

Carnegie Mellon

School of Computer Science

Deep Reinforcement Learning and Control

# Natural Policy Gradients, TRPO, PPO

CMU 10703

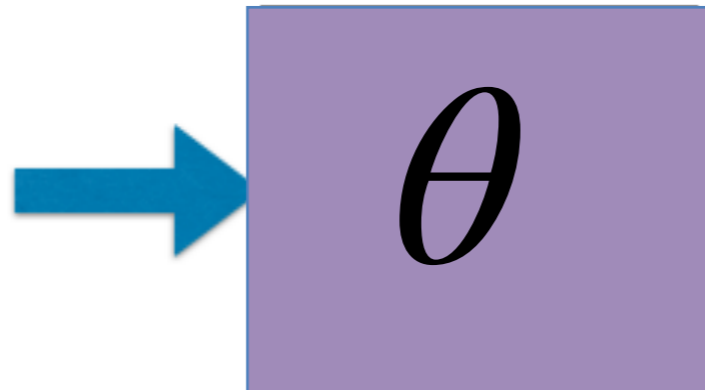
Katerina Fragkiadaki



Part of the slides adapted from John Shulman and Joshua Achiam

# Stochastic policies

## continuous actions

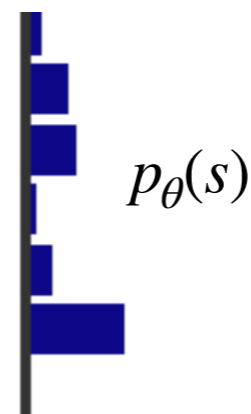
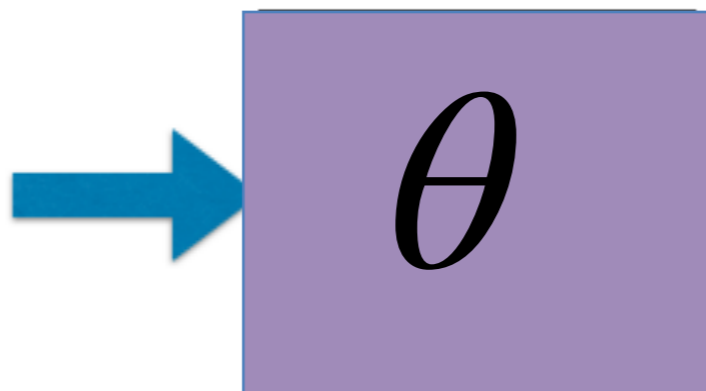


$$\mu_{\theta}(s)$$
$$\sigma_{\theta}(s)$$

usually multivariate  
Gaussian

$$a \sim \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}^2(s))$$

## discrete actions



almost always  
categorical

$$a \sim \text{Cat}(p_{\theta}(s))$$

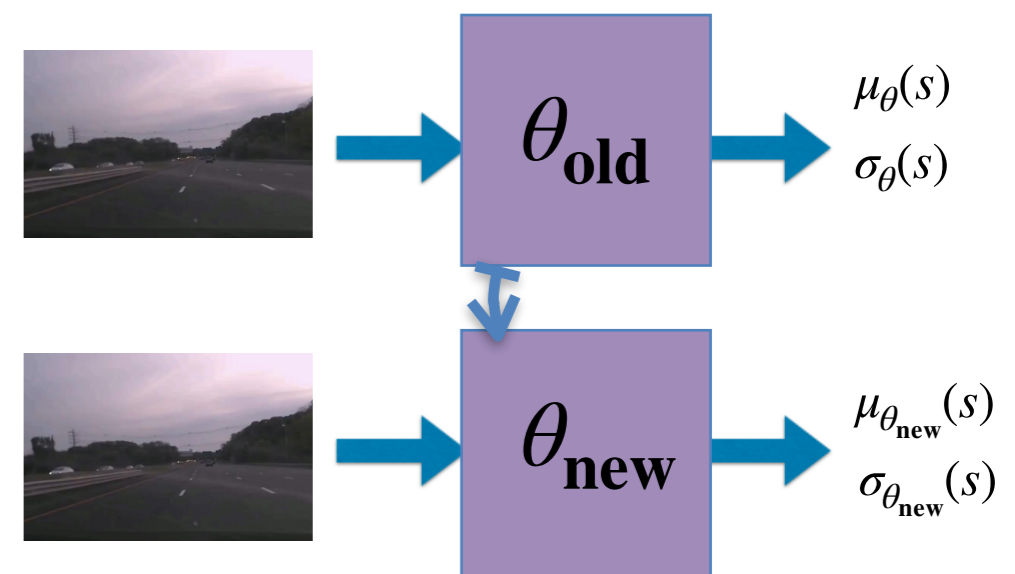
# Policy Gradients

Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta + \epsilon \cdot \hat{g}$
5. GOTO 1

This lecture is all about the stepwise



# What is the underlying objective function?

Policy gradients:

$$\hat{g} \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\alpha_t^{(i)} | s_t^{(i)}) A(s_t^{(i)}, a_t^{(i)}), \quad \tau_i \sim \pi_{\theta}$$

What is our objective? Result from differentiating the objective function:

$$J^{PG}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \log \pi_{\theta}(\alpha_t^{(i)} | s_t^{(i)}) A(s_t^{(i)}, a_t^{(i)}) \quad \tau_i \sim \pi_{\theta}$$

Is this our objective? We cannot both maximize over a variable and sample from it.

Well, we cannot optimize it too far, our advantage estimates are from samples of  $\pi_{\theta_{\text{old}}}$ . However, this constraint of “cannot optimize too far from  $\theta_{\text{old}}$ ” does not appear anywhere in the objective.

Compare to supervised learning and maximum likelihood estimation (MLE). Imagine we have access to expert actions, then the loss function we want to optimize is:

$$J^{SL}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \log \pi_{\theta}(\tilde{\alpha}_t^{(i)} | s_t^{(i)}), \quad \tau_i \sim \pi^* \quad +\text{regularization}$$

which maximizes the probability of expert actions in the training set.

Is this our SL objective?

Well, as a matter of fact, we care about test error, but this is a long story, the short answer is yes, this is good enough for us to optimize if we regularize.

# Policy Gradients

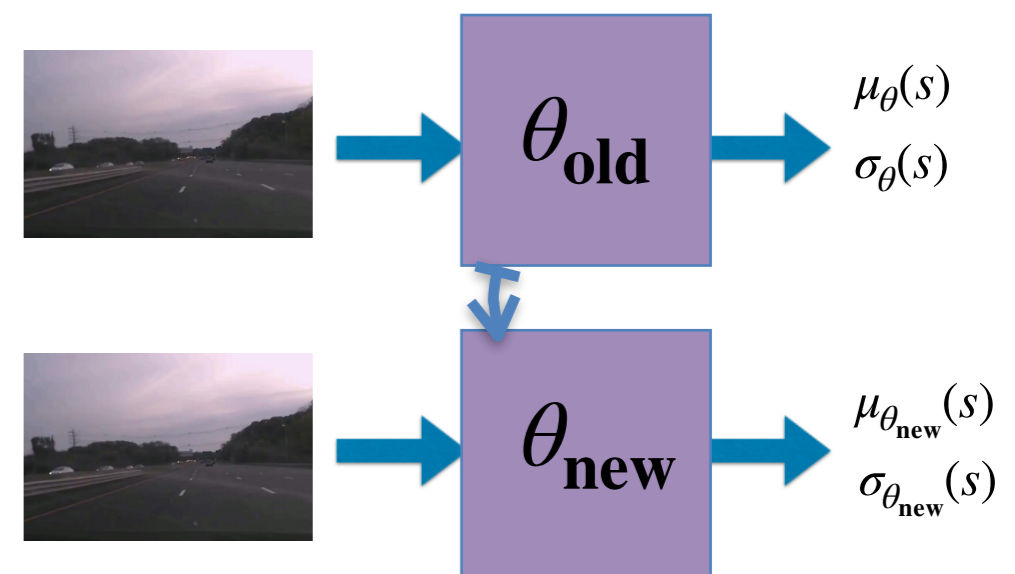
Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta + \epsilon \cdot \hat{g}$
5. GOTO 1

This lecture is all about the stepwise

It is also about writing down an objective that we can optimize with PG, and the procedure 1,2,3,4,5 will be the result of this objective maximization



# Policy Gradients

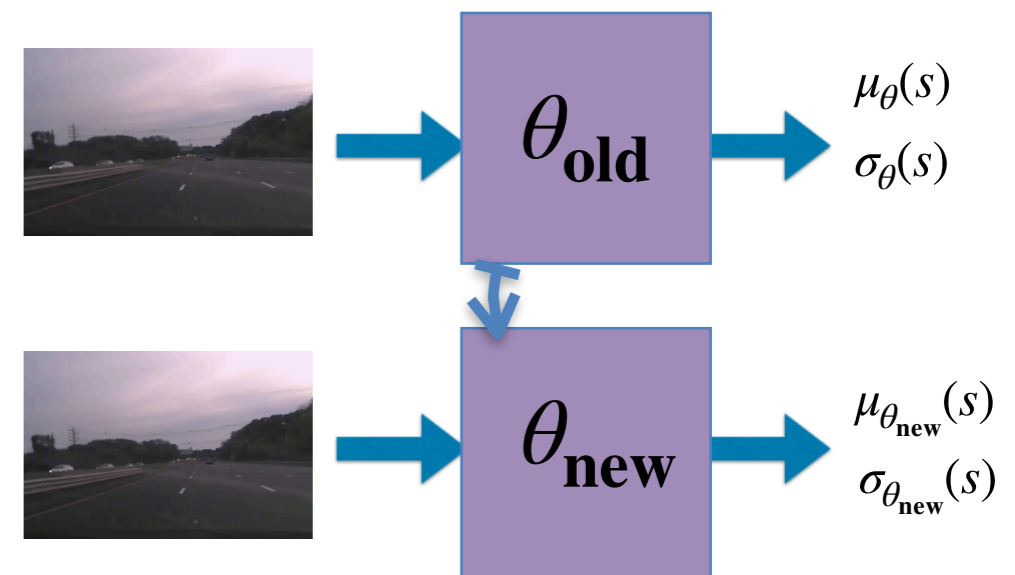
Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta + \epsilon \cdot \hat{g}$
5. GOTO 1

Two problems with the vanilla formulation:

1. Hard to choose stepwise  $\epsilon$
2. Sample inefficient: we cannot use data collected with policies of previous iterations



# Hard to choose stepsizes

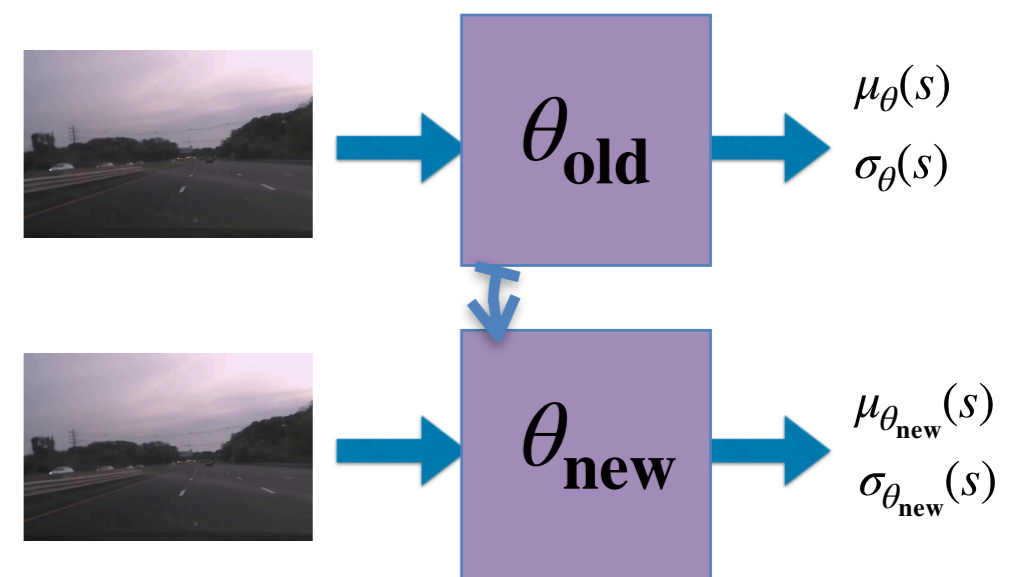
Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta + \epsilon \cdot \hat{g}$
5. GOTO 1

- Step too big  
Bad policy  $\rightarrow$  data collected under bad policy  $\rightarrow$  we cannot recover  
(in Supervised Learning, data does not depend on neural network weights)
- Step too small  
Not efficient use of experience  
(in Supervised Learning, data can be trivially re-used)

Gradient descent in parameter space does not take into account the resulting distance in the (output) policy space between  $\pi_{\theta_{old}}(s)$  and  $\pi_{\theta_{new}}(s)$





# Hard to choose stepsizes

Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta + \epsilon \cdot \hat{g}$
5. GOTO 1

Consider a family of policies with parametrization:

$$\pi_{\theta}(a) = \begin{cases} \sigma(\theta) & a = 1 \\ 1 - \sigma(\theta) & a = 2 \end{cases}$$

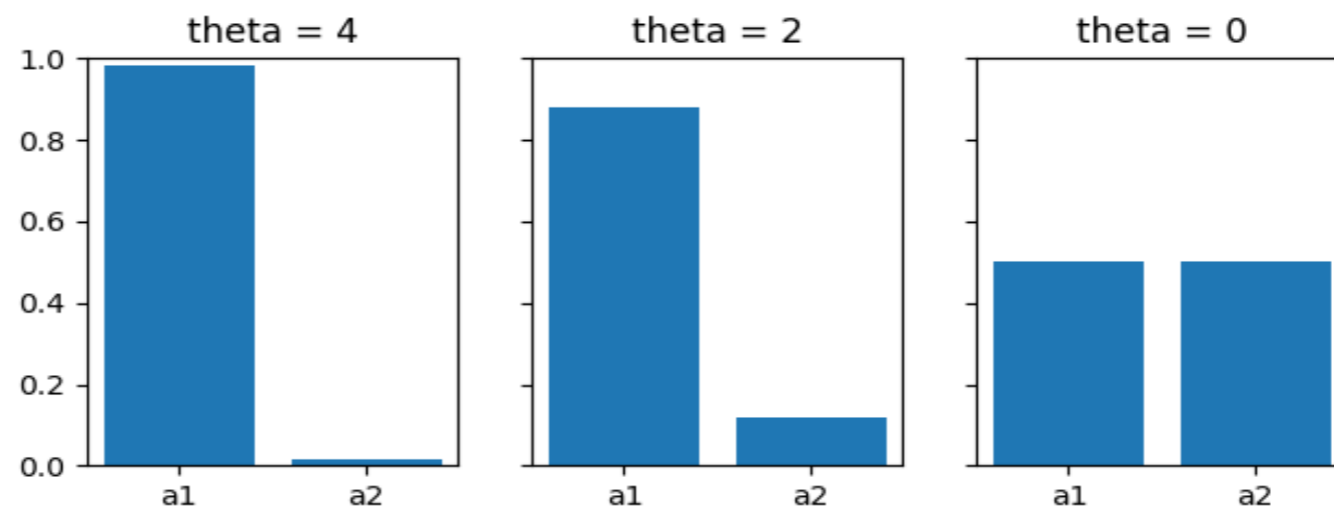


Figure: Small changes in the policy parameters can unexpectedly lead to **big** changes in the policy.

# Notation

We will use the following to denote values of parameters and corresponding policies before and after an update:

$$\theta_{old} \rightarrow \theta_{new}$$

$$\pi_{old} \rightarrow \pi_{new}$$

$$\theta \rightarrow \theta'$$

$$\pi \rightarrow \pi'$$

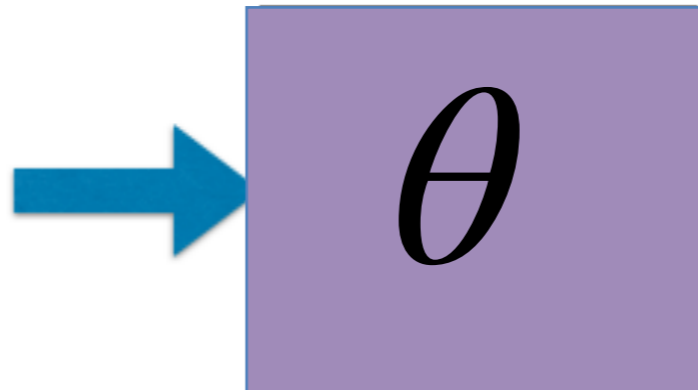
# Gradient Descent in Parameter Space

The stepwise in gradient descent results from solving the following optimization problem, e.g., using line search:

$$d^* = \arg \max_{\|d\| \leq \epsilon} J(\theta + d)$$

Euclidean distance in parameter space

It is hard to predict the result on the parameterized distribution..



$$\begin{aligned} \mu_{\theta}(s) \\ \sigma_{\theta}(s) \end{aligned}$$

SGD:  $\theta_{new} = \theta_{old} + d^*$

# Gradient Descent in Distribution Space

The stepwise in gradient descent results from solving the following optimization problem, e.g., using line search:

$$d^* = \arg \max_{\|d\| \leq \epsilon} J(\theta + d)$$

Euclidean distance in parameter space

$$\text{SGD: } \theta_{new} = \theta_{old} + d^*$$

It is hard to predict the result on the parameterized distribution.. hard to pick the threshold epsilon

**Natural gradient descent:** the stepwise in parameter space is determined by considering the KL divergence in the distributions before and after the update:

$$d^* = \arg \max_{d, \text{ s.t. } \text{KL}(\pi_\theta \| \pi_{\theta+d}) \leq \epsilon} J(\theta + d)$$

KL divergence in distribution space

Easier to pick the distance threshold!!!

$$D_{\text{KL}}(P \| Q) = \sum_i P(i) \log \left( \frac{P(i)}{Q(i)} \right)$$
$$D_{\text{KL}}(P \| Q) = \int_{-\infty}^{\infty} p(x) \log \left( \frac{p(x)}{q(x)} \right) dx$$

# Solving the KL Constrained Problem

Unconstrained penalized objective:

$$d^* = \arg \max_d J(\theta + d) - \lambda(D_{\text{KL}}[\pi_\theta \| \pi_{\theta+d}] - \epsilon)$$

First order Taylor expansion for the loss and second order for the KL:

$$\approx \arg \max_d J(\theta_{old}) + \nabla_\theta J(\theta) |_{\theta=\theta_{old}} \cdot d - \frac{1}{2} \lambda(d^\top \nabla_\theta^2 D_{\text{KL}}[\pi_{\theta_{old}} \| \pi_\theta] |_{\theta=\theta_{old}} d) + \lambda\epsilon$$

# Taylor expansion of KL

$$D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) \approx D_{\text{KL}}(p_{\theta_{old}} | p_{\theta_{old}}) + d^{\top} \nabla_{\theta} D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} + \frac{1}{2} d^{\top} \nabla_{\theta}^2 D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}}$$

$$\begin{aligned} \nabla_{\theta} D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} &= -\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta_{old}}} \log P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= -\mathbb{E}_{x \sim p_{\theta_{old}}} \nabla_{\theta} \log P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= -\mathbb{E}_{x \sim p_{\theta_{old}}} \frac{1}{P_{\theta_{old}}(x)} \nabla_{\theta} P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= \int_x P_{\theta_{old}}(x) \frac{1}{P_{\theta_{old}}(x)} \nabla_{\theta} P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= \int_x \nabla_{\theta} P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= \nabla_{\theta} \int_x P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= 0 \end{aligned}$$

$$D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) = \mathbb{E}_{x \sim p_{\theta_{old}}} \log \left( \frac{P_{\theta_{old}}(x)}{P_{\theta}(x)} \right)$$

# Taylor expansion of KL

$$D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) \approx D_{\text{KL}}(p_{\theta_{old}} | p_{\theta_{old}}) + d^{\top} \nabla_{\theta} \text{KL}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} + \frac{1}{2} d^{\top} \nabla_{\theta}^2 D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} d$$

$$\begin{aligned} \nabla_{\theta}^2 D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} &= -\mathbb{E}_{x \sim p_{\theta_{old}}} \nabla_{\theta}^2 \log P_{\theta}(x) |_{\theta=\theta_{old}} \\ &= -\mathbb{E}_{x \sim p_{\theta_{old}}} \nabla_{\theta} \left( \frac{\nabla_{\theta} P_{\theta}(x)}{P_{\theta}(x)} \right) |_{\theta=\theta_{old}} \\ &= -\mathbb{E}_{x \sim p_{\theta_{old}}} \left( \frac{\nabla_{\theta}^2 P_{\theta}(x) P_{\theta}(x) - \nabla_{\theta} P_{\theta}(x) \nabla_{\theta} P_{\theta}(x)^{\top}}{P_{\theta}(x)^2} \right) |_{\theta=\theta_{old}} \\ &= -\mathbb{E}_{x \sim p_{\theta_{old}}} \frac{\nabla_{\theta}^2 P_{\theta}(x) |_{\theta=\theta_{old}}}{P_{\theta_{old}}(x)} + \mathbb{E}_{x \sim p_{\theta_{old}}} \nabla_{\theta} \log P_{\theta}(x) \nabla_{\theta} \log P_{\theta}(x)^{\top} |_{\theta=\theta_{old}} \\ &= \mathbb{E}_{x \sim p_{\theta_{old}}} \nabla_{\theta} \log P_{\theta}(x) \nabla_{\theta} \log P_{\theta}(x)^{\top} |_{\theta=\theta_{old}} \end{aligned}$$

$$D_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) = \mathbb{E}_{x \sim p_{\theta_{old}}} \log \left( \frac{P_{\theta_{old}}(x)}{P_{\theta}(x)} \right)$$

# Fisher Information Matrix

Exactly equivalent to the Hessian of KL divergence!

$$\mathbf{F}(\theta) = \mathbb{E}_{\theta} \left[ \nabla_{\theta} \log p_{\theta}(x) \nabla_{\theta} \log p_{\theta}(x)^{\top} \right]$$

$$\mathbf{F}(\theta_{old}) = \nabla_{\theta}^2 \mathbf{D}_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}}$$

$$\begin{aligned} \mathbf{D}_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) &\approx \mathbf{D}_{\text{KL}}(p_{\theta_{old}} | p_{\theta_{old}}) + d^{\top} \nabla_{\theta} \mathbf{D}_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} + \frac{1}{2} d^{\top} \nabla_{\theta}^2 \mathbf{D}_{\text{KL}}(p_{\theta_{old}} | p_{\theta}) |_{\theta=\theta_{old}} d \\ &= \frac{1}{2} d^{\top} \mathbf{F}(\theta_{old}) d \\ &= \frac{1}{2} (\theta - \theta_{old})^{\top} \mathbf{F}(\theta_{old}) (\theta - \theta_{old}) \end{aligned}$$

Since KL divergence is roughly analogous to a distance measure between distributions, Fisher information serves as a local distance metric between distributions: how much you change the distribution if you move the parameters a little bit in a given direction.



# Solving the KL Constrained Problem

Unconstrained penalized objective:

$$d^* = \arg \max_d J(\theta + d) - \lambda(D_{\text{KL}}[\pi_\theta \| \pi_{\theta+d}] - \epsilon)$$

First order Taylor expansion for the loss and second order for the KL:

$$\approx \arg \max_d J(\theta_{old}) + \nabla_\theta J(\theta) |_{\theta=\theta_{old}} \cdot d - \frac{1}{2} \lambda (d^\top \nabla_\theta^2 D_{\text{KL}}[\pi_{\theta_{old}} \| \pi_\theta] |_{\theta=\theta_{old}} d) + \lambda \epsilon$$

Substitute for the information matrix:

$$\begin{aligned} &= \arg \max_d \nabla_\theta J(\theta) |_{\theta=\theta_{old}} \cdot d - \frac{1}{2} \lambda (d^\top \mathbf{F}(\theta_{old}) d) \\ &= \arg \min_d - \nabla_\theta J(\theta) |_{\theta=\theta_{old}} \cdot d + \frac{1}{2} \lambda (d^\top \mathbf{F}(\theta_{old}) d) \end{aligned}$$

# Natural Gradient Descent

Setting the gradient to zero:

$$\begin{aligned} 0 &= \frac{\partial}{\partial d} \left( -\nabla_{\theta} J(\theta) \Big|_{\theta=\theta_{old}} \cdot d + \frac{1}{2} \lambda (d^{\top} \mathbf{F}(\theta_{old}) d) \right) \\ &= -\nabla_{\theta} J(\theta) \Big|_{\theta=\theta_{old}} + \frac{1}{2} \lambda (\mathbf{F}(\theta_{old})) d \end{aligned}$$

$$d = \frac{2}{\lambda} \mathbf{F}^{-1}(\theta_{old}) \nabla_{\theta} J(\theta) \Big|_{\theta=\theta_{old}}$$

**The natural gradient:**

$$\tilde{\nabla} J(\theta) = \mathbf{F}^{-1}(\theta_{old}) \nabla_{\theta} J(\theta)$$

$$\theta_{new} = \theta_{old} + \alpha \cdot \mathbf{F}^{-1}(\theta_{old}) \hat{g}$$

$$D_{\text{KL}}(\pi_{\theta_{old}} | \pi_{\theta}) \approx \frac{1}{2} (\theta - \theta_{old})^{\top} \mathbf{F}(\theta_{old}) (\theta - \theta_{old})$$

$$\frac{1}{2} (\alpha g_N)^{\top} \mathbf{F}(\alpha g_N) = \epsilon$$

$$\alpha = \sqrt{\frac{2\epsilon}{(g_N^{\top} \mathbf{F} g_N)}}$$

# Natural Gradient Descent

---

## Algorithm 1 Natural Policy Gradient

---

Input: initial policy parameters  $\theta_0$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)
- and KL-divergence Hessian / Fisher Information Matrix  $\hat{H}_k$

Compute Natural Policy Gradient update:

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\epsilon}{\hat{g}_k^T \hat{H}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$$

**end for**

---

Both use samples from the current policy  $\pi_k$

# Natural Gradient Descent

---

**Algorithm 1** Natural Policy Gradient

---

Input: initial policy parameters  $\theta_0$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)
- and KL-divergence Hessian / Fisher Information Matrix  $\hat{H}_k$

Compute Natural Policy Gradient update:

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\epsilon}{\hat{g}_k^T \hat{H}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$$

**end for**

---

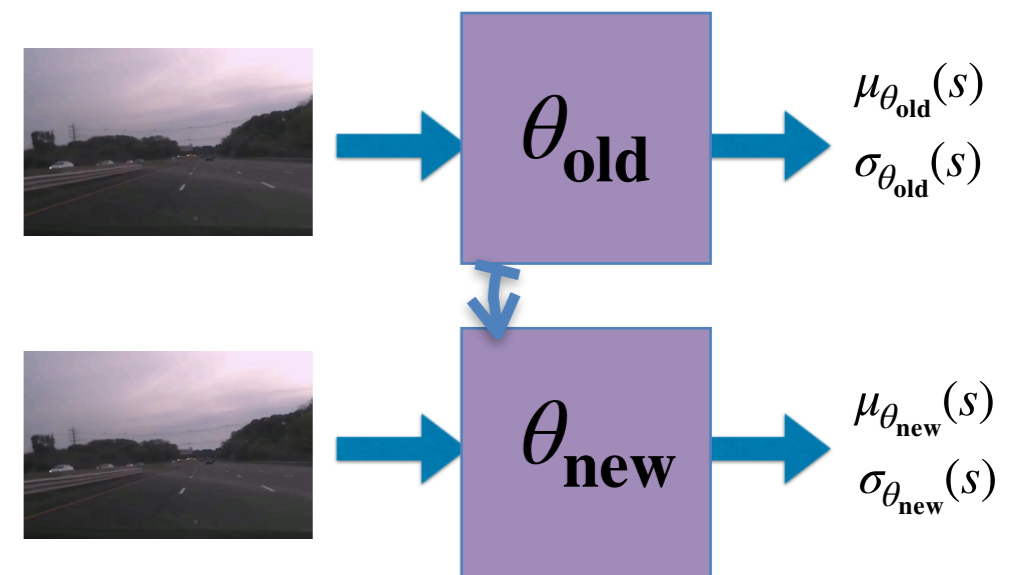
very expensive to compute for a large number of parameters!

# Policy Gradients

Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta_{old}}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta_{old} + \epsilon \cdot \hat{g}$
5. GOTO 1



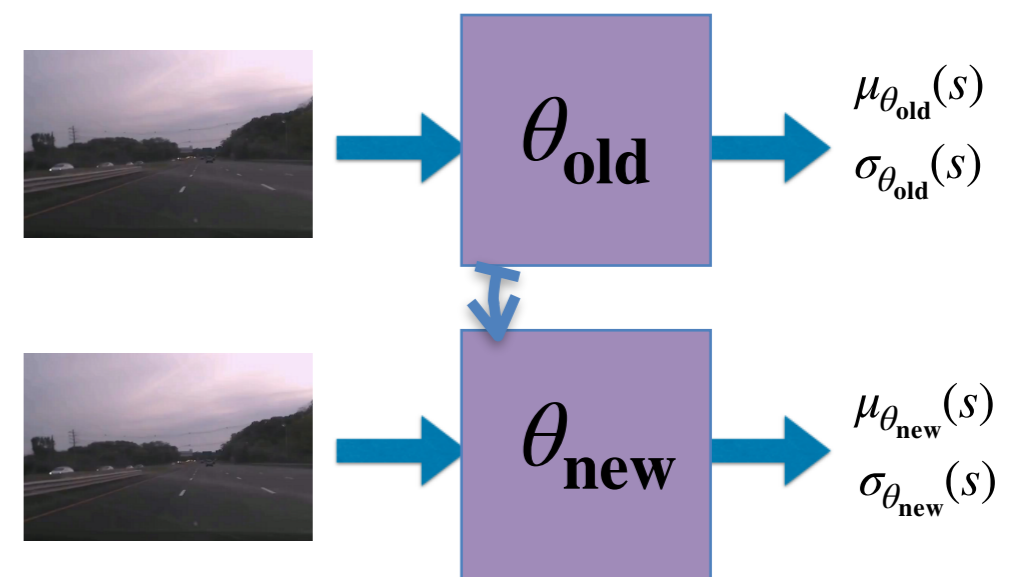
# Policy Gradients

Monte Carlo Policy Gradients (REINFORCE), gradient direction:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t \right]$

Actor-Critic Policy Gradient:  $\hat{g} = \hat{\mathbb{E}}_t \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_w(s_t) \right]$

1. Collect trajectories for policy  $\pi_{\theta_{old}}$
2. Estimate advantages  $A$
3. Compute policy gradient  $\hat{g}$
4. Update policy parameters  $\theta_{new} = \theta_{old} + \epsilon \cdot \hat{g}$
5. GOTO 1

- On policy learning can be extremely inefficient
- The policy changes only a little bit with each gradient step
- I want to be able to use earlier data..how to do that?



# Off policy learning with Importance Sampling

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau)] \\ &= \sum_{\tau} \pi_{\theta}(\tau) R(\tau) \\ &= \sum_{\tau} \pi_{\theta_{old}}(\tau) \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau) \\ &= \sum_{\tau \sim \pi_{\theta_{old}}} \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau) \\ &= \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau) \end{aligned}$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau)$$

$$\nabla_{\theta} J(\theta) |_{\theta=\theta_{old}} = \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \nabla_{\theta} \log \pi_{\theta}(\tau) |_{\theta=\theta_{old}} R(\tau)$$

<- Gradient evaluated at theta\_old is unchanged

# Off policy learning with Importance Sampling

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} [R(\tau)]$$

$$= \sum_{\tau} \pi_{\theta}(\tau) R(\tau)$$

$$= \sum_{\tau} \pi_{\theta_{old}}(\tau) \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau)$$

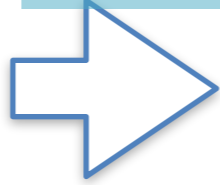
$$= \sum_{\tau \sim \pi_{\theta_{old}}} \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau)$$

$$= \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} R(\tau)$$

$$\nabla_{\theta} J(\theta) |_{\theta=\theta_{old}} = \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \nabla_{\theta} \log \pi_{\theta}(\tau) |_{\theta=\theta_{old}} R(\tau)$$

$$\frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)} = \prod_{i=1}^T \frac{\pi_{\theta}(a_i | s_i)}{\pi_{\theta_{old}}(a_i | s_i)}$$



$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \sum_{t=1}^T \prod_{t'=1}^t \frac{\pi_{\theta}(a'_{t'} | s'_{t'})}{\pi_{\theta_{old}}(a'_{t'} | s'_{t'})} \hat{A}_t$$

Now we can use data from the old policy, but the variance has increased by a lot! Those multiplications can explode or vanish!



# Trust region Policy Optimization

$$\begin{aligned} & \underset{\theta}{\text{maximize}} && \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] \\ & \text{subject to} && \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]] \leq \delta. \end{aligned}$$

- ▶ Also worth considering using a penalty instead of a constraint

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] - \beta \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]]$$

Again the KL penalized problem!

# Trust region Policy Optimization

- ▶ maximize $_{\theta} L_{\pi_{\theta_{\text{old}}}}(\pi_{\theta}) - \beta \cdot \overline{\text{KL}}_{\pi_{\theta_{\text{old}}}}(\pi_{\theta})$
- ▶ Make linear approximation to  $L_{\pi_{\theta_{\text{old}}}}$  and quadratic approximation to KL term:

$$\text{maximize}_{\theta} \quad g \cdot (\theta - \theta_{\text{old}}) - \frac{\beta}{2} (\theta - \theta_{\text{old}})^T F (\theta - \theta_{\text{old}})$$

$$\text{where } g = \frac{\partial}{\partial \theta} L_{\pi_{\theta_{\text{old}}}}(\pi_{\theta}) \Big|_{\theta=\theta_{\text{old}}}, \quad F = \frac{\partial^2}{\partial^2 \theta} \overline{\text{KL}}_{\pi_{\theta_{\text{old}}}}(\pi_{\theta}) \Big|_{\theta=\theta_{\text{old}}}$$

Exactly what we saw with natural policy gradient!

One important detail!

# Trust region Policy Optimization

Due to the quadratic approximation, the KL constraint may be violated! What if we just do a line search to find the best stepsize, making sure:

- I am improving my objective  $J(\theta)$
- The KL constraint is not violated!

$$\begin{aligned} & \underset{\theta}{\text{maximize}} && \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] \\ & \text{subject to} && \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]] \leq \delta. \end{aligned}$$

---

## Algorithm 2 Line Search for TRPO

---

Compute proposed policy step  $\Delta_k = \sqrt{\frac{2\delta}{\hat{g}_k^T \hat{H}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$

**for**  $j = 0, 1, 2, \dots, L$  **do**

    Compute proposed update  $\theta = \theta_k + \alpha^j \Delta_k$

**if**  $\mathcal{L}_{\theta_k}(\theta) \geq 0$  and  $\bar{D}_{KL}(\theta || \theta_k) \leq \delta$  **then**

        accept the update and set  $\theta_{k+1} = \theta_k + \alpha^j \Delta_k$

        break

**end if**

**end for**

---

# Trust region Policy Optimization

TRPO= NPG +LineSearch

---

## Algorithm 3 Trust Region Policy Optimization

---

Input: initial policy parameters  $\theta_0$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)
- and KL-divergence Hessian-vector product function  $f(v) = \hat{H}_k v$

Use CG with  $n_{cg}$  iterations to obtain  $x_k \approx \hat{H}_k^{-1} \hat{g}_k$

Estimate proposed step  $\Delta_k \approx \sqrt{\frac{2\delta}{x_k^T \hat{H}_k x_k}} x_k$

Perform backtracking line search with exponential decay to obtain final update

$$\theta_{k+1} = \theta_k + \alpha^j \Delta_k$$

**end for**

---

# Trust region Policy Optimization

TRPO= NPG +Line search+monotonic improvement theorem!

---

## Algorithm 3 Trust Region Policy Optimization

---

Input: initial policy parameters  $\theta_0$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)
- and KL-divergence Hessian-vector product function  $f(v) = \hat{H}_k v$

Use CG with  $n_{cg}$  iterations to obtain  $x_k \approx \hat{H}_k^{-1} \hat{g}_k$

Estimate proposed step  $\Delta_k \approx \sqrt{\frac{2\delta}{x_k^T \hat{H}_k x_k}} x_k$

Perform backtracking line search with exponential decay to obtain final update

$$\theta_{k+1} = \theta_k + \alpha^j \Delta_k$$

**end for**

---

# Relating objectives of two policies

Policy objective:

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=0}^{\infty} \gamma^t r_t$$

Policy objective can be written in terms of old one:

$$J(\pi_{\theta'}) - J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_{\theta'}} \sum_{t=0}^{\infty} \gamma^t A^{\pi_\theta}(s_t, a_t)$$

Equivalently for succinctness:

$$J(\pi') - J(\pi) = \mathbb{E}_{\tau \sim \pi'} \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t)$$

# Relating objectives of two policies

$$\begin{aligned} J(\pi') - J(\pi) &= \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \right] \\ &= \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t, s_{t+1}) + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)) \right] \\ &= J(\pi') + \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} V^\pi(s_{t+1}) - \sum_{t=0}^{\infty} \gamma^t V^\pi(s_t) \right] \\ &= J(\pi') + \mathbb{E}_{\tau \sim \pi'} \left[ \sum_{t=1}^{\infty} \gamma^t V^\pi(s_t) - \sum_{t=0}^{\infty} \gamma^t V^\pi(s_t) \right] \\ &= J(\pi') - \mathbb{E}_{\tau \sim \pi'} [V^\pi(s_0)] \\ &= J(\pi') - J(\pi) \end{aligned}$$

The initial state distribution is the same for both!

# Relating objectives of two policies

Discounted state visitation distribution:

$$\begin{aligned}d^\pi(s) &= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi) \\J(\pi') - J(\pi) &= \mathbb{E}_{\tau \sim \pi'} \sum_{t=0}^{\infty} \gamma^t A^\pi(s_t, a_t) \\&= \mathbb{E}_{s \sim d^{\pi'}, a \sim \pi'} A^\pi(s, a) \\&= \mathbb{E}_{s \sim d^{\pi'}, a \sim \pi} \left[ \frac{\pi'(a | s)}{\pi(a | s)} A^\pi(s, a) \right]\end{aligned}$$

But how are we supposed to sample states from the policy we are trying to optimize for...  
Let's use the previous policy to sample them.

$$\begin{aligned}J(\pi') - J(\pi) &\approx \mathbb{E}_{s \sim d^\pi, a \sim \pi} \frac{\pi'(a | s)}{\pi(a | s)} A^\pi(s, a) \\&= \mathcal{L}_\pi(\pi')\end{aligned}$$

It turns out we can bound this approximation error:

$$|J(\pi') - (J(\pi) + \mathcal{L}_\pi(\pi'))| \leq C \sqrt{\mathbb{E}_{s \sim d^\pi} [D_{KL}(\pi' || \pi)[s]]}$$



# Relating objectives of two policies

$$\begin{aligned}\mathcal{L}_{\pi}^{\pi'} &= \mathbb{E}_{s \sim d^{\pi}, a \sim \pi} \left[ \frac{\pi'(a | s)}{\pi(a | s)} A^{\pi}(s, a) \right] \\ &= \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \frac{\pi'(a_t | s_t)}{\pi(a_t | s_t)} A^{\pi}(s_t, a_t) \right]\end{aligned}$$

This is something we can optimize using trajectories from the old policy!

Compare to Importance Sampling:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \sum_{t=1}^T \prod_{t'=1}^t \frac{\pi_{\theta}(a_{t'} | s_{t'})}{\pi_{\theta_{old}}(a_{t'} | s_{t'})} \hat{A}_t$$

Now we do not have the product! So, the gradient will have much smaller variance! (Yes, but we have approximated, that's why!) What is the gradient?

$$\begin{aligned}\nabla_{\theta} \mathcal{L}_{\theta_k}^{\theta} |_{\theta=\theta_k} &= \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[ \sum_{t=0}^{\infty} \gamma^t \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t) |_{\theta=\theta_k}}{\pi_{\theta_k}(a_t | s_t)} A^{\pi_{\theta_k}}(s_t, a_t) \right] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[ \sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta=\theta_k} A^{\pi_{\theta_k}}(s_t, a_t) \right]\end{aligned}$$

# Monotonic Improvement Theorem

$$|J(\pi') - (J(\pi) + \mathcal{L}_\pi(\pi'))| \leq C \sqrt{\mathbb{E}_{s \sim d^\pi} [\text{KL}(\pi' | \pi)[s]]}$$

$$\Rightarrow J(\pi') - J(\pi) \geq \mathcal{L}_\pi(\pi') - C \sqrt{\mathbb{E}_{s \sim d^\pi} [\text{KL}(\pi' | \pi)[s]]}$$

Given policy  $\pi$ , we want to optimize over policy  $\pi'$  to maximize  $J(\pi')$ .

- If we maximize the RHS we are guaranteed to maximize the LHS.
- We know how to maximize the RHS. I can estimate both quantities of  $\pi'$  with sampled from  $\pi$ .
- But will I have a better policy  $\pi'$ ? (knowing that the distance of the objectives is maximized is not enough, there needs to be positive or equal to zero)

# Monotonic Improvement Theorem

Proof of improvement guarantee: Suppose  $\pi_{k+1}$  and  $\pi_k$  are related by

$$\pi_{k+1} = \arg \max_{\pi'} \mathcal{L}_{\pi_k}(\pi') - C \sqrt{\mathbb{E}_{s \sim d^{\pi_k}} [D_{KL}(\pi' || \pi_k)[s]]}.$$

- $\pi_k$  is a feasible point, and the objective at  $\pi_k$  is equal to 0.
  - $\mathcal{L}_{\pi_k}(\pi_k) \propto \mathbb{E}_{s, a \sim d^{\pi_k, \pi_k}} [A^{\pi_k}(s, a)] = 0$
  - $D_{KL}(\pi_k || \pi_k)[s] = 0$
- $\implies$  optimal value  $\geq 0$
- $\implies$  by the performance bound,  $J(\pi_{k+1}) - J(\pi_k) \geq 0$

# Approximate Monotonic Improvement

- Theory is very conservative (high value of  $C$ ) and we will use KL distance of  $\pi'$  and  $\pi$  as a constraint (trust region) as opposed to a penalty:

$$\begin{aligned} \pi_{k+1} &= \arg \max_{\pi'} \mathcal{L}_{\pi_k}(\pi') \\ \text{s.t. } & \mathbb{E}_{s \sim d^{\pi_k}} [D_{KL}(\pi' || \pi_k)[s]] \leq \delta \end{aligned}$$

# Trust region Policy Optimization

TRPO= NPG +Line search+monotonic improvement theorem!

---

## Algorithm 3 Trust Region Policy Optimization

---

Input: initial policy parameters  $\theta_0$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)
- and KL-divergence Hessian-vector product function  $f(v) = \hat{H}_k v$

Use CG with  $n_{cg}$  iterations to obtain  $x_k \approx \hat{H}_k^{-1} \hat{g}_k$

Estimate proposed step  $\Delta_k \approx \sqrt{\frac{2\delta}{x_k^T \hat{H}_k x_k}} x_k$

Perform backtracking line search with exponential decay to obtain final update

$$\theta_{k+1} = \theta_k + \alpha^j \Delta_k$$

**end for**

---

# Proximal Policy Optimization

Can I achieve similar performance without second order information (no Fisher matrix!)

- Adaptive KL Penalty

- Policy update solves unconstrained optimization problem

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

- Penalty coefficient  $\beta_k$  changes between iterations to approximately enforce KL-divergence constraint

- Clipped Objective

- New objective function: let  $r_t(\theta) = \pi_{\theta}(a_t|s_t)/\pi_{\theta_k}(a_t|s_t)$ . Then

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

where  $\epsilon$  is a hyperparameter (maybe  $\epsilon = 0.2$ )

- Policy update is  $\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$

# PPO: Adaptive KL Penalty

Input: initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

by taking  $K$  steps of minibatch SGD (via Adam)

**if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \geq 1.5\delta$  **then**

$$\beta_{k+1} = 2\beta_k$$

**else if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \leq \delta/1.5$  **then**

$$\beta_{k+1} = \beta_k/2$$

**end if**

**end for**

Don't use second order approximation for KL which is expensive, use standard gradient descent

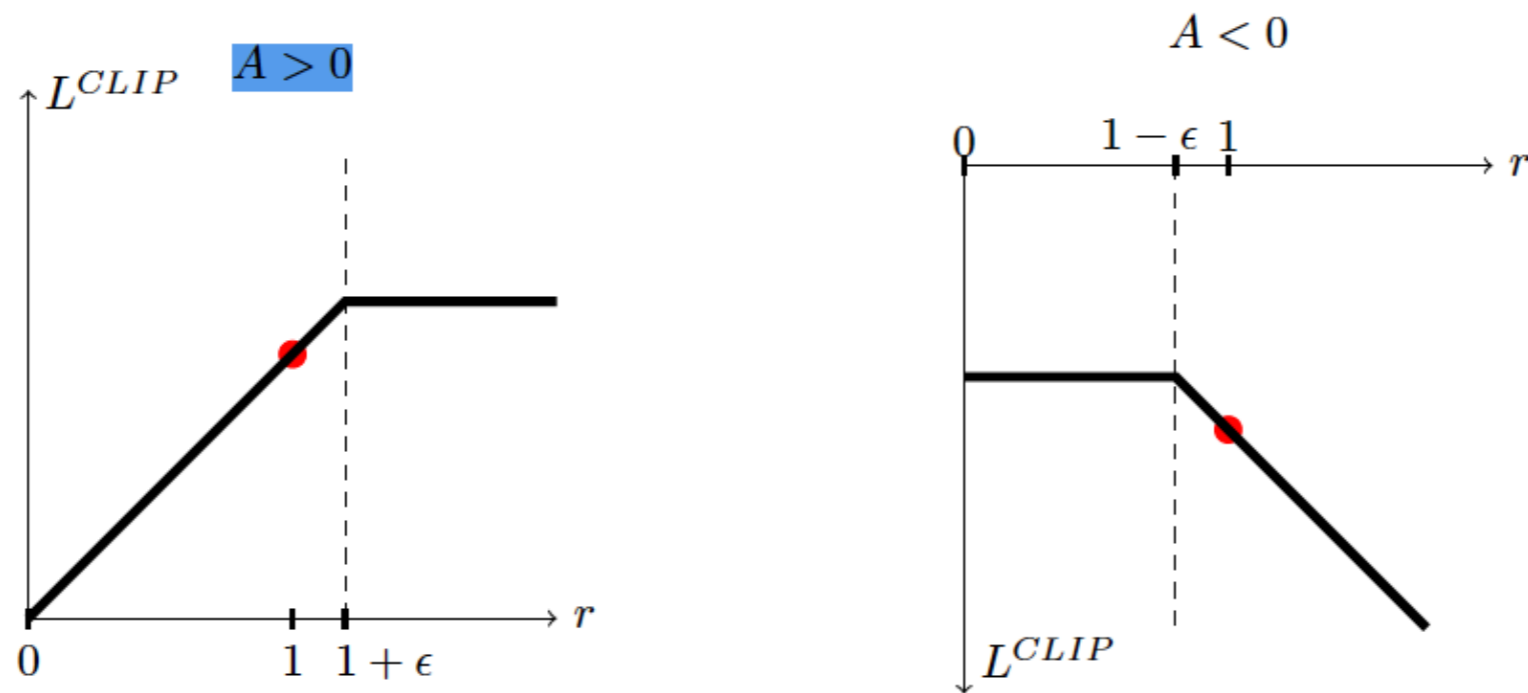
# PPO: Clipped Objective

- ▶ Recall the surrogate objective

$$L^{IS}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t \left[ r_t(\theta) \hat{A}_t \right]. \quad (1)$$

- ▶ Form a lower bound via clipped importance ratios

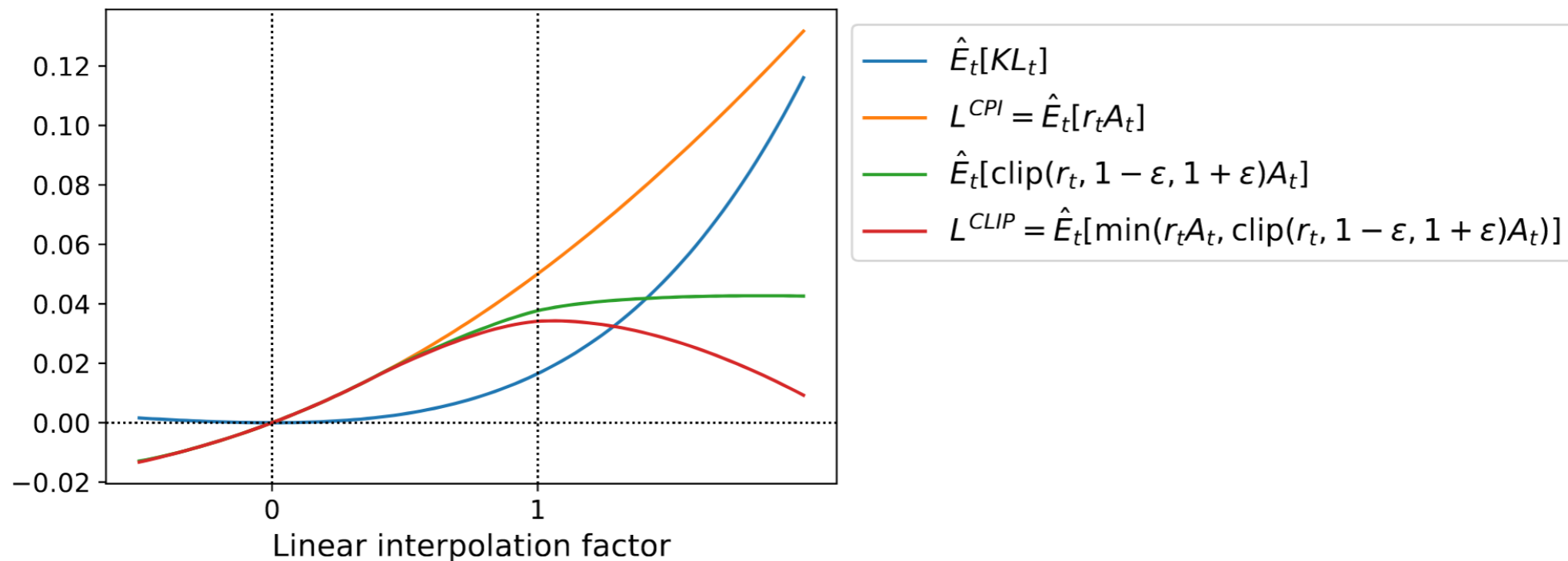
$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right] \quad (2)$$





# PPO: Clipped Objective

But *how* does clipping keep policy close? By making objective as pessimistic as possible about performance far away from  $\theta_k$ :



**Figure:** Various objectives as a function of interpolation factor  $\alpha$  between  $\theta_{k+1}$  and  $\theta_k$  after one update of PPO-Clip <sup>9</sup>

# PPO: Clipped Objective

Input: initial policy parameters  $\theta_0$ , clipping threshold  $\epsilon$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$$

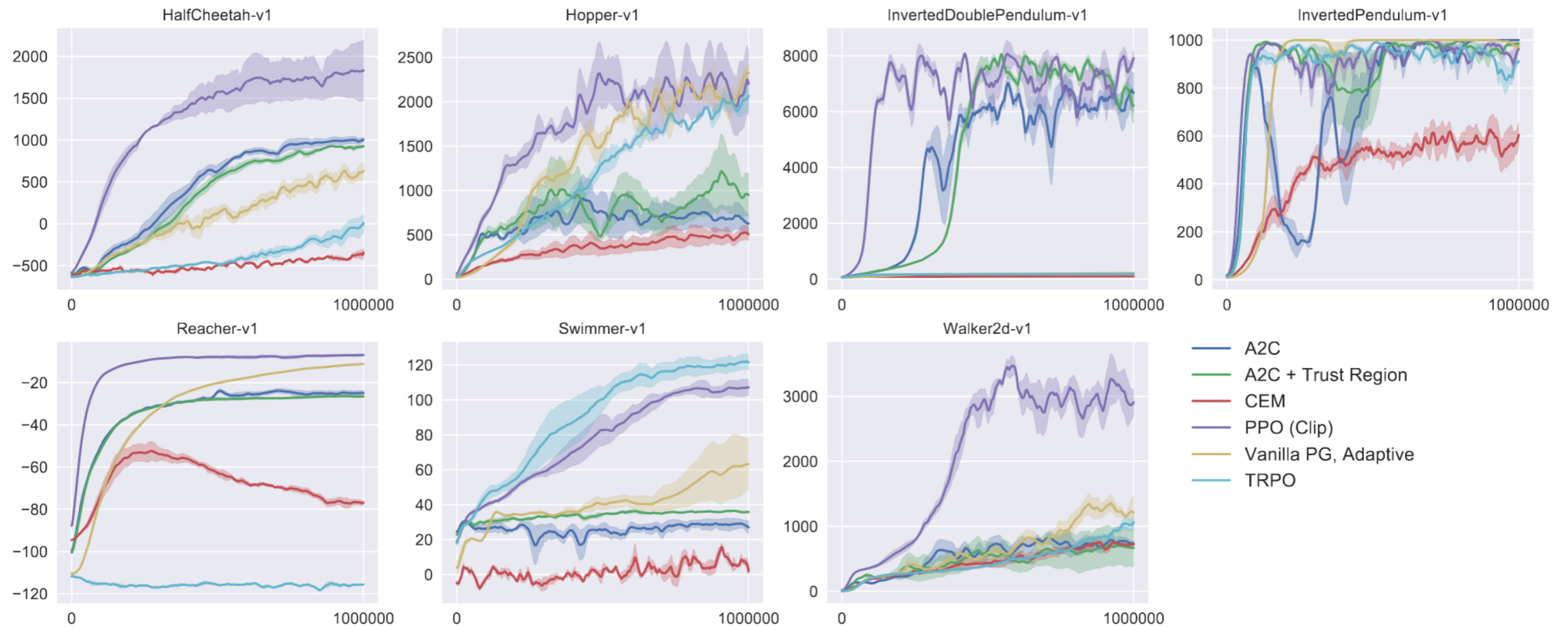
by taking  $K$  steps of minibatch SGD (via Adam), where

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

**end for**

- 
- Clipping prevents policy from having incentive to go far away from  $\theta_{k+1}$
  - Clipping seems to work at least as well as PPO with KL penalty, but is simpler to implement

# PPO: Clipped Objective



**Figure:** Performance comparison between PPO with clipped objective and various other deep RL methods on a slate of MuJoCo tasks. <sup>10</sup>

# Summary

- Gradient Descent in Parameter VS distribution space
- Natural gradients: we need to keep track of how the KL changes from iteration to iteration
- Natural policy gradients
- Clipped objective works well

## Related Readings

S. Kakade. "A Natural Policy Gradient." *NIPS*. 2001

S. Kakade and J. Langford. "Approximately optimal approximate reinforcement learning". *ICML*. 2002

J. Peters and S. Schaal. "Natural actor-critic". *Neurocomputing* (2008)

J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel. "Trust Region Policy Optimization". *ICML* (2015)

J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. "Proximal Policy Optimization Algorithms". (2017)

J. Achiam, D. Held, A. Tamar, P. Abbeel "Constrained Policy Optimization". (2017)