \pi_{\theta_t}(a \mid s)}{\pi_{\theta_t}(a \mid s)} \left[ \frac{\nabla_\theta \pi_{\theta_t}(a \mid s)}{\pi_{\theta_t}(a \mid s)} \right]^\top$$

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

It's called fisher Information Matrix!

# 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 much simplified constrained optimization:

# Summary so far:

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

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

$x\,x^\top$

↑ PSD

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

This leads to the following much simplified constrained optimization:

PG

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

$$\text{s.t. } (\theta - \theta_t)^\top F_{\theta_t}(\theta - \theta_t) \leq \delta$$

PSD

# Outlines

✅ 1. Quick intro on KL-divergence

✅ 2. A Trust-Region Formulation for Policy Optimization

$$\underset{s\,a}{E}\left[\nabla_\theta \ln \pi(s|s)\,\nabla_\theta \ln \pi(a|s)^T\right]$$

is PSD

3. Algorithm: Natural Policy Gradient

**Put everything together, we get:**

At iteration t, we update to $\theta_{t+1}$ via:

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

$$\text{s.t. } (\theta - \theta_t)^{\top} F_{\theta_t}(\theta - \theta_t) \leq \delta$$

**Put everything together, we get:**

At iteration t, we update to $\theta_{t+1}$ via:

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

$$\text{s.t. } (\theta - \theta_t)^{\top} F_{\theta_t}(\theta - \theta_t) \leq \delta$$

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

# Put everything together, we get:

At iteration t, we update to $\theta_{t+1}$ via:

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

$$\text{s.t. } (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t) \leq \delta$$

invertible

$F_{\theta_t} + \lambda I$

$\uparrow$

$1e^{-7}$

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

Indeed this gives us:

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

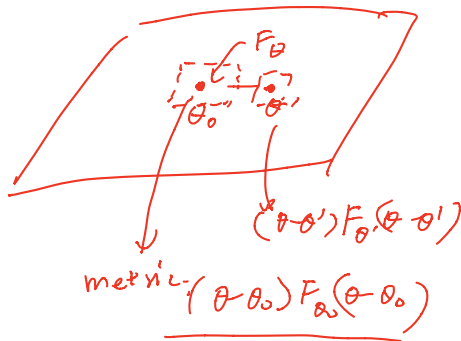# Put everything together, we get:

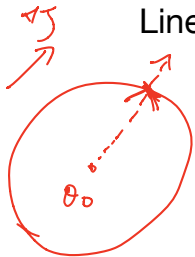At iteration t, we update to $\theta_{t+1}$ via:

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

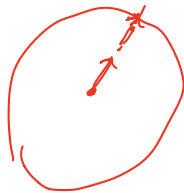$$\text{s.t. } (\theta - \theta_t)^\top F_{\theta_t}(\theta - \theta_t) \leq \delta$$

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

Indeed this gives us:

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

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

(handwritten annotations:)

$\max_{\theta} \nabla J^\top (\theta - \theta_v)$

$(\theta - \theta_o) // \nabla J$

$F_\theta$

$(\theta - \theta') F_\theta (\theta - \theta')$

metric $(\theta - \theta_o) F_\theta (\theta - \theta_o)$

$\Delta$

# Algorithm: Natural Policy Gradient

Initialize $\theta_0$

For t = 0, …

# Algorithm: Natural Policy Gradient

Initialize $\theta_0$

For t = 0, …

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

# Algorithm: Natural Policy Gradient

Initialize $\theta_0$

For t = 0, …

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

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

*finite # of samples*

# Algorithm: Natural Policy Gradient

Initialize $\theta_0$

For t = 0, …

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

Natural Gradient

# Algorithm: Natural Policy Gradient

Initialize $\theta_0$

For t = 0, …

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

$$\nabla_\theta J(\pi_\theta) \Big|_{\theta = \theta_t}$$

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

Tune $\quad 1e^{-2}$

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

# Algorithm: Natural Policy Gradient

Initialize $\theta_0$

For t = 0, …

    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 HW2 on Cartpole)

# Summary for today:

Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

# Summary for today:

Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

# Summary for today:

Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

$$\max_\theta \nabla_\theta J(\pi_{\theta_t})^\top (\theta - \theta_t)$$
$$\text{s.t. } (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t) \le \delta$$

symmetric

# Summary for today:

Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

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

$$\text{s.t. } (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t) \leq \delta$$

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

# Summary for today:

Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

$$\theta_{t+1} = \theta_t + \eta \underbrace{F_{\theta_t}^{-1} \nabla_\theta J(\pi_{\theta_t})}_{\text{NPG}} \longleftarrow \begin{array}{l} \max_\theta \nabla_\theta J(\pi_{\theta_t})^\top (\theta - \theta_t) \\ \text{s.t. } (\theta - \theta_t)^\top F_{\theta_t} (\theta - \theta_t) \leq \delta \end{array}$$

(Exercise: work out the arg max)
$$\theta$$

$$F^{-1} g$$

$$F = \begin{bmatrix} \sigma_1 & \sigma_2 & \\ & & \ddots & \\ & & & \boxed{\sigma_d} \end{bmatrix} \underset{\sim 0}{}$$

$$= \begin{bmatrix} \frac{1}{\sigma_1} g_1 \\ \vdots \\ \boxed{\frac{1}{\sigma_d} g_d} \end{bmatrix}$$

$$\sigma_d = 0.1$$

$$KL(P|Q) = \mathop{E}_{x \sim P} \ln \frac{P(x)}{Q(x)}$$

$$KL(U|Q) = \mathop{E}_{x \sim U} \ln \frac{U(x)}{Q(x)}$$

$$= \underline{\mathop{E}_{x \sim U} \ln U(x)} - \underline{\mathop{E}_{x \sim U} \ln Q(x)}$$

$$KL(Q|U) = \mathop{E}_{x \sim Q} \ln Q(x) - \boxed{\underline{\mathop{E}_{x \sim Q} \ln U(x)}}$$

$$= \mathop{E}_{x \sim Q} \ln Q(x) - \ln \frac{1}{|X|}$$

$$KL(P|Q) = \mathop{E}_{x \sim P} \ln \frac{P(x)}{Q(x)}$$



$Q$     $P$

$x^*$     $x$