

Modern Policy Gradient Methods

Miguel Calvo-Fullana and Santiago Paternain
Electrical and Systems Engineering, University of Pennsylvania
{spater,cfullana}@seas.upenn.edu

November 21 —December 5, 2019

Classical Policy Gradients

Monotonic Improvement Guarantees

Natural Policy Gradients and Trust Region Policy Optimization (TRPO)

Proximal Policy Optimization (PPO)

- ▶ A Markov Decision Process is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R}, P)$
- ▶ P is a Markov transition probability if for any $\mathcal{S}' \subseteq \mathcal{S}$ and $\mathcal{R}' \subseteq \mathcal{R}$

$$P[\mathcal{S}_{t+1} \in \mathcal{S}', R_{t+1} \in \mathcal{R}' | \mathcal{S}_t, A_t, \dots, \mathcal{S}_0, A_0] = P[\mathcal{S}_{t+1} \in \mathcal{S}', R_{t+1} \in \mathcal{R}' | \mathcal{S}_t, A_t]$$
- ▶ We select the actions based on parameterized policies $\pi_\theta(a|s)$
 - ⇒ We cannot work with general continuous functions
 - ⇒ Parameterization is necessary
- ▶ Find the best policy within the functions that our parameterization defines
 - ⇒ “Best” is defined by the value function

$$v_{\pi_\theta}(s) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} | \mathcal{S}_0 = s \right]$$

- ⇒ Recall that the expectation is with respect to all the rewards seen
- ⇒ We write \mathbb{E}_{π_θ} to denote the we are following the policy $\pi_\theta(a|s)$

- ▶ Expanding the expectation, the value function is written as

$$v_{\pi_{\theta}}(s) = \sum_{k=0}^{\infty} \int_{\mathcal{R}^k \mathcal{A}^k \mathcal{S}^{k-1}} \gamma^k r_k \prod_{j=0}^{k-1} p(s_{j+1}, r_{j+1} | s_j, a_j) \pi_{\theta}(a_j | s_j) ds_k da_{k-1} dr_k$$

- ▶ Where the policy is **parameterized** by $\theta \in \mathbb{R}^d$
- ▶ We will update the parameters via gradient ascent
 - ⇒ This gradient is given by the policy gradient theorem
- ▶ The gradient of v with respect to $\theta \in \mathbb{R}^d$ is given by

$$\nabla_{\theta} v(\theta) = \left(\frac{\partial v(\theta)}{\partial \theta_1}, \dots, \frac{\partial v(\theta)}{\partial \theta_d} \right)^T$$

- ▶ To find the maximum, we update the parameters θ via

$$\theta_{k+1} = \theta_k + \alpha_k \nabla_{\theta} v(\theta), \text{ with } \alpha_k > 0$$

- ▶ Policy gradient for Episodic Tasks

$$\nabla_{\theta} v(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) G_t \right]$$

- ▶ Policy gradient for Continuing Tasks

$$\nabla_{\theta} v_{\pi_{\theta}}(s_i) = (1 - \gamma)^{-1} \mathbb{E}_{S \sim \rho_{\theta}, A \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(A | S) q_{\pi_{\theta}}(S, A)]$$

$$\rho_{\theta}(s_i, s') = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(S_0 = s_i \rightarrow S_t = s')$$

- ▶ There are **several issues** with the classical policy gradient
 - ⇒ We have attempted to address some of them
 - ⇒ Via baselines and actor-critic methods
- ▶ There are still issues that we have not tackled
- ▶ Modern methods for policy gradients attempt to tackle these issues
- ▶ We measure distance in the space of parameters
 - ⇒ Parameters that are close can generate very different distributions
 - ⇒ Especially problematic with deep neural networks
 - ⇒ Makes the choice of **step size difficult**
 - ⇒ We can attempt to characterize distances in the space of distributions
- ▶ We can be more efficient by **reusing data** from previous trajectories
 - ⇒ We have studied some of this for off-policy methods
 - ⇒ We can attempt to do so via **importance sampling**

- ▶ The Kullback-Leibler (KL) divergence is an **asymmetric** measure of difference between probability distributions
- ▶ For discrete probability distributions P and Q it is defined as

$$D_{\text{KL}}(P\|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

- ▶ Properties of the KL divergence
 - $\Rightarrow D_{\text{KL}}(P\|Q) \geq 0$
 - $\Rightarrow D_{\text{KL}}(P\|P) = 0$
 - $\Rightarrow D_{\text{KL}}(P\|Q) \neq D_{\text{KL}}(Q\|P)$
- ▶ The KL divergence is **not a metric** as it is not symmetric
 - \Rightarrow Nonetheless, it is a good measure of closeness between distributions

- ▶ Importance sampling is a technique for estimating a probability distribution
⇒ With only samples from a different distribution
- ▶ For distributions P and Q and an arbitrary function f we have that

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

- ▶ We denote the ratio $P(x)/Q(x)$ as the **importance sampling weight**
- ▶ The variance of the stochastic approximation estimator is then given by

$$\begin{aligned} \text{Var} \left(\frac{P(x)}{Q(x)} f(x) \right) &= \mathbb{E}_{x \sim Q} \left[\left(\frac{P(x)}{Q(x)} f(x) \right)^2 \right] - \mathbb{E}_{x \sim Q} \left[\frac{P(x)}{Q(x)} f(x) \right]^2 \\ &= \mathbb{E}_{x \sim P} \left[\frac{P(x)}{Q(x)} f(x)^2 \right] - \mathbb{E}_{x \sim P} [f(x)]^2 \end{aligned}$$

- ▶ If the importance sampling weight $P(x)/Q(x)$ is large
⇒ The variance of the estimator can blow up

- ▶ Consider a trajectory $\tau = (S_0, A_0, S_1, R_1, A_1, S_2, R_2 \dots, S_T, R_T)$
- ▶ Recall that the distribution of the trajectory is given by

$$p(\tau) = p(S_0) \prod_{t=0}^{T-1} \pi_{\theta}(A_t|S_t)p(S_{t+1}, R_{t+1}|S_t, A_t)$$

- ▶ Which depends on the policy π_{θ} with parametrization θ
- ▶ Further Recall that the policy gradient for finite horizon is given by

$$\nabla_{\theta} v(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(A_t|S_t) G_t \right]$$

- ▶ We want to reuse information from a different trajectory
 ⇒ Trajectory τ' generated from a policy $\pi_{\theta'}$
- ▶ Via importance sampling we can write

$$\nabla_{\theta} v(\theta) = \mathbb{E}_{\tau'} \left[\frac{p(\tau)}{p(\tau')} \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) G_t \right]$$

- ▶ By taking a closer look at the ratio $p(\tau)/p(\tau')$ we have that

$$\frac{p(\tau)}{p(\tau')} = \frac{p(S_0) \prod_{t=0}^{T-1} \pi_{\theta}(A_t | S_t) p(S_{t+1}, R_{t+1} | S_t, A_t)}{p(S_0) \prod_{t=0}^{T-1} \pi_{\theta'}(A_t | S_t) p(S_{t+1}, R_{t+1} | S_t, A_t)} = \prod_{t=0}^{T-1} \frac{\pi_{\theta}(A_t | S_t)}{\pi_{\theta'}(A_t | S_t)}$$

- ▶ **Small differences can multiply** and become quite large
 ⇒ Not the most stable way to reuse information from a different policy

Classical Policy Gradients

Monotonic Improvement Guarantees

Natural Policy Gradients and Trust Region Policy Optimization (TRPO)

Proximal Policy Optimization (PPO)

- ▶ First step to tackle the previous issues
 - ⇒ Characterize the **difference between two policies**
- ▶ Recall that the advantage is defined as the difference between q and v

$$a(S_t, A_t) = q(S_t, A_t) - v(S_t)$$

- ▶ It is a normalization with respect to the state
 - ⇒ How much an action can improve over the value of the current state
 - ⇒ Or the advantage of choosing a specific action
- ▶ Given policies π and π' and respective value functions $v(\pi)$ and $v(\pi')$

$$v(\pi') - v(\pi) = \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t a_{\pi}(S_t, A_t) \right]$$

- ▶ We relate the return of a policy in terms of the advantage over another
- ▶ We call this the **relative policy performance identity**

- Note that the advantage is given by

$$a_{\pi}(S_t, A_t) = R_{t+1} + \gamma v_{\pi}(S_{t+1}) - v_{\pi}(S_t)$$

- Then we can rewrite the term on the right hand side

$$\begin{aligned} \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t a_{\pi}(S_t, A_t) \right] &= \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t \left(R_{t+1} + \gamma v_{\pi}(S_{t+1}) - v_{\pi}(S_t) \right) \right] \\ &= v(\pi') + \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^{t+1} v_{\pi}(S_{t+1}) - \sum_{t=0}^{\infty} \gamma^t v_{\pi}(S_t) \right] \\ &= v(\pi') - \mathbb{E}_{\pi'} [v_{\pi}(S_0)] \\ &= v(\pi') - v(\pi) \end{aligned}$$

- ▶ We can use this identity to find a policy π' that maximizes $v(\pi')$

$$\max_{\pi'} v(\pi') = \max_{\pi'} v(\pi') - v(\pi)$$

- ▶ By the relative policy performance identity this is equivalent to

$$\max_{\pi'} v(\pi') = \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t a_{\pi}(S_t, A_t) \right]$$

- ▶ This allows us to **assess the quality of policy π' via the advantage of π**
- ▶ This expression still depends on trajectories generated from the policy π'
⇒ We want to **remove this dependence**

- ▶ We try to remove the dependence on trajectories sampled from π'
- ▶ Let us define the **discounted state distribution** d_π as

$$d_\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(S_t = s | \pi)$$

- ▶ Then we can rewrite the relative performance identity as

$$\begin{aligned} v(\pi') - v(\pi) &= \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t a_\pi(S_t, A_t) \right] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{S \sim d_{\pi'}, A \sim \pi'} [a_\pi(S, A)] \end{aligned}$$

- ▶ We can remove the dependence $a \sim \pi'$ via importance sampling

$$v(\pi') - v(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{S \sim d_{\pi'}, A \sim \pi} \left[\frac{\pi'(A|S)}{\pi(A|S)} a_\pi(S, A) \right]$$

- ▶ All that is left is to take care of the states $s \sim d_{\pi'}$

- ▶ How can we get rid of the dependency on $s \sim d_{\pi'}$?
- ▶ Let us assume for a second that $d_{\pi'} \approx d_{\pi}$
- ▶ We then have the local approximation $L_{\pi}(\pi')$ as follows

$$\begin{aligned}
 v(\pi') - v(\pi) &\approx \frac{1}{1-\gamma} \mathbb{E}_{S \sim d_{\pi}, A \sim \pi} \left[\frac{\pi'(A|S)}{\pi(A|S)} a_{\pi}(S, A) \right] \\
 &\triangleq L_{\pi}(\pi')
 \end{aligned}$$

- ▶ Under these conditions, we have the following **approximation guarantee**

$$|v(\pi') - (v(\pi) + L_{\pi}(\pi'))| \leq C \sqrt{\mathbb{E}_{S \sim d_{\pi}} [D_{\text{KL}}(\pi' || \pi)]}$$

- ▶ Where C is a system-dependent constant
- ▶ The approximation is good if the **policies are close in KL divergence**

- ▶ The local approximation is given by

$$L_{\pi}(\pi') = \frac{1}{1 - \gamma} \mathbb{E}_{S \sim d_{\pi}, A \sim \pi} \left[\frac{\pi'(A|S)}{\pi(A|S)} a_{\pi}(S, A) \right]$$

- ▶ We can rewrite it as

$$L_{\pi}(\pi') = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \frac{\pi'(A_t|S_t)}{\pi(A_t|S_t)} a_{\pi}(S_t, A_t) \right]$$

- ▶ We can use this approximation to **optimize with respect to π'**
 - ⇒ Using **only trajectories generated from π**
 - ⇒ Valid as long as policies are close in KL divergence
- ▶ Compared with the importance sampling of the policy gradient
 - ⇒ The ratio $\pi'(A_t|S_t)/\pi(A_t|S_t)$ **only appears per time instance**
 - ⇒ It is not multiplied over the whole trajectory

- ▶ We have another version of the relative performance bound

$$v(\pi') \geq L_{\pi}(\pi') - CD_{\text{KL}}^{\max}(\pi' \parallel \pi)$$

- ▶ Where $C = \frac{4\epsilon\gamma}{(1-\gamma)^2}$ and $D_{\text{KL}}^{\max}(\pi' \parallel \pi) = \max_s D_{\text{KL}}(\pi' \parallel \pi)$
- ▶ If we maximize the right hand side of this bound

$$\pi_{k+1} = \underset{\pi}{\operatorname{argmax}} \left[L_{\pi_k}(\pi) - CD_{\text{KL}}^{\max}(\pi \parallel \pi_k) \right]$$

- ▶ We are **guaranteed** to generate a **monotonically improving sequence**

$$v(\pi_0) \leq v(\pi_1) \leq v(\pi_2) \leq \dots$$

- ▶ To see this let $M_k(\pi) = L_{\pi_k}(\pi) - CD_{\text{KL}}^{\max}(\pi \parallel \pi_k)$ then

$$\begin{aligned} v(\pi_{k+1}) &\geq M_k(\pi_{k+1}) \\ v(\pi_k) &= M_k(\pi_k) \\ v(\pi_{k+1}) - v(\pi_k) &\geq M_k(\pi_{k+1}) - M_k(\pi_k) \end{aligned}$$

Input: Policy π_θ

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

 Generate an episode following $\pi_\theta : S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$

 Compute values of the advantage function $a_{\pi_k}(S_t, A_t)$

 Update policy according to

$$\pi_{k+1} = \operatorname{argmax}_{\pi} \left[L_{\pi_k}(\pi) - CD_{\text{KL}}^{\max}(\pi \| \pi_k) \right]$$

 Where $C = \frac{4\epsilon\gamma}{(1-\gamma)^2}$ and $D_{\text{KL}}^{\max}(\pi' \| \pi) = \max_S D_{\text{KL}}(\pi' \| \pi)$

end

Algorithm 1: Policy Iteration with Guaranteed Improvement

- ▶ The previous algorithm is a **Majorization-Minimization (MM) algorithm**
 - ⇒ In our case a Minorization-Maximization
 - ⇒ We construct a lower bound on the original objective function
 - ⇒ We maximize this lower bound
- ▶ This approach **guarantees monotonic improvement**
- ▶ However in practice computing $D_{\text{KL}}^{\max}(\pi' \parallel \pi)$ is complicated
 - ⇒ Requires evaluating the KL divergence at all states
 - ⇒ Computationally expensive or simply impossible
- ▶ Substitute it for the **average KL divergence** over states

$$v(\pi') \geq L_{\pi}(\pi') - C \mathbb{E}_{S \sim d_{\pi}} [D_{\text{KL}}(\pi' \parallel \pi)]$$

Classical Policy Gradients

Monotonic Improvement Guarantees

Natural Policy Gradients and Trust Region Policy Optimization (TRPO)

Proximal Policy Optimization (PPO)

- ▶ From now on we will consider parameterized policies π_θ
- ▶ From previously we have the following maximization

$$\theta_{k+1} = \operatorname{argmax}_\theta L_{\theta_k}(\theta) - C \mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)]$$

- ▶ The value of $C = \frac{4\epsilon\gamma}{(1-\gamma)^2}$ is given by the relative performance bound
 - ⇒ In practice this value is too conservative
- ▶ Instead we propose to solve the following optimization problem

$$\begin{aligned} & \underset{\theta}{\text{maximize}} && L_{\theta_k}(\theta) \\ & \text{subject to} && \mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)] \leq \delta \end{aligned}$$

- ▶ Where we have control over a parameter δ
 - ⇒ This is called a **trust region method**
 - ⇒ δ controls the extend of the region of trust

- ▶ The **natural policy gradient** is a special case of the previous problem
 - ⇒ Linearly approximate the objective function $L_{\theta_k}(\theta)$
 - ⇒ Quadratically approximate the constraint $\mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)]$
- ▶ Done via **Taylor's expansion** resulting in the following approximations

$$L_{\theta_k}(\theta) \approx L_{\theta_k}(\theta_k) + g^T(\theta - \theta_k)$$

$$\mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)] \approx \frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k)^T$$

- ▶ Where $g \triangleq \nabla_{\theta} L_{\theta_k}(\theta_k) |_{\theta=\theta_k}$ and $H \triangleq \nabla_{\theta}^2 \mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)] |_{\theta=\theta_k}$
- ▶ We can then write the **approximate optimization problem**

$$\begin{aligned} & \underset{\theta}{\text{maximize}} && L_{\theta_k}(\theta_k) + g^T(\theta - \theta_k) \\ & \text{subject to} && \frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k)^T \leq \delta \end{aligned}$$

- ▶ The gradient ascent update on which is given by

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

- ▶ The term $g \triangleq \nabla_{\theta} L_{\theta_k}(\theta_k) |_{\theta=\theta_k}$ corresponds to the **policy gradient**

$$\begin{aligned} \nabla_{\theta} L_{\theta_k}(\theta_k) |_{\theta=\theta_k} &= \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=0}^{\infty} \gamma^t \frac{\nabla_{\theta} \pi_{\theta}(A_t | S_t) |_{\theta=\theta_k}}{\pi_{\theta_k}(A_t | S_t)} a_{\pi_{\theta_k}}(S_t, A_t) \right] \\ &= \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) |_{\theta=\theta_k} a_{\pi_{\theta_k}}(S_t, A_t) \right] \end{aligned}$$

- ▶ The Hessian matrix $H \triangleq \nabla_{\theta}^2 \mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta || \theta_k)] |_{\theta=\theta_k}$ is given by

$$H = \mathbb{E}_{\pi_{\theta_k}} \left[\nabla_{\theta} \log \pi_{\theta}(A|S) |_{\theta=\theta_k} \nabla_{\theta} \log \pi_{\theta}(A|S) |_{\theta=\theta_k}^T \right]$$

- ▶ Called the **Fisher information matrix**
- ▶ Important property is that $H^{-1}g$ is **invariant to parametrization**

Input: Policy π_θ , KL divergence target δ

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

Generate an episode following $\pi_{\theta_k} : S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$

Estimate policy gradient g_k

$$g_k = \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) \Big|_{\theta=\theta_k} a_{\pi_{\theta_k}}(S_t, A_t) \right]$$

Estimate Hessian H_k

$$H_k = \mathbb{E}_{\pi_{\theta_k}} \left[\nabla_{\theta} \log \pi_{\theta}(A|S) \Big|_{\theta=\theta_k} \nabla_{\theta} \log \pi_{\theta}(A|S) \Big|_{\theta=\theta_k}^T \right]$$

Compute natural policy gradient update

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{g_k^T H_k^{-1} g_k}} H_k^{-1} g_k$$

end

Algorithm 2: Natural Policy Gradient

- ▶ The previous approach has a severe downfall
 - ⇒ We **need to compute the inverse H^{-1}**
 - ⇒ This can be problematic as it **does not scale well**
- ▶ Instead of finding H^{-1} use the **Conjugate Gradient** (CG) to compute $H^{-1}g$
- ▶ The conjugate gradient method can be used to solve for x in $Ax = b$
 - ⇒ Widely used iterative method to solve linear system of equations
 - ⇒ Finds it via projections to the Krylov subspace $\text{span}\{b, Ab, A^2b, \dots\}$
- ▶ Conjugate gradient **only needs to evaluate Hessian-vector products**

- ▶ We want to ensure that the KL divergence constraint is satisfied
 - ⇒ KL might not be satisfied due to the quadratic approximation
 - ⇒ At some iterations, the trust region given by δ can be too large
- ▶ We include a line search step
 - ⇒ Backtracking line search with exponential decay
 - ⇒ Enforce improvement in the approximation $L_{\theta_k}(\theta) \geq 0$
 - ⇒ Ensure the KL divergence constraint $\mathbb{E}_{S \sim d_{\pi_{\theta_k}}} [D_{\text{KL}}(\theta || \theta_k)] \leq \delta$ is met
- ▶ Trust Region Policy Optimization (TRPO) consists of
 - ⇒ The natural policy gradient
 - ⇒ Efficient Hessian-vector product computation via conjugate gradient
 - ⇒ A line search to enforce the KL divergence

Input: Policy π_θ , KL divergence target δ

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

Generate an episode following $\pi_{\theta_k} : S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$

Estimate policy gradient g_k

$$g_k = \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) \Big|_{\theta=\theta_k} a_{\pi_{\theta_k}}(S_t, A_t) \right]$$

Use Conjugate Gradient to estimate $x_k = H_k^{-1} g_k$

Conduct a line search with the proposed update

$$\Delta_k = \sqrt{\frac{2\delta}{g_k^T x_k}} x_k$$

Update parameters after line search

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

end

Algorithm 3: Trust Region Policy Optimization (TRPO)

Input: Update step $\Delta_k = \sqrt{\frac{2\delta}{\mathbf{g}_k^T \mathbf{H}_k^{-1} \mathbf{g}_k}} \mathbf{H}_k^{-1} \mathbf{g}_k$, step size decay $\alpha \in (0, 1)$

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

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

if $L_{\theta_k}(\theta) \geq 0$ **and** $\mathbb{E}_{S \sim d_{\pi_{\theta_k}}} [D_{\text{KL}}(\theta \| \theta_k)] \leq \delta$ **then**

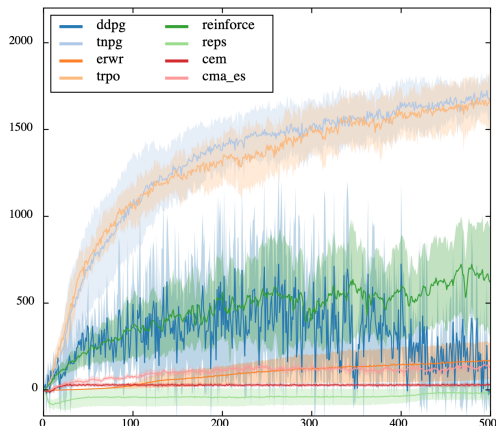
 Set $\theta_{k+1} = \theta_k + \alpha^j \Delta_k$

 Stop

end

end

Algorithm 4: Line Search for TRPO



► Comparison of various RL methods on a MuJoCo walker task.¹

¹Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel “Benchmarking Deep Reinforcement Learning for Continuous Control”, in *International Conference on Machine Learning*, pp. 1329-1338, 2016.

Classical Policy Gradients

Monotonic Improvement Guarantees

Natural Policy Gradients and Trust Region Policy Optimization (TRPO)

Proximal Policy Optimization (PPO)

- ▶ Methods based on **natural gradients are computationally expensive**
- ▶ Even with the modifications made to TRPO it is still expensive
 - ⇒ They are fundamentally second order methods
- ▶ Can we obtain similar performance with a first-order method?
 - ⇒ **Avoid Hessian-related computations**
- ▶ Proximal Policy Optimization (PPO) attempt to do so
 - ⇒ It is a first order approach based on heuristics
 - ⇒ Matches or surpasses TRPO performance in practice
- ▶ Proximal Policy Optimization has two variants
 - ⇒ Proximal Policy Optimization with KL penalty
 - ⇒ Proximal Policy Optimization with clipped objective

- ▶ Recall the previous unconstrained optimization problem

$$\theta_{k+1} = \operatorname{argmax}_{\theta} L_{\theta_k}(\theta) - \beta_k \mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)]$$

- ▶ Choosing a fixed value of β_k is complicated
- ▶ PPO with KL penalty attempts to **adjust the value of β_k dynamically**
 - ⇒ This is done via a heuristic check at each iteration
- ▶ The policy adapts the value β_k
 - ⇒ Individual iterations can violate KL constraints as β_k adapts
 - ⇒ This is not much of an issue in practice as adaptation occurs fast

Input: Policy π_θ , KL divergence target δ

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

Generate an episode following $\pi_\theta : S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$

Update policy according to

$$\theta_{k+1} = \operatorname{argmax}_\theta L_{\theta_k}(\theta) - \beta_k \mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta \| \theta_k)]$$

if $\mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta_{k+1} \| \theta_k)] \geq 1.5\delta$ **then**

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

end

if $\mathbb{E}_{S \sim d_{\theta_k}} [D_{\text{KL}}(\theta_{k+1} \| \theta_k)] \leq \delta/1.5$ **then**

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

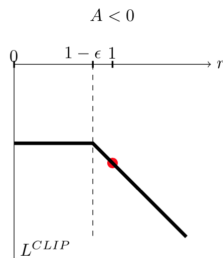
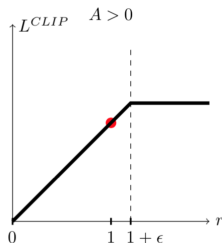
end

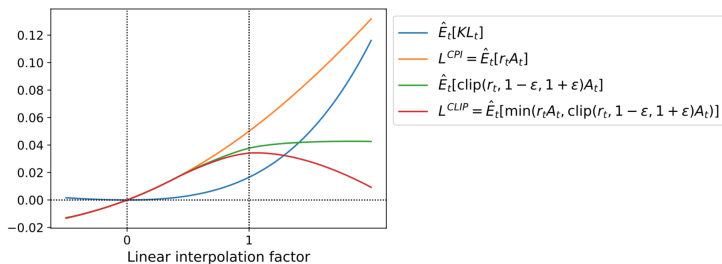
end

Algorithm 5: PPO with KL Penalty

- ▶ Clipping PPO modifies the objective to **penalize policies moving far away**
 ⇒ Extremely simple to implement
- ▶ The update is given by the maximization of the objective function

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \left[\mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=0}^{\infty} \min \left(\frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_k}(A_t|S_t)} a_{\pi_{\theta_k}}(S_t, A_t), \right. \right. \right. \\ \left. \left. \left. \operatorname{clip} \left(\frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_k}(A_t|S_t)}, 1 - \epsilon, 1 + \epsilon \right) a_{\pi_{\theta_k}}(S_t, A_t) \right) \right] \right]$$





- ▶ Objective function vs interpolation factor between θ_k and θ_{k+1}
- ▶ The objective is **penalized as the policy moves away from θ_k**

Input: Policy π_θ , Clipping value ϵ

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

 Generate an episode following $\pi_{\theta_k} : S_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$

 Update policy according to

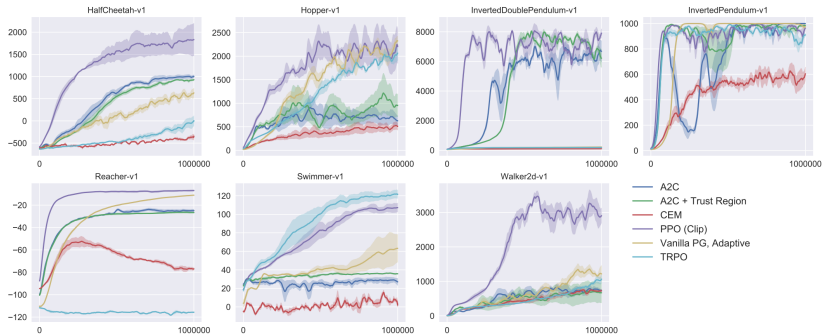
$$\theta_{k+1} = \underset{\theta}{\operatorname{argmax}} L_{\theta_k}^{\text{CLIP}}(\theta)$$

 Where $L_{\theta_k}^{\text{CLIP}}$ is given by

$$L_{\theta_k}^{\text{CLIP}} = \sum_{t=0}^{T-1} \min \left(\frac{\pi_\theta(A_t|S_t)}{\pi_{\theta_k}(A_t|S_t)} a_{\pi_{\theta_k}}(S_t, A_t), \operatorname{clip} \left(\frac{\pi_\theta(A_t|S_t)}{\pi_{\theta_k}(A_t|S_t)}, 1 - \epsilon, 1 + \epsilon \right) a_{\pi_{\theta_k}}(S_t, A_t) \right)$$

end

Algorithm 6: PPO with Clipped Objective



► Performance of PPO on various MuJoCo tasks.²

²J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. "Proximal policy optimization algorithms", *arXiv preprint arXiv:1707.06347*, 2017.

- ▶ These methods are based on the **KL divergence between policies**
 - ⇒ This is an asymmetric measure and hence not a metric
- ▶ **TRPO** provides a **fundamentally solid** approach
 - ⇒ However it is **computationally expensive**
- ▶ **PPO** is based on intuitive **heuristics**
 - ⇒ In practice **works surprisingly well**
- ▶ Still, there is an important question that we have not answered
 - ⇒ What is the right measure of similarity between two policies?
- ▶ We should use a true metric (Wasserstein distance)
 - ⇒ Wasserstein Reinforcement Learning