

# **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

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

The most commonly used formulation:

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

## Recap Policy Gradient

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

The most commonly used formulation:

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

Algorithm: Stochastic Gradient Ascent

# Recap on Conservative Policy Iteration

For  $t = 0 \dots$

1. Greedy Policy Selector:

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

2. Incremental Update:

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

# Recap on Conservative Policy Iteration

For  $t = 0 \dots$

1. Greedy Policy Selector:

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

2. Incremental Update:

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

Q: Why this is incremental? In what sense?

Q: Can we get monotonic policy improvement?

## 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

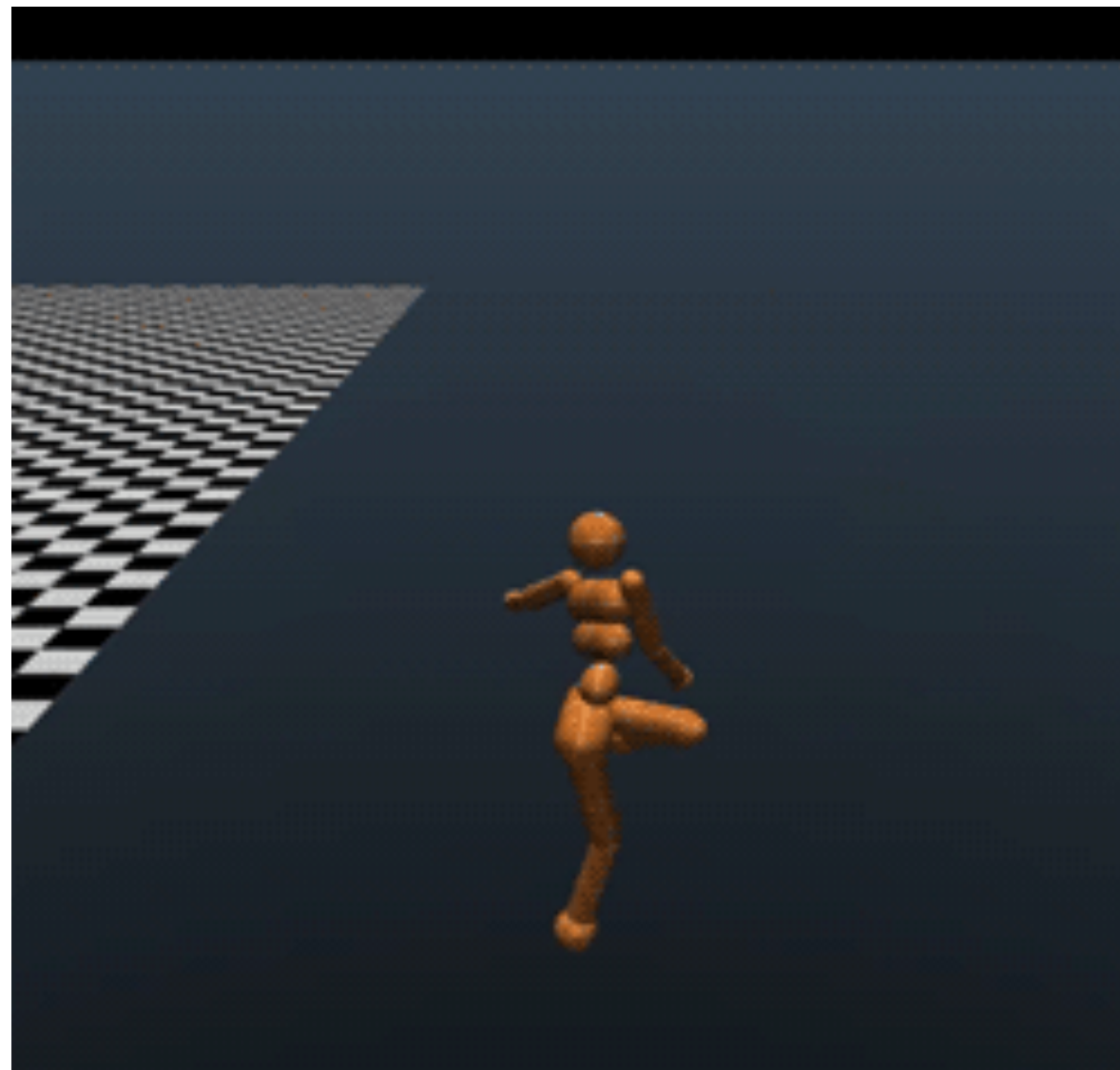
# Interesting videos from the today's algorithm

**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



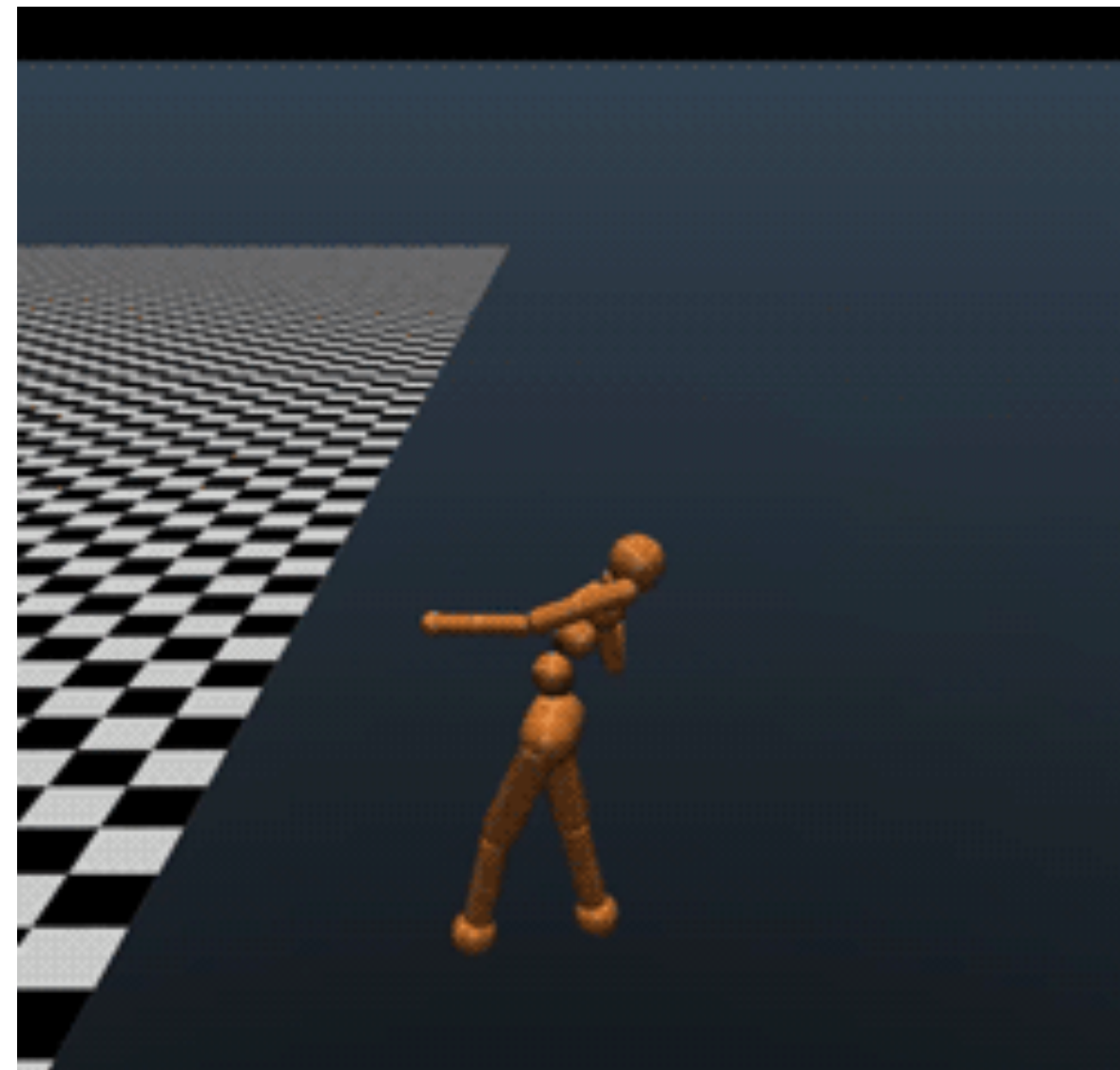
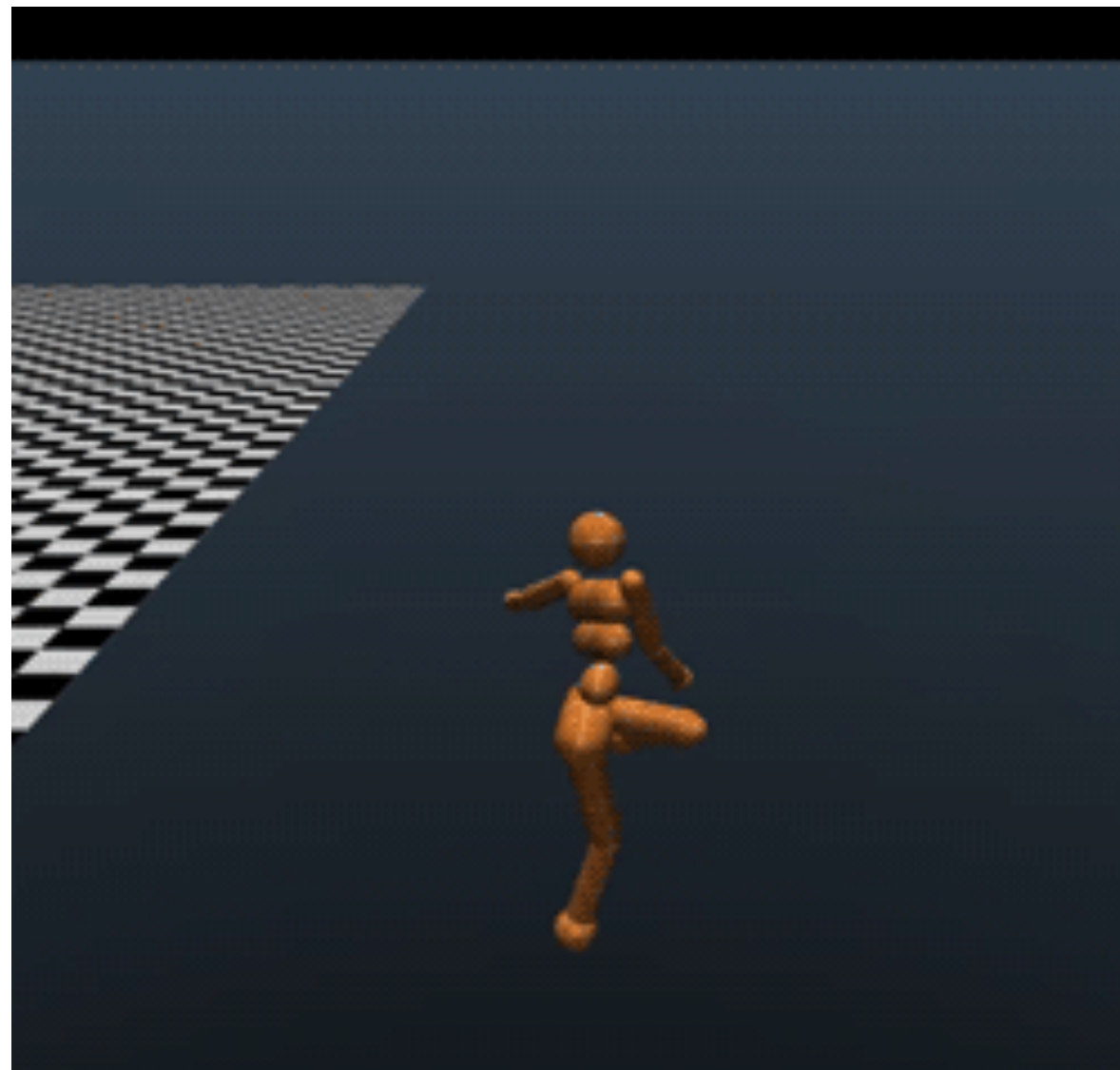
# Interesting videos from the today's algorithm

**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



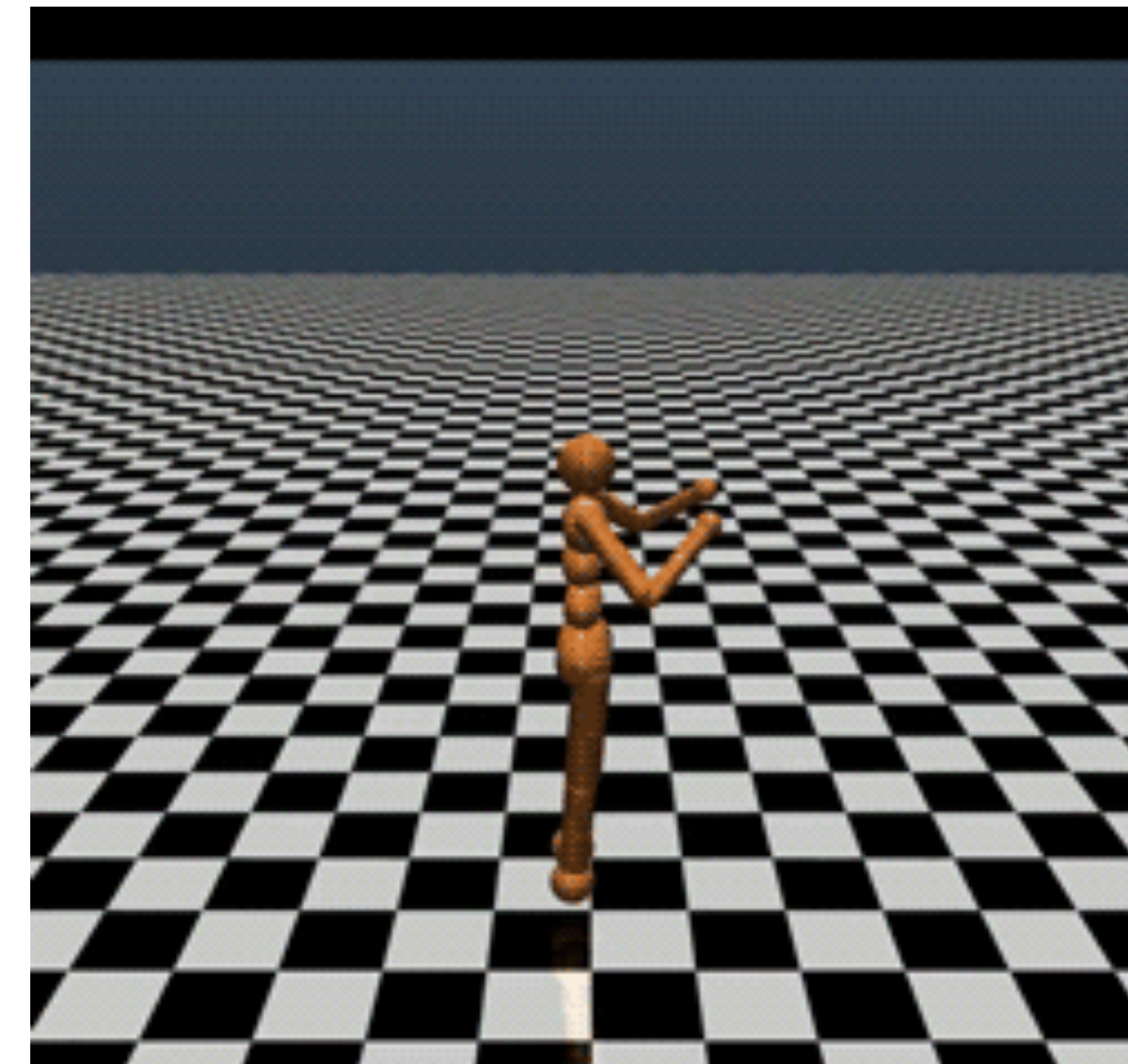
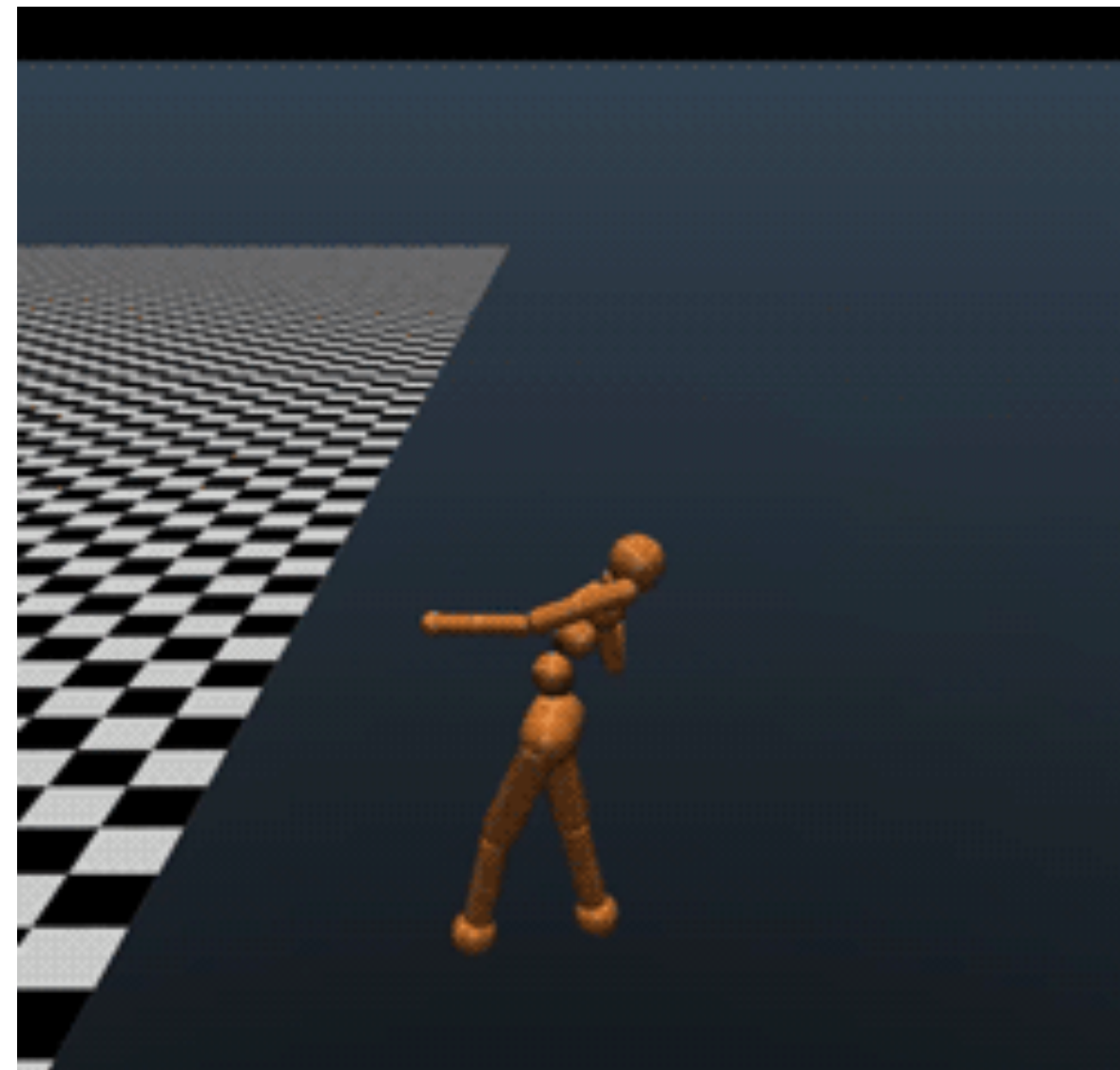
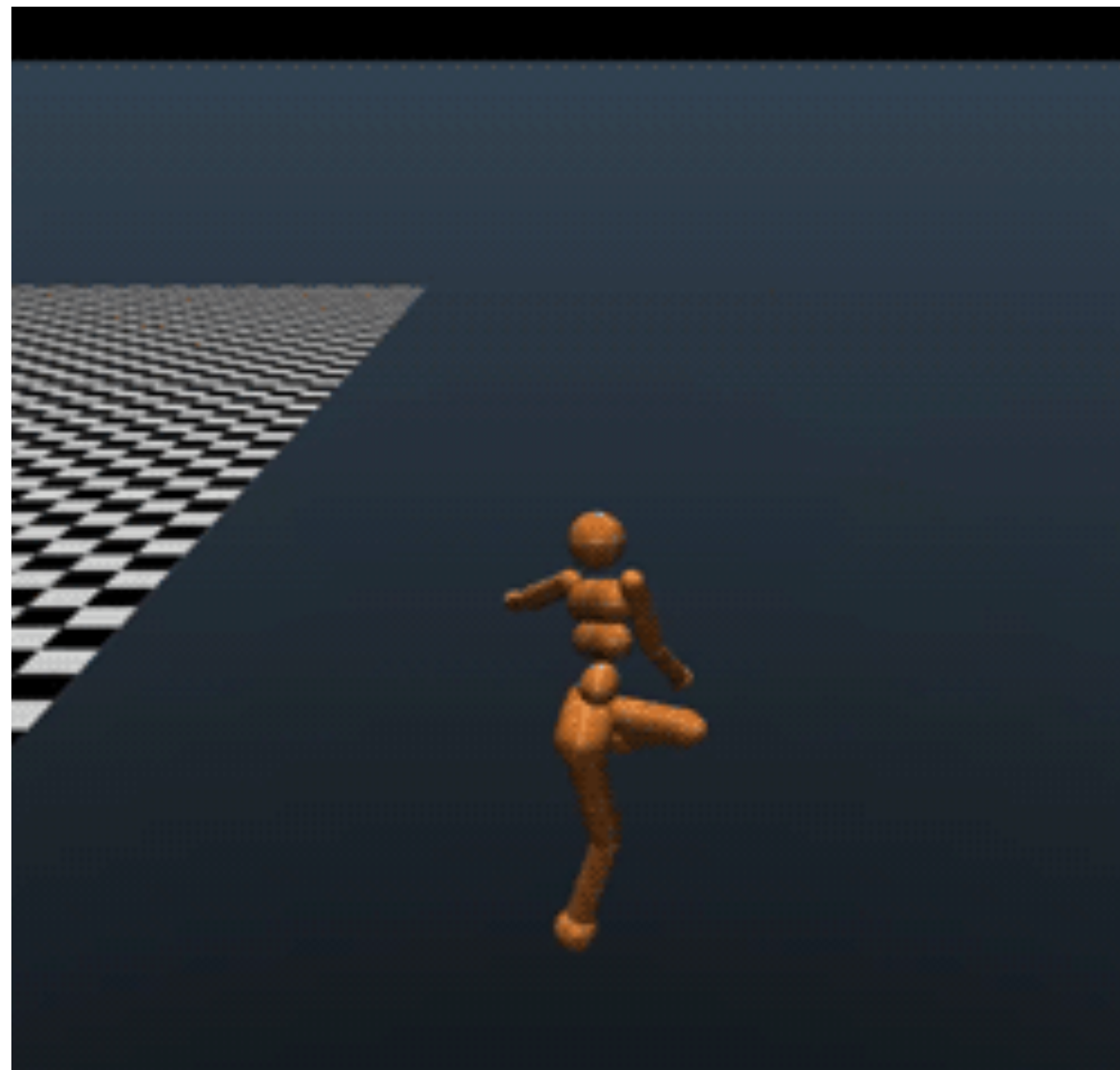
# Interesting videos from the today's algorithm

**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



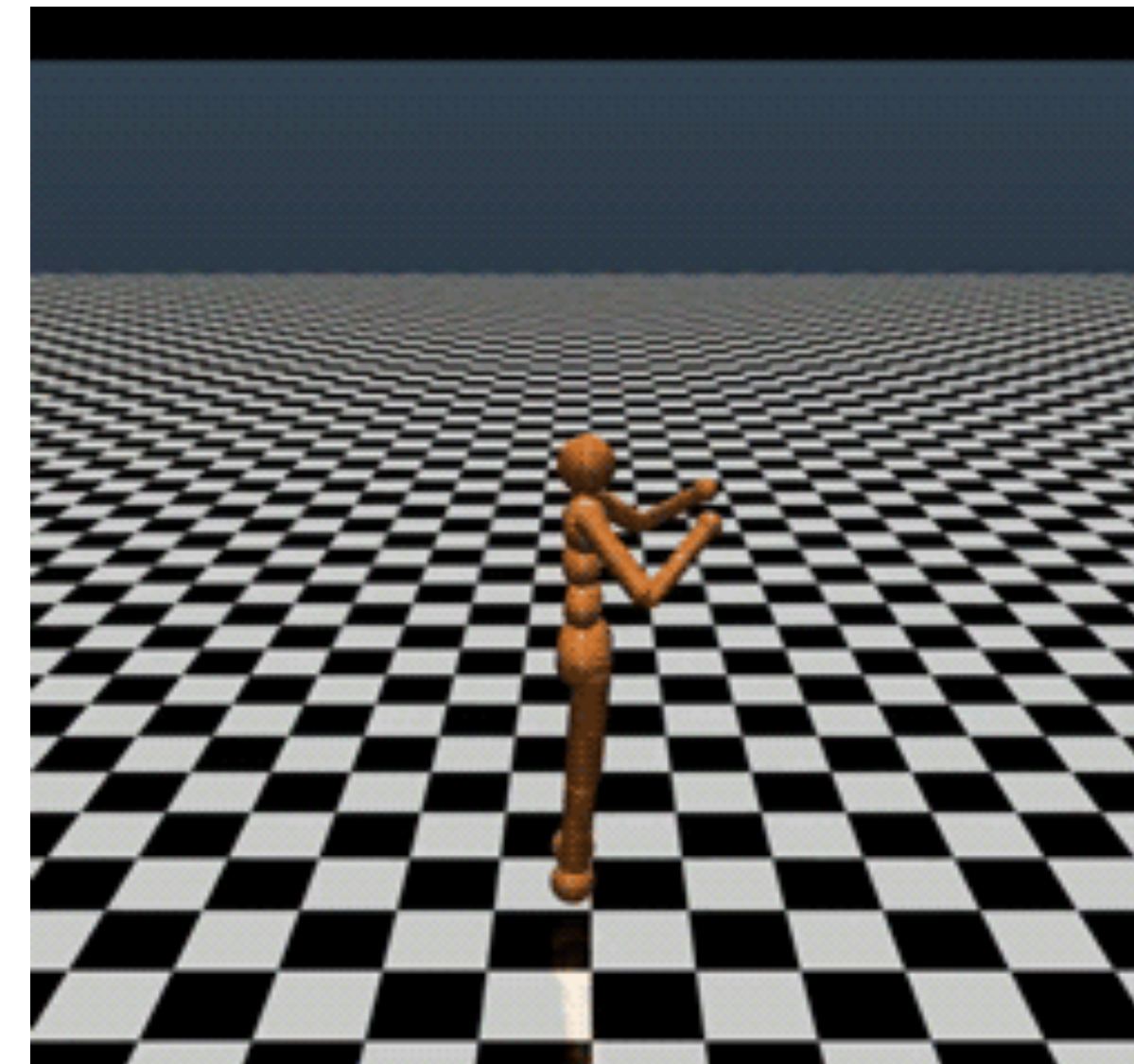
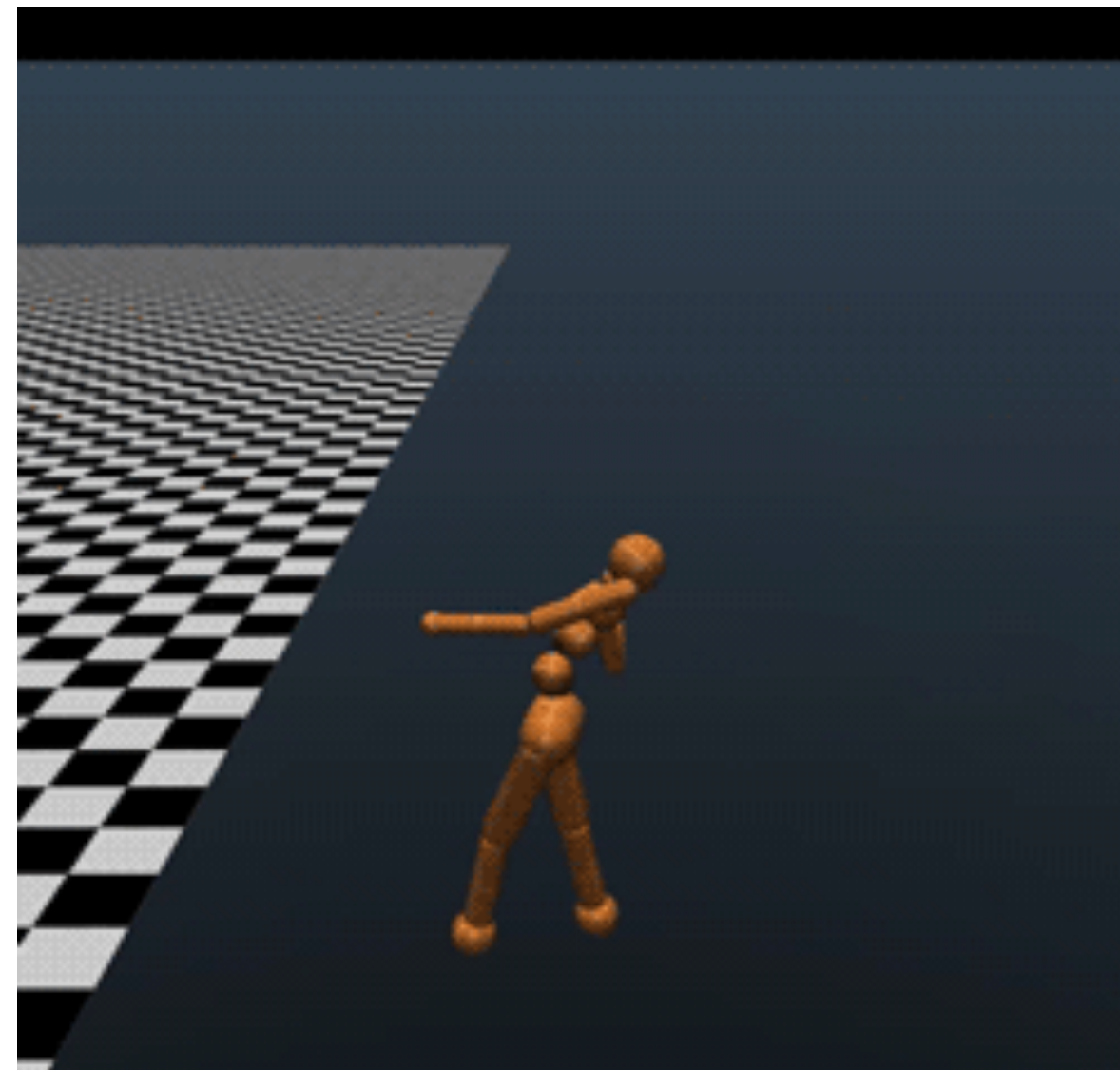
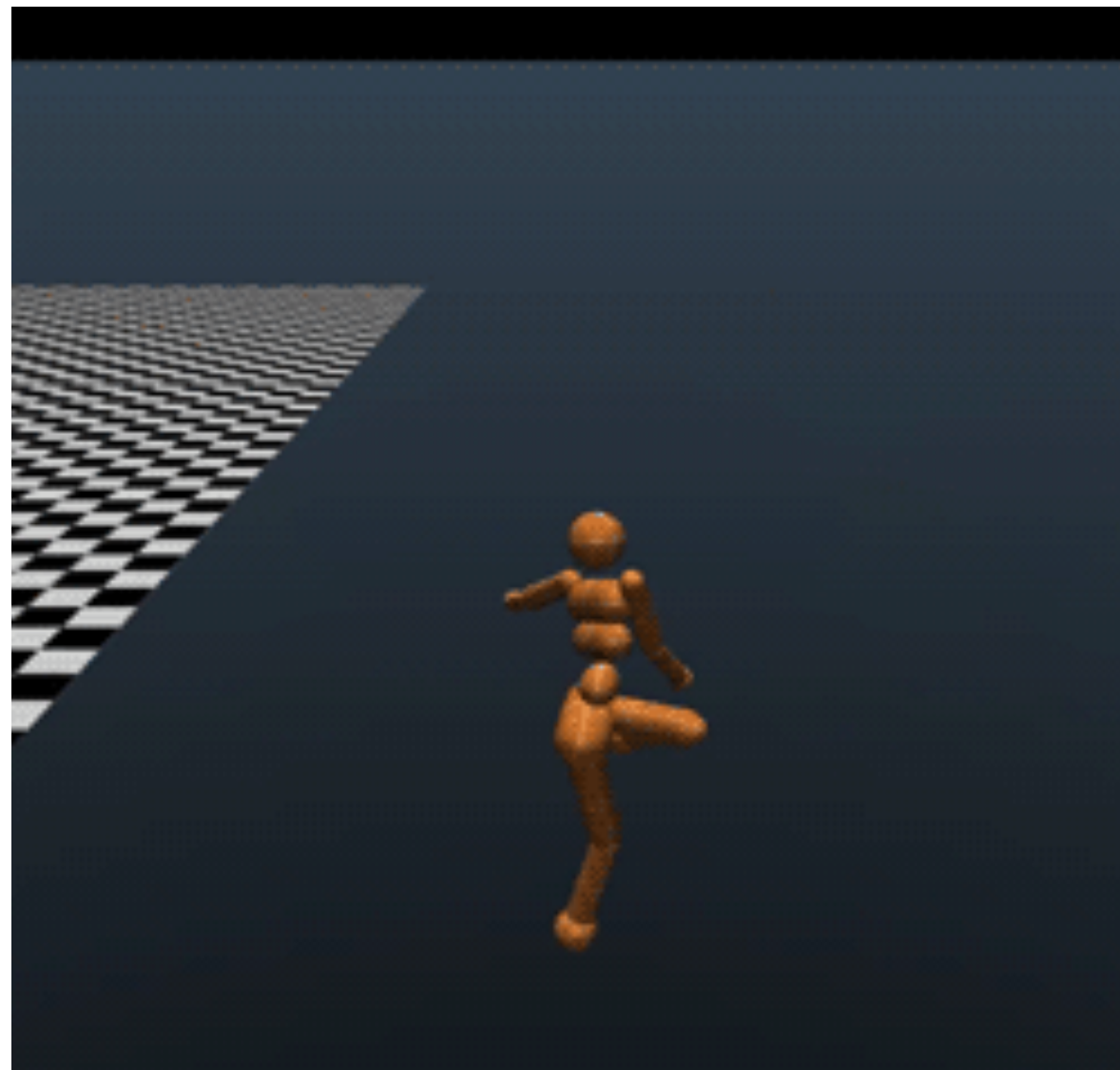
# Interesting videos from the today's algorithm

**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:

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

## 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:

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

## 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  $P = Q$

# 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:

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

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



## A trust region formulation for policy update:

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

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

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

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:

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

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

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:

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

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

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

# 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

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

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

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

# 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:

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

**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  $\nabla_{\theta}^2 \ell(\theta_t)$  has a nice form!

**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:

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

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

Change from trajectory distribution to state-action distribution:

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

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

Change from trajectory distribution to state-action distribution:

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

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



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

Change from trajectory distribution to state-action distribution:

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

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

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

Change from trajectory distribution to state-action distribution:

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

$$\begin{aligned} \nabla_{\theta} \ell(\theta) \big|_{\theta=\theta_t} &= \mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \left( -\nabla_{\theta} \ln \pi_{\theta}(a_h \mid s_h) \big|_{\theta=\theta_t} \right) \\ &= -\mathbb{E}_{s \sim d_{\mu}^{\pi_{\theta_t}}} \sum_a \pi_{\theta_t}(a \mid s) \frac{\nabla_{\theta} \pi_{\theta_t}(a \mid s)}{\pi_{\theta_t}(a \mid s)} = \mathbf{0} \end{aligned}$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

It's called fisher Information Matrix!

## Summary so far:

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

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

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



## Summary so far:

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

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

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

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:

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

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

This leads to the following much simplified constrained optimization:

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

# 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:

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:

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

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:

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

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  $\theta_0$

For  $t = 0, \dots$



# Algorithm: Natural Policy Gradient

Initialize  $\theta_0$

For  $t = 0, \dots$

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

# Algorithm: Natural Policy Gradient

Initialize  $\theta_0$

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  $\theta_0$

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  $\theta_0$

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

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

# Algorithm: Natural Policy Gradient

Initialize  $\theta_0$

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 HW2 on Cartpole)

# Summary for today:

## 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 today:

## 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 today:

## Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

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



# Summary for today:

## Trust Region Policy Optimization and NPG

At iteration t:

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

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

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 today:

## Trust Region Policy Optimization and NPG

At iteration t:

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

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

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

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

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  $\arg \max_{\theta}$ )