

Advance Policy Gradient

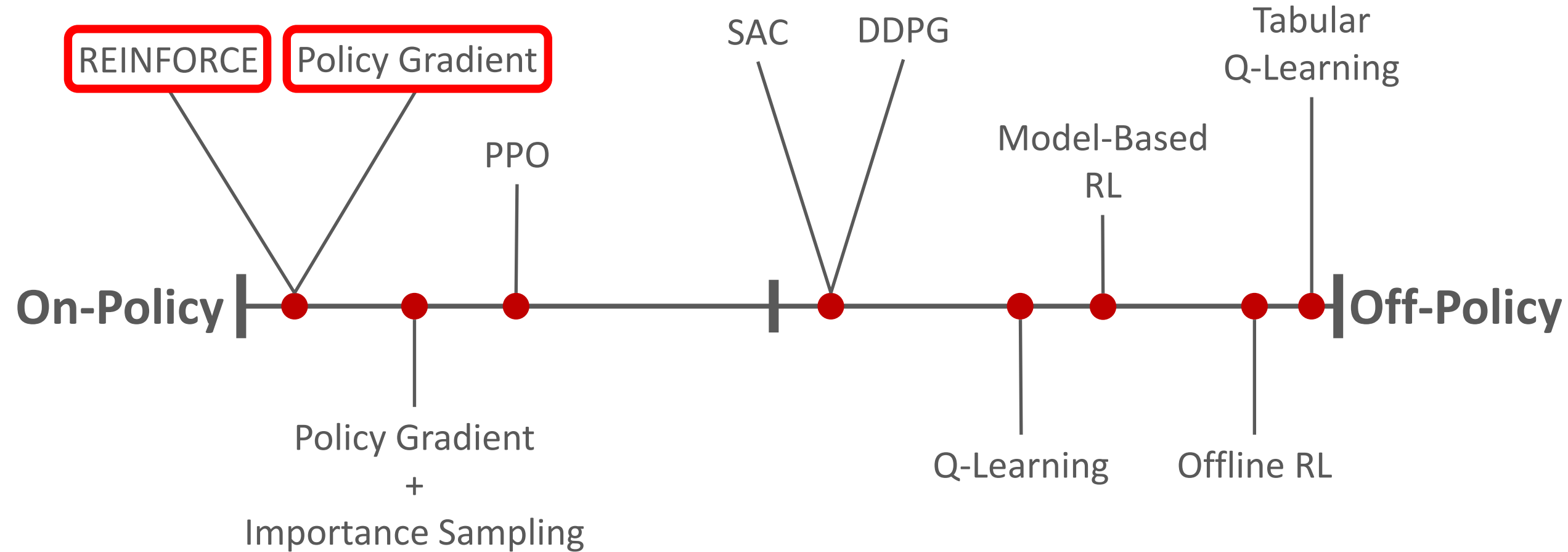
CMPT 729 G100

Jason Peng

Overview

- Off-Policy Policy Gradient
- Constrained Policy Optimization
- Proximal Policy Optimization

On-Policy vs Off-Policy



REINFORCE

ALGORITHM: REINFORCE

1: $\theta \leftarrow$ initialize policy parameters

2: **while** not done **do**

3: Sample trajectories $\{\tau^i\}$ from policy $\pi_\theta(\mathbf{a}|\mathbf{s})$

4: Estimate policy gradient

$$\nabla_\theta J(\pi_\theta) \approx \frac{1}{N} \sum_i R(\tau^i) \sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)$$

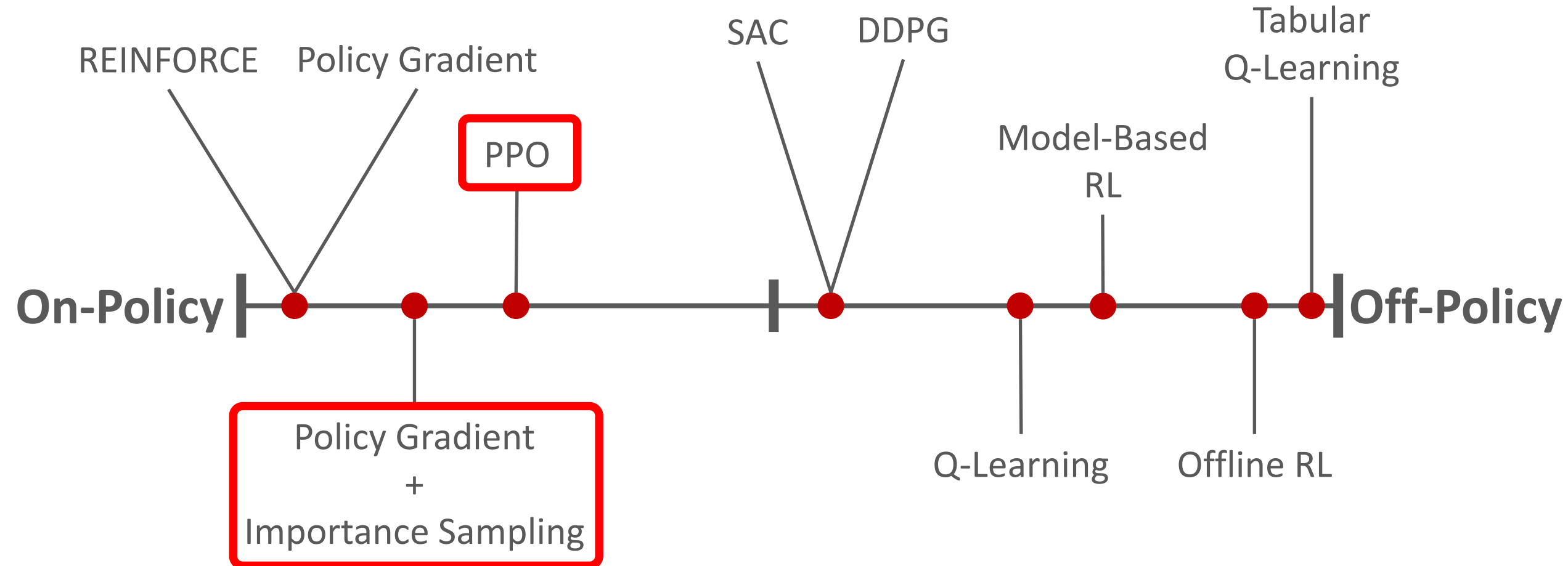
5: Update policy $\theta \leftarrow \theta + \alpha \nabla_\theta J(\pi_\theta)$

6: **end while**

7: return policy π_θ

Perform just one grad update,
then throw out data

On-Policy vs Off-Policy



Off-Policy REINFORCE

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\pi)} [\nabla_{\pi} \log p(\tau|\pi) R(\tau)]$$

Must be from
current policy

- Off-Policy Reinforce: can we estimate $\nabla_{\pi} J(\pi)$ using data from another policy $\mu(\mathbf{a}|\mathbf{s})$?

Importance Sampling

- Want to estimate $\mathbb{E}_{x \sim p(x)} [f(x)]$, but only have data $x \sim q(x)$

$$\begin{aligned}\mathbb{E}_{x \sim p(x)} [f(x)] &= \sum_x p(x) f(x) \\ &= \sum_x \frac{q(x)}{\underline{q(x)}} p(x) f(x) \\ &= 1\end{aligned}$$

Importance Sampling

- Want to estimate $\mathbb{E}_{x \sim p(x)} [f(x)]$, but only have data $x \sim q(x)$

$$\begin{aligned}\mathbb{E}_{x \sim p(x)} [f(x)] &= \sum_x p(x) f(x) \\ &= \sum_x \frac{q(x)}{q(x)} p(x) f(x) \\ &= \sum_x q(x) \frac{p(x)}{q(x)} f(x) = \mathbb{E}_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right]\end{aligned}$$

“Importance Sampling”
weight

Off-Policy REINFORCE

$$\begin{aligned}\nabla_{\pi} J(\pi) &= \mathbb{E}_{\tau \sim p(\tau|\pi)} [\nabla_{\pi} \log p(\tau|\pi) R(\tau)] \\ &= \sum_{\tau} p(\tau|\pi) \nabla_{\pi} \log p(\tau|\pi) R(\tau)\end{aligned}$$

$\mu(\mathbf{a}|\mathbf{s})$: behavior policy

$$= \sum_{\tau} \frac{p(\tau|\mu)}{\underbrace{p(\tau|\mu)}_{= 1}} p(\tau|\pi) \nabla_{\pi} \log p(\tau|\pi) R(\tau)$$

Off-Policy REINFORCE

$$\begin{aligned}\nabla_{\pi} J(\pi) &= \mathbb{E}_{\tau \sim p(\tau|\pi)} [\nabla_{\pi} \log p(\tau|\pi) R(\tau)] \\ &= \sum_{\tau} p(\tau|\pi) \nabla_{\pi} \log p(\tau|\pi) R(\tau)\end{aligned}$$

$\mu(\mathbf{a}|\mathbf{s})$: behavior policy

$$\begin{aligned}&= \sum_{\tau} \frac{p(\tau|\mu)}{p(\tau|\mu)} p(\tau|\pi) \nabla_{\pi} \log p(\tau|\pi) R(\tau) \\ &= \sum_{\tau} p(\tau|\mu) \frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]\end{aligned}$$

Off-Policy REINFORCE

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]$$

Data sampled
according to μ

Importance Sampling

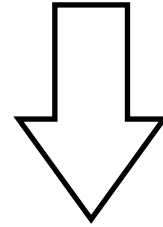
$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]$$

“Importance Sampling”
weight

Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]$$

If $p(\tau|\mu) = p(\tau|\pi)$:



$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\pi)} [\nabla_{\pi} \log p(\tau|\pi) R(\tau)]$$

Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]$$

< 1

If $p(\tau|\pi) < p(\tau|\mu)$:

- Down-weight likelihood of trajectory

Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{\underline{p(\tau|\mu)}} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]$$

> 1

If $p(\tau|\pi) > p(\tau|\mu)$:

- Up-weight likelihood of trajectory

Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right]$$

$$\begin{aligned} \frac{p(\tau|\pi)}{p(\tau|\mu)} &= \frac{\cancel{p(\mathbf{s}_0)} \prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t) \cancel{p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)}}{\cancel{p(\mathbf{s}_0)} \prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t) \cancel{p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)}} \\ &= \frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \end{aligned}$$

Importance Sampling

$$\begin{aligned}\nabla_{\pi} J(\pi) &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right] \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \nabla_{\pi} \log p(\tau|\pi) \right] \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \left(\frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \right) \left(\sum_{t=0}^{T-1} \nabla_{\pi} \log \pi(\mathbf{a}_t|\mathbf{s}_t) \right) \right]\end{aligned}$$

Importance Sampling

$$\begin{aligned}\nabla_{\pi} J(\pi) &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right] \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \nabla_{\pi} \log p(\tau|\pi) \right] \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \left(\frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \right) \left(\sum_{t=0}^{T-1} \nabla_{\pi} \log \pi(\mathbf{a}_t|\mathbf{s}_t) \right) \right]\end{aligned}$$

Importance Sampling

$$\begin{aligned}\nabla_{\pi} J(\pi) &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[\frac{p(\tau|\pi)}{p(\tau|\mu)} \nabla_{\pi} \log p(\tau|\pi) R(\tau) \right] \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \nabla_{\pi} \log p(\tau|\pi) \right] \\ &= \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \left(\frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \right) \left(\sum_{t=0}^{T-1} \nabla_{\pi} \log \pi(\mathbf{a}_t|\mathbf{s}_t) \right) \right]\end{aligned}$$

Importance Sampling

$$\mathcal{T} = \left\{ \begin{array}{l} \boxed{\text{s}_0 \quad \text{a}_0 \quad r_0} + \boxed{\log \mu(\mathbf{a}_0 | \mathbf{s}_0)} \\ \boxed{\text{s}_1 \quad \text{a}_1 \quad r_1} + \boxed{\log \mu(\mathbf{a}_1 | \mathbf{s}_1)} \\ \vdots \\ \boxed{\text{s}_T} \end{array} \right.$$

$\left(\frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t | \mathbf{s}_t)} \right)$

Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \left(\frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \right) \left(\sum_{t=0}^{T-1} \nabla_{\pi} \log \pi(\mathbf{a}_t|\mathbf{s}_t) \right) \right]$$

- Can estimate gradient from arbitrary distribution, as long as $\mu(\mathbf{a}|\mathbf{s}) > 0$ for all actions (e.g. Gaussian distribution)
- Never used in practice

Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\tau \sim p(\tau|\mu)} \left[R(\tau) \left(\frac{\prod_{t=0}^{T-1} \pi(\mathbf{a}_t|\mathbf{s}_t)}{\prod_{t=0}^{T-1} \mu(\mathbf{a}_t|\mathbf{s}_t)} \right) \left(\sum_{t=0}^{T-1} \nabla_{\pi} \log \pi(\mathbf{a}_t|\mathbf{s}_t) \right) \right]$$

- Can estimate gradient from arbitrary distribution, as long as $\mu(\mathbf{a}|\mathbf{s}) > 0$ for all actions (e.g. Gaussian distribution)
- Never used in practice
 - Very high variance if $\pi \neq \mu$
 - Importance sampling weights very quickly vanish or explode

Reward-to-Go Policy Gradient

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underbrace{(Q^{\pi}(\mathbf{s}, \mathbf{a}) - V^{\pi}(\mathbf{s}))}]$$

“advantage”

$$A^{\pi}(\mathbf{s}, \mathbf{a})$$

Reward-to-Go Policy Gradient

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a})]$$

$\mu(\mathbf{a}|\mathbf{s})$: behavior policy

$$\begin{aligned} \nabla_{\pi} J(\pi) &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[\frac{\mu(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right] \\ &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \underline{\mu(\mathbf{a}|\mathbf{s})}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right] \end{aligned}$$

Reward-to-Go Policy Gradient

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a})]$$

$\mu(\mathbf{a}|\mathbf{s})$: behavior policy

$$\begin{aligned} \nabla_{\pi} J(\pi) &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[\frac{\mu(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right] \\ &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right] \end{aligned}$$

single-step
lower variance

Reward-to-Go Policy Gradient

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a})]$$

$\mu(\mathbf{a}|\mathbf{s})$: behavior policy

$$\begin{aligned} \nabla_{\pi} J(\pi) &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[\frac{\mu(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right] \\ &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right] \end{aligned}$$

What about the
state distribution?

Reward-to-Go Policy Gradient

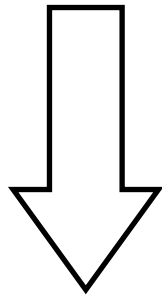
$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim \underline{d_{\pi}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right]$$

Computing the IS weights
for $d_{\pi}(\mathbf{s})$ is intractable.

$$\frac{\cancel{d_{\pi}(\mathbf{s})}}{\cancel{d_{\mu}(\mathbf{s})}}$$

Reward-to-Go Policy Gradient

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim \underline{d_{\pi}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right]$$

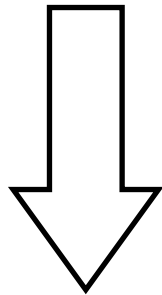


$$\nabla_{\pi} J(\pi) \approx \mathbb{E}_{\mathbf{s} \sim \underline{d_{\mu}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\pi}(\mathbf{s}, \mathbf{a}) \right]$$

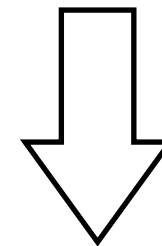
Ok, if $\mu \approx \pi$?

Reward-to-Go Policy Gradient

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim \underline{d_{\pi}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underline{A^{\pi}(\mathbf{s}, \mathbf{a})} \right]$$



$$\nabla_{\pi} J(\pi) \approx \mathbb{E}_{\mathbf{s} \sim \underline{d_{\mu}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underline{A^{\pi}(\mathbf{s}, \mathbf{a})} \right]$$



$$\approx \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underline{A^{\mu}(\mathbf{s}, \mathbf{a})} \right]$$

Policy Gradient + Importance Sampling

$$\nabla_{\pi} J^{\mu}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

Surrogate objective:

$$J^{\mu}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$

Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \underline{A^\mu(\mathbf{s}, \mathbf{a})} \right]$$

Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[\underline{Q^\pi(\mathbf{s}, \mathbf{a})} \right]$$

Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$

Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$

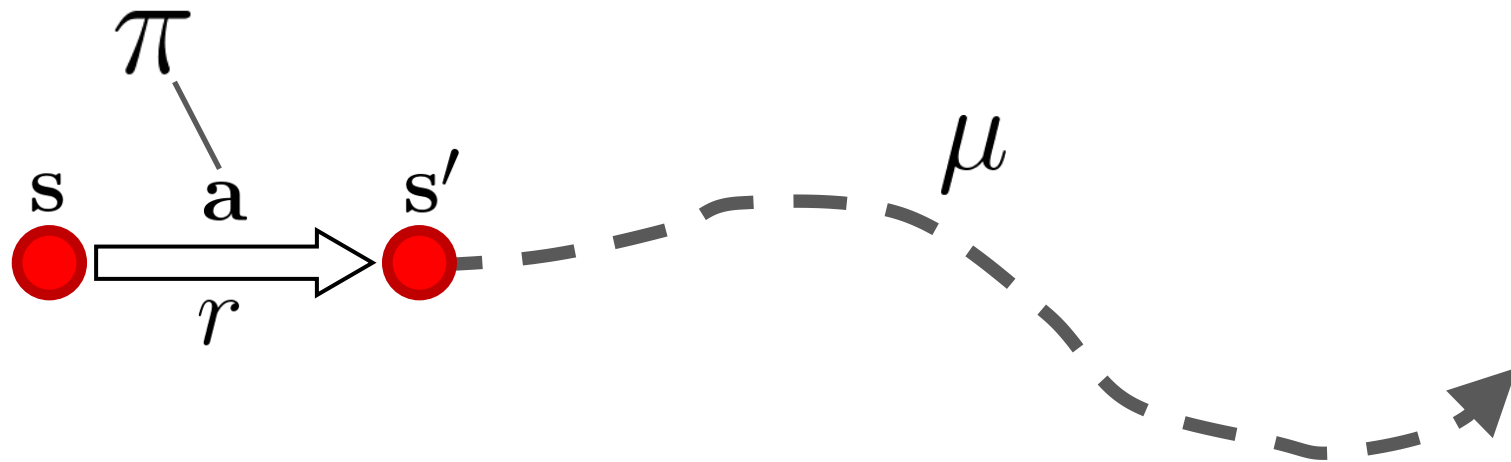
Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$



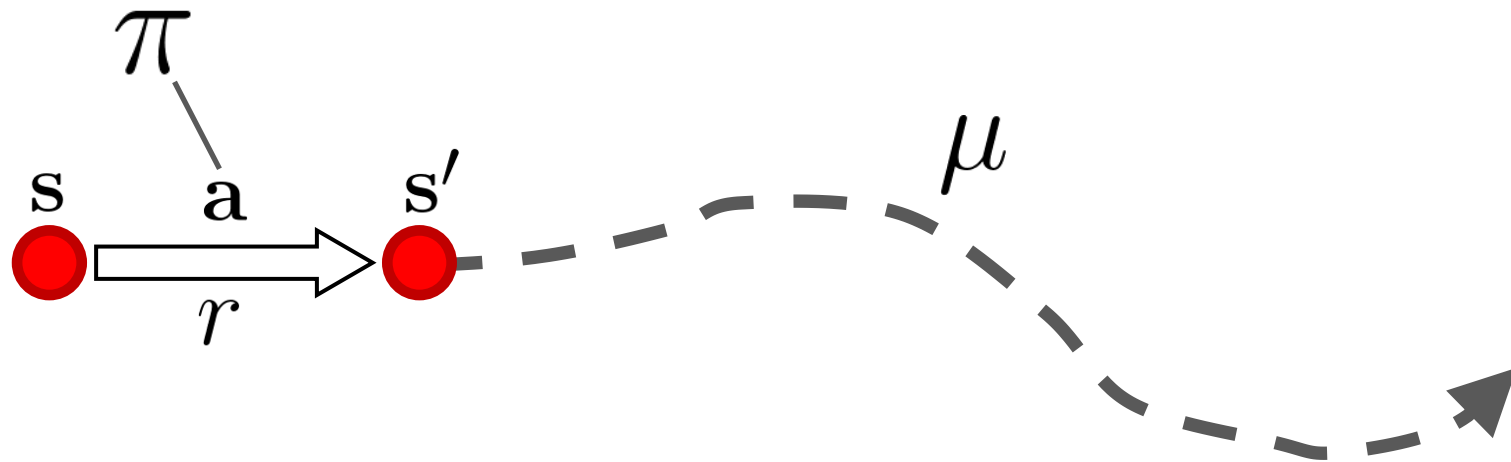
Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$



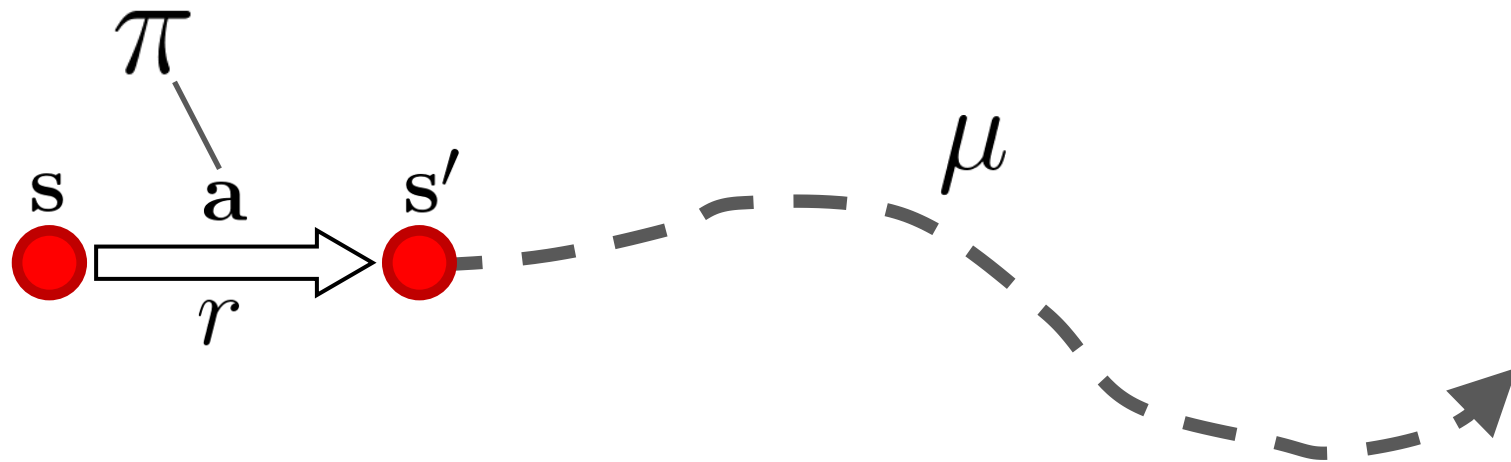
Surrogate Objective

Policy Gradient + Importance Sampling:

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

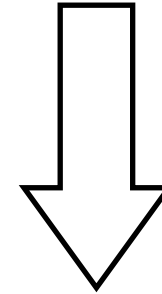
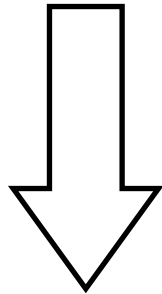
Soft Actor-Critic:

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$



Policy Gradient + Importance Sampling

$$\nabla_{\pi} J(\pi) = \mathbb{E}_{\mathbf{s} \sim \underline{d_{\pi}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underline{A^{\pi}(\mathbf{s}, \mathbf{a})} \right]$$



$$\nabla_{\pi} J^{\mu}(\pi) = \mathbb{E}_{\mathbf{s} \sim \underline{d_{\mu}(\mathbf{s})}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} \nabla_{\pi} \log \pi(\mathbf{a}|\mathbf{s}) \underline{A^{\mu}(\mathbf{s}, \mathbf{a})} \right]$$

Ok, if $\mu \approx \pi$?

Surrogate Objective

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

Reasonable if π is *close* to μ

$$D_{\text{KL}}^{\max}(\mu, \pi) = \max_{\mathbf{s}} D_{\text{KL}}(\mu(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))$$

Surrogate Objective

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

If $D_{\text{KL}}^{\text{max}}(\mu, \pi) \leq \epsilon$,

$$J(\pi) \geq J^\mu(\pi) - \underline{C}\epsilon$$

constant

Surrogate Objective

$$J^\mu(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\mu(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^\mu(\mathbf{s}, \mathbf{a}) \right]$$

If $D_{\text{KL}}^{\max}(\mu, \pi) \leq \epsilon$,

$$J(\pi) \geq J^\mu(\pi) - C\epsilon$$

The surrogate objective is a lower bound on the real objective for sufficiently small ϵ !

Constrained Optimization

$$\arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

$$\text{s.t. } \underline{D_{\text{KL}}^{\text{max}}(\mu, \pi)} \leq \epsilon \quad \text{“Trust region”}$$

$$D_{\text{KL}}^{\text{max}}(\mu, \pi) = \max_{\mathbf{s}} D_{\text{KL}}(\mu(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))$$

Constrained Optimization

$$\arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

$$\text{s.t. } \underline{D_{\text{KL}}^{\text{max}}(\mu, \pi)} \leq \epsilon$$

$$D_{\text{KL}}^{\text{max}}(\mu, \pi) = \max_{\mathbf{s}} D_{\text{KL}}(\mu(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))$$

Hard to compute

Constrained Optimization

$$\arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

$$\text{s.t. } \underline{D_{\text{KL}}^{\text{mean}}(\mu, \pi)} \leq \epsilon$$

$$D_{\text{KL}}^{\text{mean}}(\mu, \pi) = \mathbb{E}_{\mathbf{s} \sim d^{\mu}(\mathbf{s})} [D_{\text{KL}}(\mu(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))]$$

Constrained Optimization

$$\arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $D_{\text{KL}}^{\text{mean}}(\mu, \pi) \leq \epsilon$

How do we pick μ ?

- In practice, collect data using current policy $\mu = \pi^k$

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Trust Region Policy Optimization

[Schulman et al. 2015]

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Trust Region Policy Optimization

[Schulman et al. 2015]

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Trust Region Policy Optimization

[Schulman et al. 2015]

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

- 1: $\pi_0 \leftarrow$ initialize policy
 - 2: **for** iteration $k = 0, \dots, n - 1$ **do**
 - 3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$
 - 4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$
 - 5: Fit value function $V^k(\mathbf{s})$
 - 6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}
 - 7: Update policy:
$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$
$$\text{s.t. } \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$$
 - 8: **end for**
 - 9: return policy π^n
-

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

1: $\pi_0 \leftarrow$ initialize policy

2: **for** iteration $k = 0, \dots, n - 1$ **do**

3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$

4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$

5: Fit value function $V^k(\mathbf{s})$

6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}

7: Update policy:

$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$

8: **end for**

9: return policy π^n

Still need to collect a new batch of data every iteration

Constrained Policy Optimization

ALGORITHM: Constrained Policy Optimization

- 1: $\pi_0 \leftarrow$ initialize policy
 - 2: **for** iteration $k = 0, \dots, n - 1$ **do**
 - 3: Sample trajectories τ^i from policy $\pi^k(\mathbf{a}|\mathbf{s})$
 - 4: Store trajectories in dataset $\mathcal{D} = \{\tau^i\}$
 - 5: Fit value function $V^k(\mathbf{s})$
 - 6: Calculate advantage $A^k(\mathbf{s}, \mathbf{a})$ for every (\mathbf{s}, \mathbf{a}) in \mathcal{D}
 - 7: Update policy:
$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\pi^k(\mathbf{a}|\mathbf{s})} A^k(\mathbf{s}, \mathbf{a}) \right]$$
$$\text{s.t. } \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [D_{\text{KL}}(\pi^k(\cdot|\mathbf{s}) || \pi(\cdot|\mathbf{s}))] \leq \epsilon$$
 - 8: **end for**
 - 9: return policy π^n
-

Update policy with multiple grad steps



Constrained Optimization

$$\arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$
$$\text{s.t. } D_{\text{KL}}^{\text{mean}}(\mu, \pi) \leq \epsilon$$

How do we solve this?


Trust Region Policy Optimization (TRPO):

- Linear approximation of objective
- Quadratic approximation of constraint
- Solve with conjugate gradient method

Lagrangian

$$\arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

s.t. $D_{\text{KL}}^{\text{mean}}(\mu, \pi) \leq \epsilon$



Lagrangian

$$\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda \left(\underline{D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon} \right)$$

Lagrangian

$$\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)$$

“Lagrange multiplier”

Lagrangian

$$\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda \left(D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon \right)$$

$\lambda \rightarrow \infty$

> 0
constraint violated

Lagrangian

$$\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda \left(D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon \right)$$

$\lambda \rightarrow 0$

< 0
constraint satisfied

Lagrangian

$$\underbrace{\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)}_{\mathcal{L}(\pi, \lambda)}$$

Dual gradient descent:

- Maximize $\mathcal{L}(\pi, \lambda)$ wrt π
- Update λ : $\lambda \leftarrow \max(0, \lambda + \alpha \underbrace{(D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)}_{= -\nabla_{\lambda} \mathcal{L}(\pi, \lambda)})$

Lagrangian

$$\underbrace{\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)}_{\mathcal{L}(\pi, \lambda)}$$

Dual gradient descent:

- Maximize $\mathcal{L}(\pi, \lambda)$ wrt π
- Update λ : $\lambda \leftarrow \max(0, \lambda + \underline{\alpha} (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon))$
stepsize

Lagrangian

$$\underbrace{\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)}_{\mathcal{L}(\pi, \lambda)}$$

Dual gradient descent:

- Maximize $\mathcal{L}(\pi, \lambda)$ wrt π
- Update λ : $\lambda \leftarrow \max(0, \lambda + \alpha (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon))$

gradient descent

Lagrangian

$$\underbrace{\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)}_{\mathcal{L}(\pi, \lambda)}$$

Dual gradient descent:

- Maximize $\mathcal{L}(\pi, \lambda)$ wrt π
- Update λ : $\lambda \leftarrow \max(0, \lambda + \alpha (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon))$

Lagrangian

$$\underbrace{\arg \max_{\pi} \min_{\lambda \geq 0} \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\mu(\mathbf{a}|\mathbf{s})} A^{\mu}(\mathbf{s}, \mathbf{a}) \right] - \lambda (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon)}_{\mathcal{L}(\pi, \lambda)}$$

Dual gradient descent:

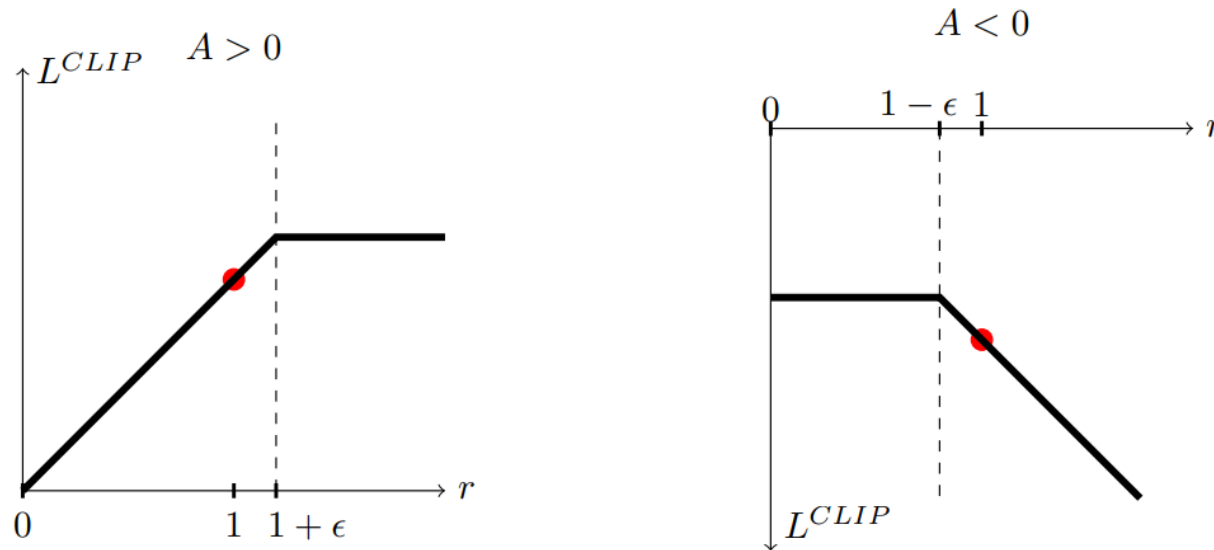
- Maximize $\mathcal{L}(\pi, \lambda)$ wrt π
 - Update λ : $\lambda \leftarrow \max(0, \lambda + \alpha (D_{\text{KL}}^{\text{mean}}(\mu, \pi) - \epsilon))$
- } Proximal Policy Optimization (PPO)

Proximal Policy Optimization (PPO)

In practice:

- Most PPO implementations use a clipping objective:

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \underline{\text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)} \hat{A}_t) \right]$$



Proximal Policy Optimization (PPO)

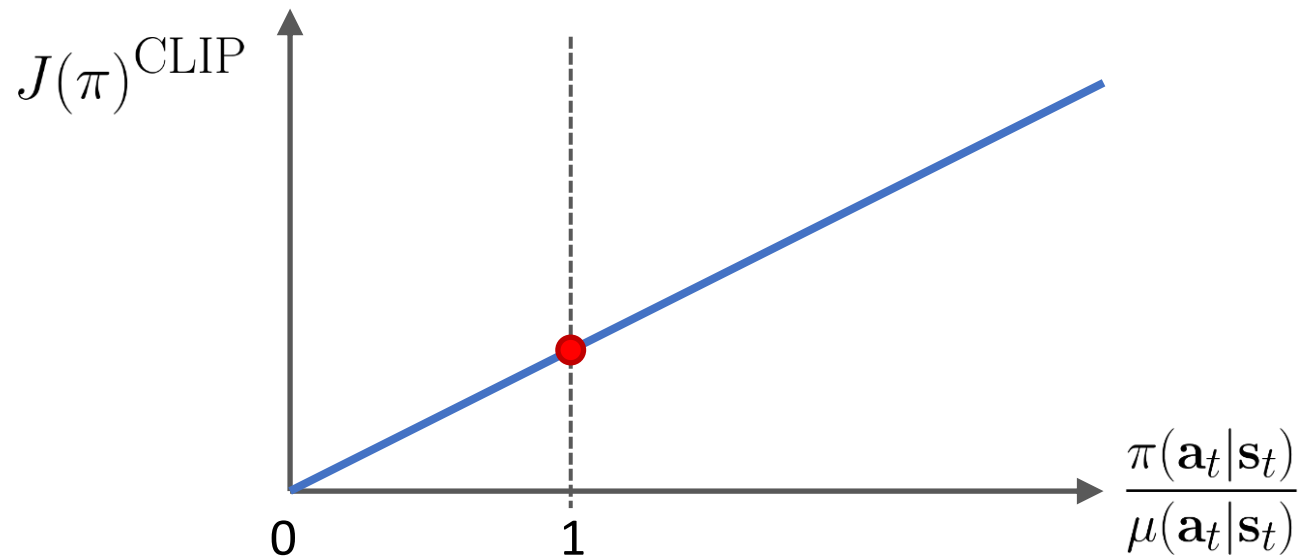
$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)} A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

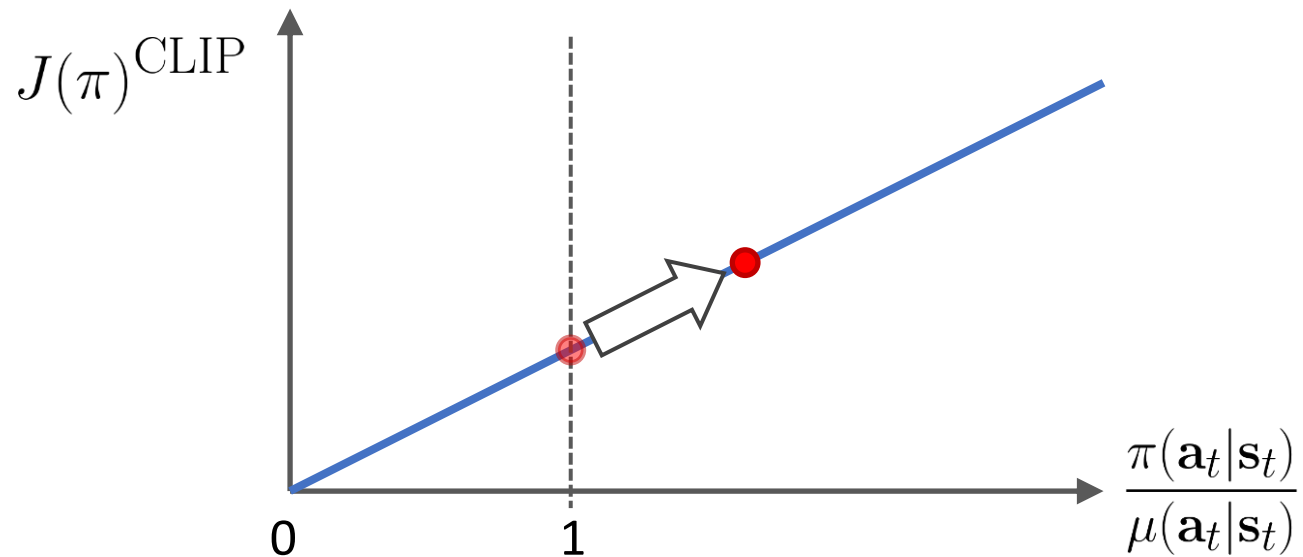
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{> 0} \right]$$



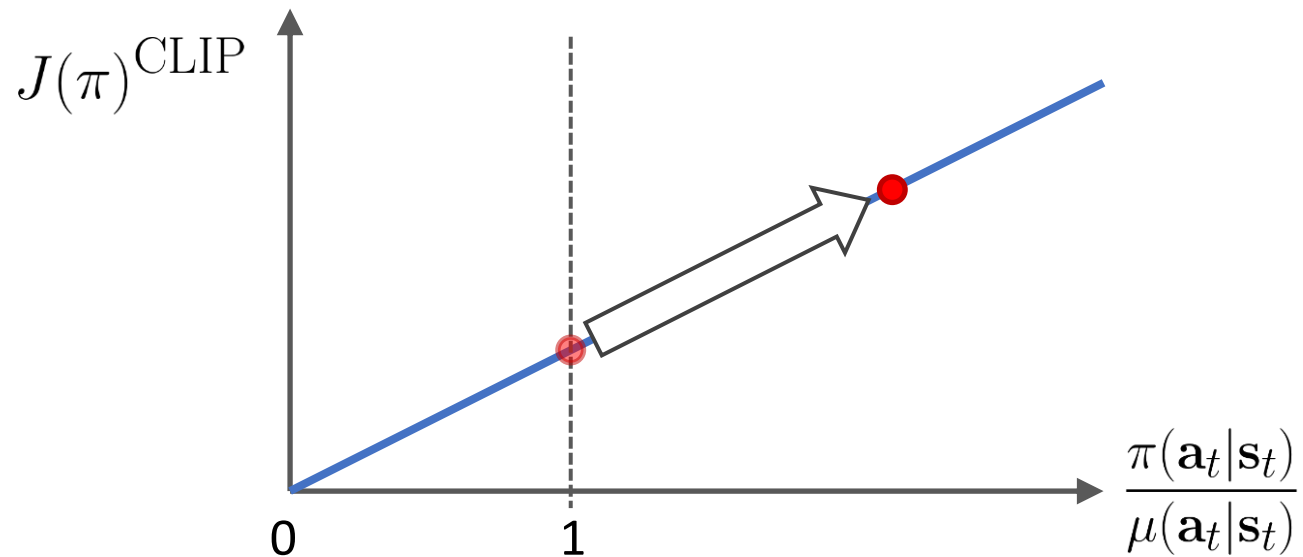
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{> 0} \right]$$



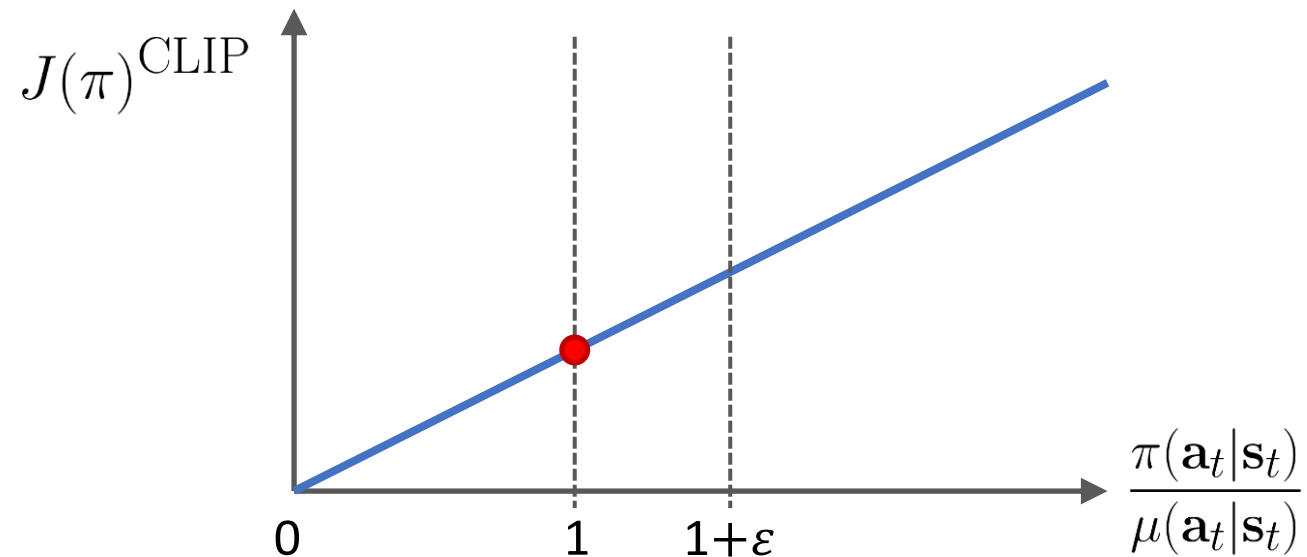
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{> 0} \right]$$



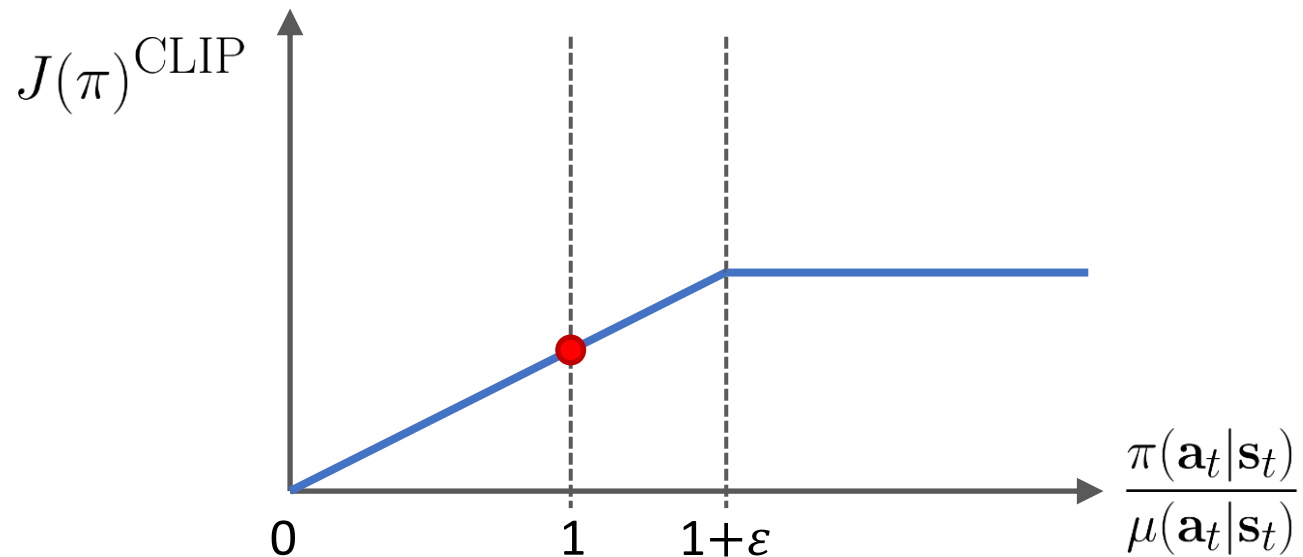
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{> 0} \right]$$



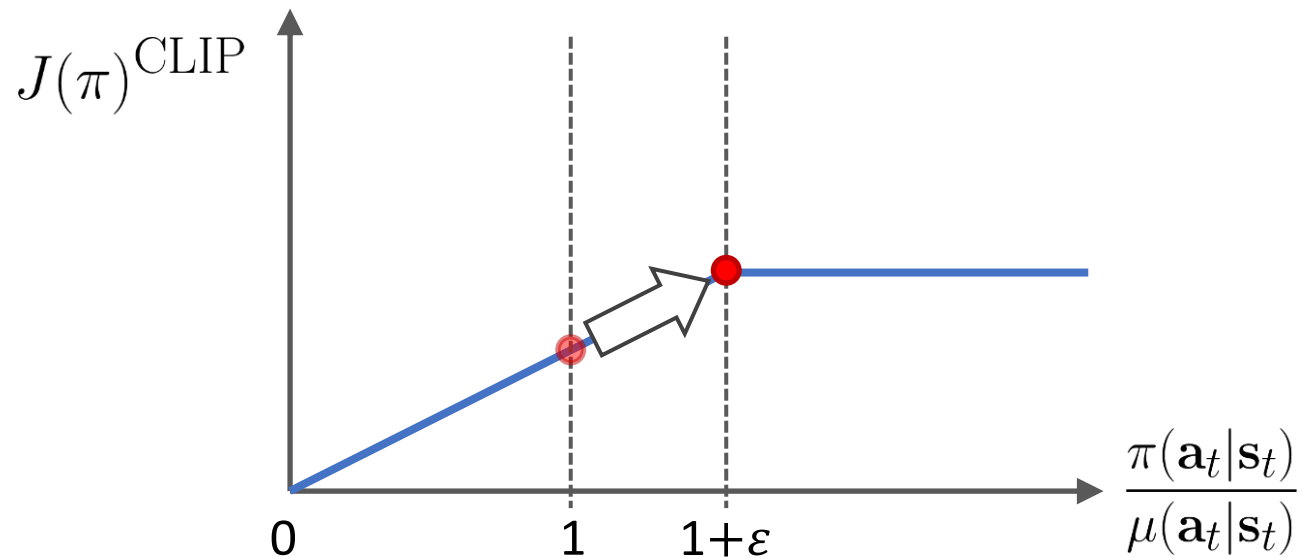
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{> 0} \right]$$



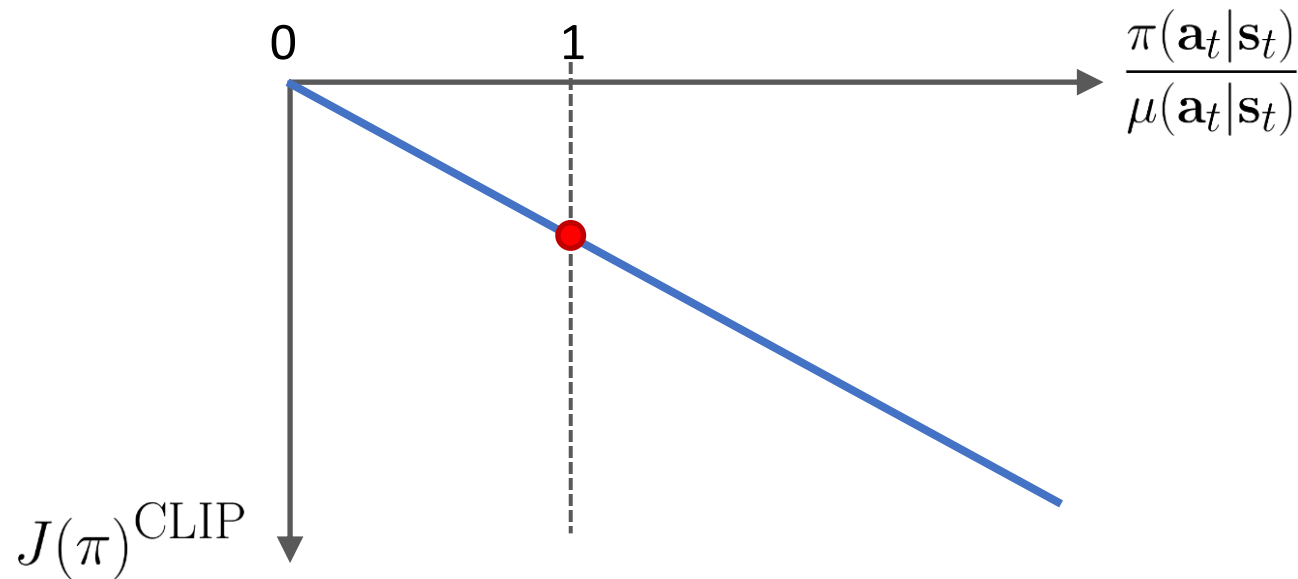
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{> 0} \right]$$



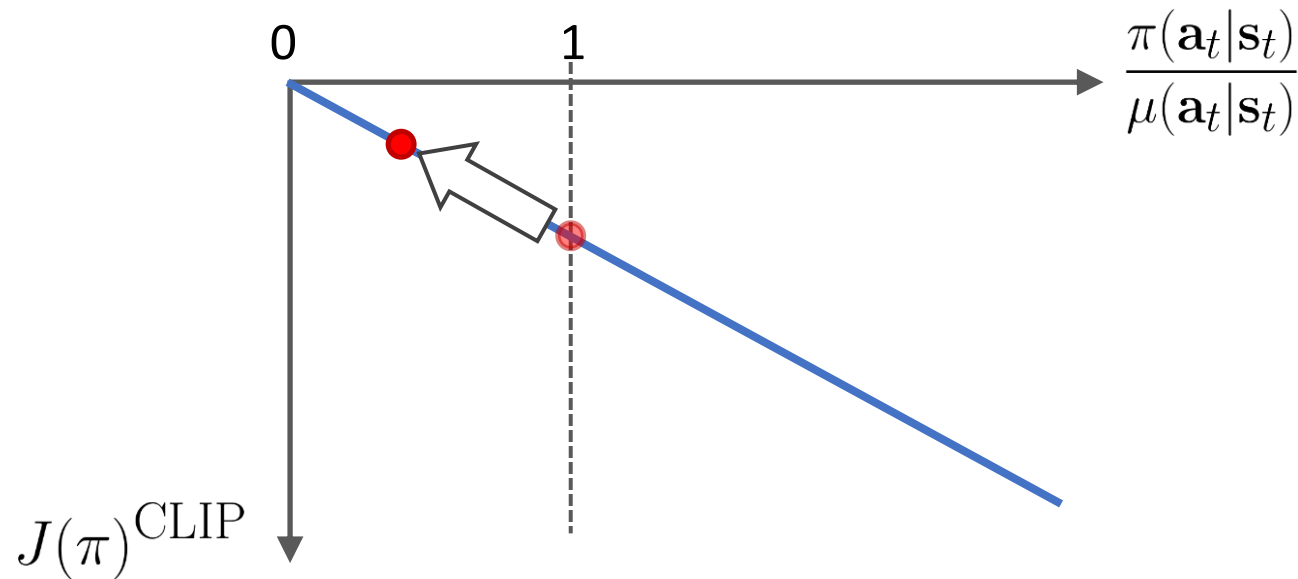
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{< 0} \right]$$



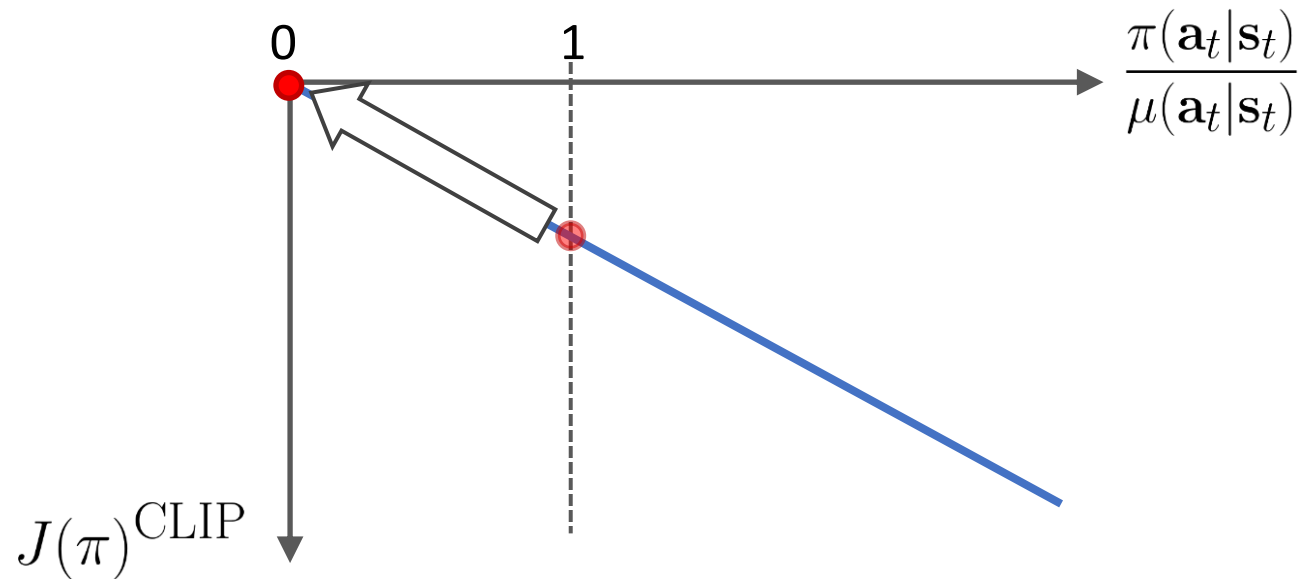
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{< 0} \right]$$



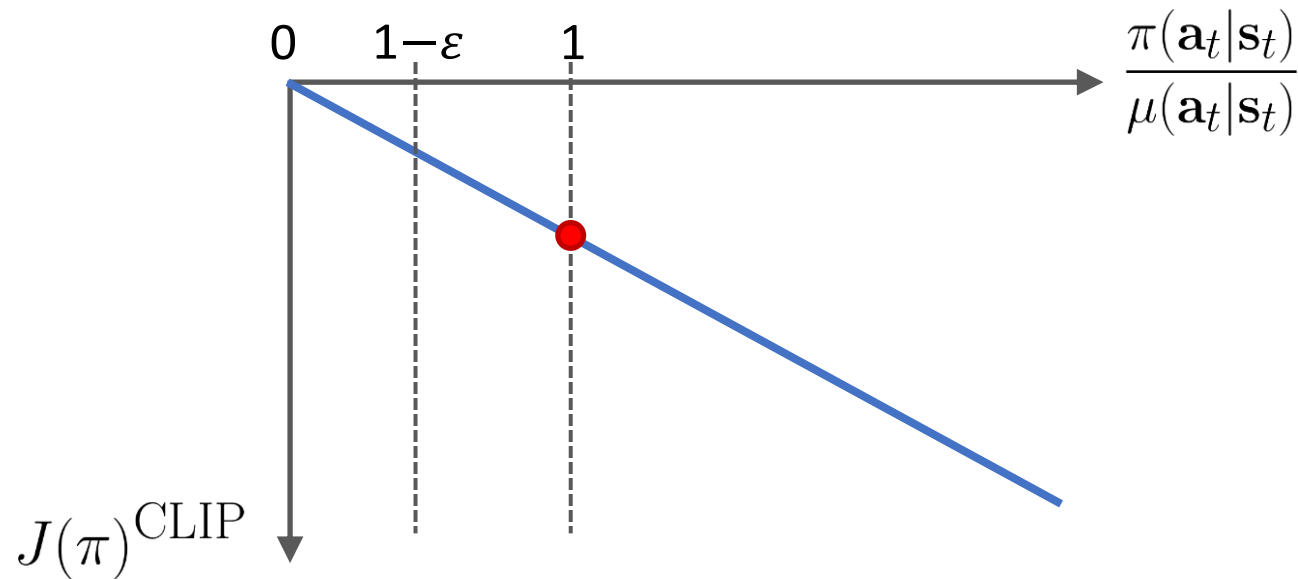
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{< 0} \right]$$



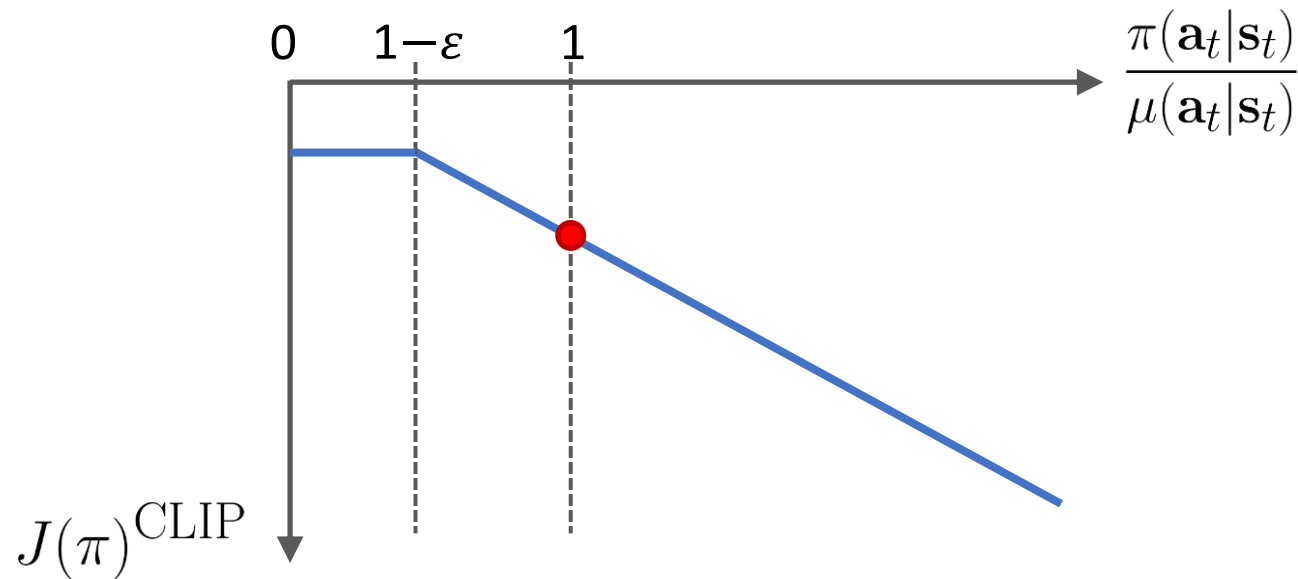
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{< 0} \right]$$



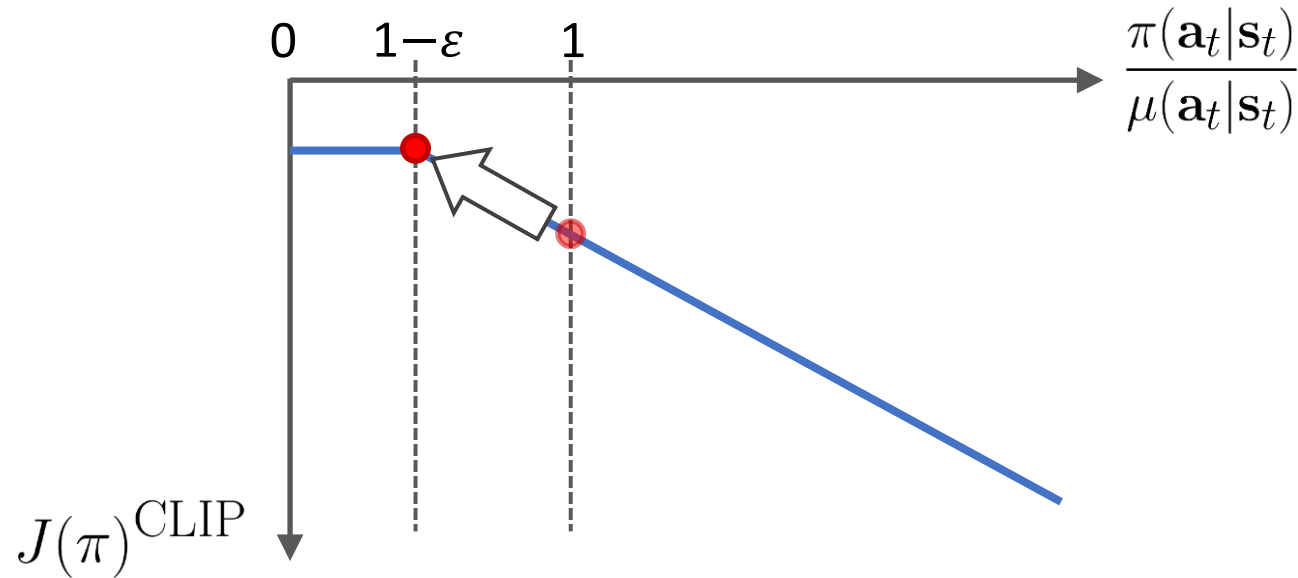
Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{< 0} \right]$$



Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \frac{A^{\mu}(\mathbf{s}, \mathbf{a})}{< 0} \right]$$

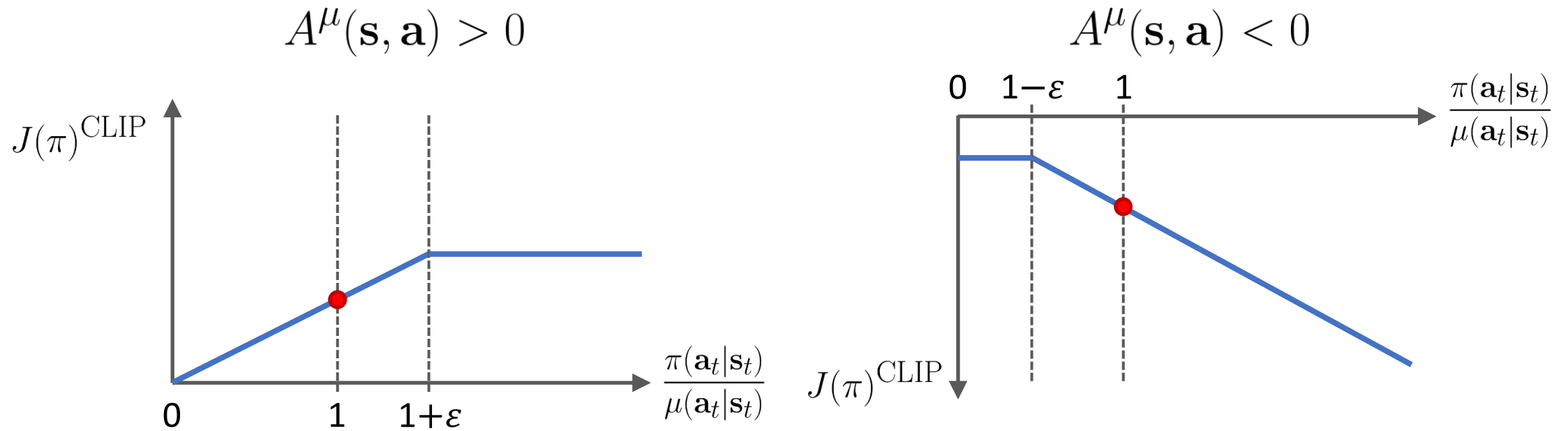


Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$

Proximal Policy Optimization (PPO)

$$J(\pi)^{\text{CLIP}} = \mathbb{E}_{\mathbf{s} \sim d_{\mu}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} \left[\text{clip} \left(\frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\mu(\mathbf{a}_t|\mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) A^{\mu}(\mathbf{s}, \mathbf{a}) \right]$$



Robotic Locomotion



Learning Robust Perceptive Locomotion for Quadrupedal Robots in the Wild
[Miki et al. 2022]

Dota



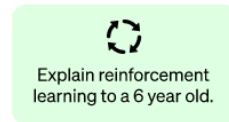
Dota 2 with Large Scale Deep Reinforcement Learning
[OpenAI et al. 2019]

ChatGPT

Step 1

Collect demonstration data and train a supervised policy.

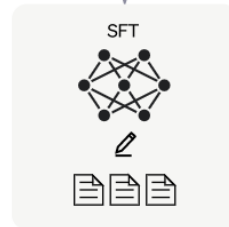
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



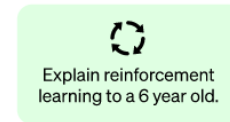
This data is used to fine-tune GPT-3.5 with supervised learning.



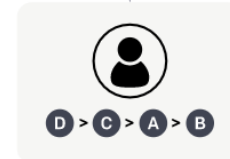
Step 2

Collect comparison data and train a reward model.

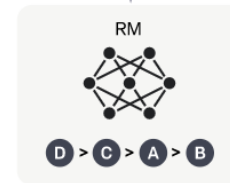
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

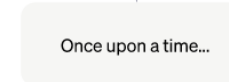
A new prompt is sampled from the dataset.



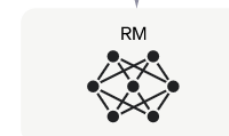
The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



[OpenAI 2022]

Summary

- Off-Policy Policy Gradient
- Constrained Policy Optimization
- Proximal Policy Optimization