

Trust Region Policy Optimization

Tri Nguyen

Summer Reading 2021 - Oregon State University

September 3, 2021

Roadmap

- ▶ The method belong to policy gradient class, where policy is parametrized
- ▶ Use another objective function (avoid using the trick of Policy Gradient)
- ▶ Propose an optimization method to solve that objective function *approximately*
- ▶ Experimental result

The proposed method is pleasingly complicated :)))

Trust Region Policy Optimization

John Schulman
Sergey Levine
Philipp Moritz
Michael Jordan
Pieter Abbeel

JOSCHU@EECS.BERKELEY.EDU
SLEVINE@EECS.BERKELEY.EDU
PCMORITZ@EECS.BERKELEY.EDU
JORDAN@CS.BERKELEY.EDU
PABBEEL@CS.BERKELEY.EDU

University of California, Berkeley, Department of Electrical Engineering and Computer Sciences

- ▶ Kakade, Sham M. "A natural policy gradient." Advances in neural information processing systems 14 (2001). (910 citations)
- ▶ **Schulman, John, et al. "Trust region policy optimization." International conference on machine learning. PMLR, 2015. (3826 citations)**
- ▶ Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347. (5066 citations)

Notation

- ▶ $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the transition probability distribution.
- ▶ $r : \mathcal{S} \rightarrow \mathbb{R}$ is the reward function.
- ▶ $\rho_0 : \mathcal{S} \rightarrow \mathbb{R}$ is the distribution of the initial state s_0 .
- ▶ $\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} [\sum_{t=1}^{\infty} \gamma^t r(s_t)]$ is the expected discounted reward, where

$$s_0 \sim \rho_0(s_0), a_t \sim \pi(a_t | s_t), s_{t+1} \sim P(s_{t+1} | s_t, a_t).$$

- ▶ $Q_{\pi}(s_t, a_t) = \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} [\sum_{\ell=0}^{\infty} \gamma^{\ell} r(s_{t+1})]$
- ▶ $V_{\pi}(s_t) = \mathbb{E}_{a_t, s_{t+1}, \dots} [\sum_{\ell=0}^{\infty} \gamma^{\ell} r(s_{t+1})]$
- ▶ $A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(s)$ is the advantage function.

Starting Point

- ▶ Kakade & Langford (2002) showed a relation between any policy π and $\tilde{\pi}$.

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] \quad (1)$$

Proof.

Let $\tau \sim \tilde{\pi}$ be a trajectory sampled using $\tilde{\pi}$.

$$\begin{aligned} & \mathbb{E}_{\tau | \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] \\ &= \mathbb{E}_{\tau | \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t) + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t)) \right] \\ &= \mathbb{E}_{\tau | \tilde{\pi}} \left[r(s_0) + \gamma V_{\pi}(s_1) - V_{\pi}(s_0) + \gamma(r(s_1) + \gamma V_{\pi}(s_2) - V_{\pi}(s_1)) + \gamma^2(\cdot) + \dots \right] \\ &= \mathbb{E}_{\tau | \tilde{\pi}} \left[-V_{\pi}(s_0) + \sum_{t=0}^{\infty} \gamma^t r(s_t) \right] = -\mathbb{E}_{s_0} [V_{\pi}(s_0)] + \mathbb{E}_{\tau | \tilde{\pi}} \gamma^t r(s_t) \\ &= -\eta_{\pi} + \eta_{\tilde{\pi}} \end{aligned}$$

- ▶ Define ρ_π is the discounted visitation frequencies, $\rho_\pi = \sum_{t=0}^{\infty} \gamma^t P(s_t = s)$. Note that $s_0 \sim \rho_0$, the others depend on π and the environment.
- ▶ Rewrite (1)

$$\begin{aligned}
 \eta(\tilde{\pi}) &= \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] \\
 &= \eta(\pi) + \sum_{t=0}^{\infty} \sum_s P(s_t = s | \tilde{\pi}) \sum_a \tilde{\pi}(a|s) \gamma^t A_\pi(s_t, a_t) \\
 &= \eta(\pi) + \sum_s \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \tilde{\pi}) \sum_a \tilde{\pi}(a|s) A_\pi(s_t, a_t) \\
 &= \eta(\pi) + \sum_s \rho_{\tilde{\pi}(s)} \sum_a \tilde{\pi}(a|s) A_\pi(s_t, a_t)
 \end{aligned}$$

- ▶ If the blue term is nonnegative at every state, then $\tilde{\pi}$ is better or equal π
- ▶ In deterministic setting, it reduces to policy improvement, i.e., $\tilde{\pi}(s) = \arg \max_a A_\pi(s, a)$
- ▶ Maximizing the RHS respect to parameters of $\tilde{\pi}$ would result the best policy

The first approximation

- ▶ Recall

$$\eta(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s_t, a_t)$$

- ▶ Define

$$L_{\pi}(\tilde{\pi}) := \eta(\pi) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s_t, a_t)$$

then $L_{\pi}(\tilde{\pi}) \approx \eta(\tilde{\pi})$ locally in a sense that

$$L_{\pi_{\theta_0}}(\pi_{\theta_0}) = \eta(\pi_{\theta_0}), \quad \text{and} \quad \nabla_{\theta} L_{\pi_{\theta_0}}(\pi_{\theta}) |_{\theta=\theta_0} = \nabla_{\theta} \eta(\pi_{\theta}) |_{\theta=\theta_0}$$

- ▶ The first equality holds since $\underbrace{\sum_s \rho_{\pi}(s) \sum_a \pi(a|s)}_{\mathbb{E}} (Q_{\pi}(a, s) - V_{\pi}(s)) = 0$
- ▶ An improvement is guaranteed if using the following updating rule

$$\pi_{\text{new}}(a|s) = (1 - \alpha)\pi_{\text{old}}(a|s) + \alpha\pi'(a|s),$$

where $\pi' = \arg \max_{\pi} L_{\pi_{\text{old}}}(\pi)$ and it is bounded by

$$\eta(\pi_{\text{new}}) \geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{2\epsilon\gamma}{(1-\gamma)^2} \alpha^2$$

- ▶ Maximizing $L_{\pi_{\text{old}}}(\pi_{\theta})$ respect to θ is guaranteed to improve over π_{old}

New lower bound

- ▶ Define $D_{\text{TV}}(p||q) = \frac{1}{2} \sum_{i=1} |p_i - q_i|$ (called total variation divergence),
and
- ▶ $D_{\text{TV}}^{\max}(\pi, \tilde{\pi}) = \max_s D_{\text{TV}}(\pi(a|s), \tilde{\pi}(a|s))$

Theorem

Let $\alpha = D_{\text{TV}}^{\max}(\pi_{\text{old}}, \pi_{\text{new}})$. Then the following bound holds

$$\eta(\pi_{\text{new}}) \geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{4\epsilon\gamma}{(1-\gamma)^2} \alpha^2,$$

where $\epsilon = \max_{a,s} |A_{\pi}(s, a)|$

The improvement is guaranteed to general stochastic policy.

Algorithm

- ▶ Define $D_{\text{KL}}^{\max}(\pi, \tilde{\pi}) = \max_s D_{\text{KL}}(\pi(a|s), \tilde{\pi}(a|s))$
- ▶ Since $D_{\text{TV}}(p||q)^2 \leq D_{\text{KL}}(p||q)$

$$\eta(\pi_{\text{new}}) \geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{4\epsilon\gamma}{(1-\gamma)^2} D_{\text{KL}}^{\max}(\pi_{\text{old}}, \pi_{\text{new}}) \quad (2)$$

Algorithm 1 Policy iteration algorithm guaranteeing non-decreasing expected return η

Initialize π_0 .

for $i = 0, 1, 2, \dots$ until convergence **do**

 Compute all advantage values $A_{\pi_i}(s, a)$.

 Solve the constrained optimization problem

$$\pi_{i+1} = \arg \max_{\pi} [L_{\pi_i}(\pi) - CD_{\text{KL}}^{\max}(\pi_i, \pi)]$$

 where $C = 4\epsilon\gamma/(1-\gamma)^2$

$$\text{and } L_{\pi_i}(\pi) = \eta(\pi_i) + \sum_s \rho_{\pi_i}(s) \sum_a \pi(a|s) A_{\pi_i}(s, a)$$

end for

- ▶ The algorithm design is a type of minorization-minimization, where M_i is the surrogate function
- ▶ It is slow if C large
- ▶ Optimization problem:

$$\max_{\theta} L_{\theta_{\text{old}}}(\theta)$$

$$\text{subject to } D_{\text{KL}}^{\max}(\theta_{\text{old}}, \theta) \leq \delta$$

Policy improvement guarantee:

- ▶ Let $M_i(\pi) = L_{\pi_i}(\pi) - CD_{\text{KL}}^{\max}(\pi_i, \pi)$

$$\begin{aligned} \eta(\pi_{i+1}) &\geq M_i(\pi_{i+1}) \quad (\text{by 2}) \\ &\geq M_i(\pi_i) \quad (\text{by updating rule}) \\ &= \eta(\pi_i) \end{aligned}$$

The second approximation

Let $\bar{D}_{\text{KL}}^\rho = \mathbb{E}_{s \sim \rho} [D_{\text{KL}}](\pi_{\theta_1}(\cdot|s) || \pi_{\theta_2}(\cdot|s))$

- ▶ Relax constrain to $\bar{D}_{\text{KL}}^{\rho_{\text{old}}}(\theta_{\text{old}}, \theta) \leq \delta$
- ▶ Rewrite the objective function in expectation form

$$\begin{aligned} & \arg \max_{\theta} \sum_s \rho_{\theta_{\text{old}}}(s) \sum_a \pi_{\theta}(a|s) A_{\theta_{\text{old}}}(s, a) \\ &= \arg \max_{\theta} \left(\sum_s \rho_{\theta_{\text{old}}}(s) \sum_a \pi_{\theta}(a|s) Q_{\theta_{\text{old}}}(s, a) - \sum_s \rho_{\theta_{\text{old}}}(s) \sum_a \pi_{\theta}(a|s) V_{\theta_{\text{old}}}(s) \right) \\ &= \arg \max_{\theta} \sum_s \rho_{\theta_{\text{old}}}(s) \mathbb{E}_{a \sim q} \frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{\text{old}}} \\ &= \arg \max_{\theta} \mathbb{E}_{s \sim \rho_{\text{old}}, a \sim q} \frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{\text{old}}} \end{aligned}$$

- ▶ Final optimization problem

$$\begin{aligned} & \min_{\theta} \mathbb{E}_{s \sim \rho_{\text{old}}, a \sim q} \left[\frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{\text{old}}} \right] \\ & \text{subject to } \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}} [D_{\text{KL}}(\pi_{\text{old}}(\cdot|s) || \pi_{\text{new}}(\cdot|s))] \leq \delta \end{aligned}$$

where q is behaviour policy, $q = \pi_{\theta_{\text{old}}}$

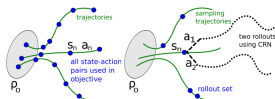
Practical algorithm - The third approximation

$$\min_{\theta} \mathbb{E}_{s \sim \rho_{\text{old}}, a \sim q} \left[\frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{\text{old}}} \right] \quad (3)$$

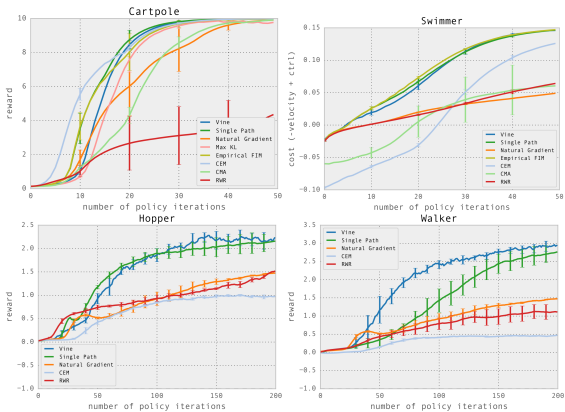
subject to $\mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}} [D_{\text{KL}}(\pi_{\text{old}}(\cdot|s) || \pi_{\text{new}}(\cdot|s))] \leq \delta$

Repeat the following steps:

- ▶ Use Monte Carlo simulation to collect trajectories.
 - ▶ Single path: generate 1 episode, then move to step 2
 - ▶ Vine: generate a number of trajectories, then choose a subset of states and samples actions from these state to generate new branching trajectories.
- ▶ Construct the estimated objective and constraint of Problem (3)
- ▶ Approximately solve this optimization problem using conjugate gradient algorithm followed by a line search.



Experiment result



Experiment result

	<i>B. Rider</i>	<i>Breakout</i>	<i>Enduro</i>	<i>Pong</i>	<i>Q*bert</i>	<i>Seaquest</i>	<i>S. Invaders</i>
Random	354	1.2	0	-20.4	157	110	179
Human	7456	31.0	368	-3.0	18900	28010	3690
Deep Q Learning	4092	168.0	470	20.0	1952	1705	581
UCC-1	5702	380	741	21	20025	2995	692
TRPO - single path	1425.2	10.8	534.6	20.9	1973.5	1908.6	568.4
TRPO - vine	859.5	34.2	430.8	20.9	7732.5	788.4	450.2

Table 1. Performance comparison for vision-based RL algorithms on the Atari domain. Our algorithms (bottom rows) were run once on each task, with the same architecture and parameters. Performance varies substantially from run to run (with different random initializations of the policy), but we could not obtain error statistics due to time constraints.

Conclusion

- ▶ Theorem 1 justifies for the surrogate objective function.
- ▶ Proposing using hard constraint instead of using penalty objective.
- ▶ Single path or vine using MC for estimating the minimization problem.
- ▶ Using conjugate gradient method for search direction, and line search to ensure the current step satisfies constraint.

Optimization method

$$\begin{aligned} & \max_{\theta} L_{\theta_{old}}(\theta) \\ & \text{subject to } \overline{D}_{\text{KL}}(\theta_{old}, \theta) \leq \delta \end{aligned}$$

Step 1: find the optimal direction to update

- ▶ The fisher information matrix is defined as

$$\mathbf{F} = \mathbb{E}_{x \sim p(x; \theta)} [\nabla \log p(x; \theta) \nabla \log p(x; \theta)^T]$$

- ▶ Fact: $\mathbf{F} = -\mathbf{H}_{\theta}[\log p(x; \theta)]$

- ▶ A derived fact: $\mathbf{F} = \mathbf{H}_{D_{\text{KL}}}$

- ▶ A derived approximation: $D_{\text{KL}}(p_{\theta} || p_{\theta'}) \approx \frac{1}{2}(\theta' - \theta)^T \mathbf{F}(\theta' - \theta)$

Then we can derive the optimal direction by

- ▶ Linear approximation of the objective:

$$L_{\theta_{old}}(\theta) \approx L_{\theta_{old}} + \nabla_{\theta} L_{\theta_{old}}(\theta) |_{\theta=\theta_{old}} [\theta - \theta_{old}]$$

- ▶ Quadratic approximation of the constrain:

$$D_{\text{KL}}(\theta_{old}, \theta) \approx \frac{\lambda}{2}(\theta - \theta_{old})^T \mathbf{F}(\theta - \theta_{old}),$$

where \mathbf{F} is Fisher information matrix.

- ▶ Lagrangian form

$$f(\theta) := L_{\theta_{old}} + \nabla_{\theta} L_{\theta_{old}}(\theta) |_{\theta=\theta_{old}} [\theta - \theta_{old}] + \frac{\lambda}{2}(\theta - \theta_{old})^T \mathbf{F}(\theta - \theta_{old})$$

Find optimal direction

- ▶ To find optimal f^* , we can find θ^* such as $\nabla f(\theta^*) = 0$,

$$0 = \nabla_{\theta} L_{\text{old}}(\theta) |_{\theta=\theta_{\text{old}}} + \lambda \mathbf{F}(\theta^* - \theta_{\text{old}}) \Leftrightarrow \theta^* = \theta_{\text{old}} - \lambda \mathbf{F}^{-1} \nabla_{\theta} L_{\theta_{\text{old}}}(\theta) |_{\theta=\theta_{\text{old}}}$$

- ▶ Therefore, the optimal direction is $\mathbf{y} = \mathbf{F}^{-1} \nabla_{\theta} L_{\theta_{\text{old}}}(\theta) |_{\theta=\theta_{\text{old}}}$
- ▶ Conjugate gradient can be used to solve $\mathbf{F}\mathbf{y} = \nabla_{\theta} L_{\theta_{\text{old}}}(\theta) |_{\theta=\theta_{\text{old}}}$

Connection to other methods

- ▶ Standard policy gradient

$$\begin{aligned} & \max_{\theta} [\nabla_{\theta} L_{\theta_{\text{old}}}(\theta) |_{\theta=\theta_{\text{old}}} \cdot (\theta - \theta_{\text{old}})] \\ & \text{subject to } \frac{1}{2} \|\theta - \theta_{\text{old}}\|^2 \leq \delta \end{aligned}$$

- ▶ Natural policy gradient (Kakade, 2002)

$$\begin{aligned} & \max_{\theta} [\nabla_{\theta} L_{\theta_{\text{old}}}(\theta) |_{\theta=\theta_{\text{old}}} \cdot (\theta - \theta_{\text{old}})] \\ & \text{subject to } \frac{1}{2} (\theta_{\text{old}} - \theta)^T \mathbf{F} (\theta_{\text{old}} - \theta) \leq \delta \end{aligned}$$

Proximal policy optimization

“Proximal policy optimization algorithms” (PPO) improved upon this by using only first-order derivative.

- ▶ Recall the objective in TRPO is

$$\mathbb{E} \left[\frac{\pi_{\theta}(\mathbf{a}|\mathbf{s})}{\pi_{\theta_{\text{old}}}(\mathbf{a}|\mathbf{s})} Q_{\theta_{\text{old}}} \right] = \mathbb{E}[r(\theta) Q_{\theta_{\text{old}}}]$$

- ▶ In PPO, the objective is

$$\mathbb{E} [\min(r(\theta) Q_{\theta_{\text{old}}}, \text{clip}(r(\theta), 1 - \epsilon, 1 + \epsilon) Q_{\theta_{\text{old}}})]$$

PPO experiment result

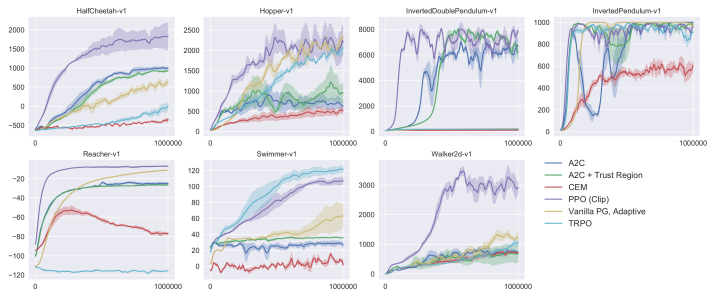


Figure 3: Comparison of several algorithms on several MuJoCo environments, training for one million timesteps.