

Trust Region Policy Optimization

June 11, 2016

TRPO is a gradient-based algorithm for policy optimization. It does not necessarily follow known patterns of reinforcement learning (i.e. policy iteration, policy gradient or actor-critic), but I hear it is wonderfully easy to use, and gives you great results.

1 Kakade & Langford's Policy Improvement

TRPO builds upon Kakade & Langford's result on policy improvement. This work starts on the following identity, which expresses the expected discount cost (or reward) η of a modified policy $\tilde{\pi}$ using expected costs of the original policy π .

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, s_1, a_1, \dots} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right]$$

where $A_{\pi}(s_t, a_t)$ is the usual advantage function, defined as $Q_{\pi}(s_t, a_t) - V(s_t, a_t)$. It should be noted that $a_t \sim \tilde{\pi}(a_t|s_t)$, which means the trajectory follows $\tilde{\pi}$, not π . I recognize this identity from when we proved the correctness (or convergence?) of policy iteration: if we can keep all advantage functions nonpositive, we can guarantee improvement.

There's a problem: with function approximators, you cannot guarantee $\sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$ will be non-positive for all states. We overcome this by forming a local approximation:

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

The difference here is that the state distribution function ρ uses π - which means we follow π , and only use $\tilde{\pi}$ for weighting advantage function values. Now, suppose our policies are parameterized by a parameter vector θ . Then $L_{\pi_{\theta_0}}(\pi_{\theta})$ is a good linear approximation of $\eta(\pi_{\theta})$ when θ is close to θ_0 - the value matches (obviously $L_{\pi_{\theta_0}}(\pi_{\theta_0}) = \eta(\pi_{\theta_0})$), and the gradients match as well.

Kakade & Langford now creates a new policy π_{new} that is close to π_{old} which minimizes $L_{\pi_{old}}(\pi_{new})$ which will hopefully minimize $\eta(\pi_{new})$. Unfortunately, they only had guarantees for π_{new} which is in the form of a mixed policy:

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

This is of course not desirable as policies are bloated over time and the changes we can make are pretty small (?).

2 TRPO

TRPO realizes that K&L's guarantee works when we plug in $D_{TV}^{\max}(\pi_{old}, \pi_{new})$ as α . D_{TV} is a distance measure between probability distributions, which is upper bounded by KL divergence. The paper proves:

$$\eta(\tilde{\pi}) \leq L_{\pi}(\tilde{\pi}) + C \cdot D_{KL}^{\max}(\pi, \tilde{\pi})$$

for "some" C (they have the formula for that, I don't find it important..). D_{KL}^{\max} denotes the maximum KL divergence of $\pi(\cdot|s)$ and $\tilde{\pi}(\cdot|s)$ among all states. Now, we can employ some majorization-minimization argument (reminds me of FISTA?) to prove that minimizing the RHS of above inequality improves η as well!

The actual TRPO algorithm is an approximation of minimizing the above quantity. The problems they were trying to solve was the step sizes were too small, and D_{KL}^{\max} is not always well defined. There are two approximations they used to deal with these:

1. Simply minimize $L_{\pi}(\tilde{\pi})$ with a constraint that D_{KL}^{\max} is lower than a certain threshold.
2. Use average KL divergence (weight each state by the distribution function ρ).

The second approximation allows us to estimate KL divergence from the samples we actually get! Tada!
Here is the outline of the TRPO algorithm.

- We sample (from a simulator or the real environment) one or more rollout sets using π . Now, pick some state-action pairs, and the rewards observed, and use constrained optimization to find the next generation policy that minimizes $L_{\pi}(\tilde{\pi})$ with the KL divergence constraint.
- When calculating $L_{\pi}(\tilde{\pi})$, use importance sampling (this comes up quite often!) to adjust differences between π and $\tilde{\pi}$. (Actually, we don't even have to rollout using π - we just need its advantage functions. We can take a different sampling distribution $q(a|s)$.. for discrete action spaces they say uniform q works better than taking π .)
- Replace advantage functions with action values (this only shifts the objective function by the value function, which ends up being a constant after expectations are taken of).

With these changes, at each iteration, we solve:

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

Now we can estimate all of these from samples (Q , take s from samples gained using $\pi_{\theta_{old}}$, etc).

2.1 Sampling Strategies

The paper introduces two strategies: single-path (simulate a single trajectory), vine (simulate multiple trajectories, sample random states along those trajectories, sample a number of short rollouts from those to estimate rewards). Of course, vine works better usually.

One interesting tidbit mentioned: using a common random number source across different short rollouts in vine reduces variance. I am not entirely sure what that means...

2.2 Solving the Optimization Problem

They recommend using conjugate gradient algorithm followed by a line search. They talk a little about calculating Hessian analytically, which I guess makes sense, except that I don't know how to do that with general policies... I should look into actual implementations..

I looked at performance graphs. They are good, obviously, and holy crap, CEM is so good for a gradient free method.