



Paper Reading

2018-10-27

TRPO

PPO

公式预警



Trust Region Policy Optimization(TRPO)---Review: MDP

MDP: $\langle S, A, P, r, \gamma, \rho_0 \rangle$

S : is the finite set of states

A : is the finite set of actions

$P : S \times A \times S \rightarrow R$: is the transition probability distribution

$r : S \rightarrow R$: is the reward function

$\rho_0 : S \rightarrow R$: is the distribution of the initial state s_0

$$Q_\pi(s_t, a_t) = E_{s_{t+1}, a_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right]$$

$$V_\pi(s_t) = E_{a_t, s_{t+1}, a_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right]$$

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(s_t)$$

$$\eta(\pi) = E_{s_0, a_0, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$



Trust Region Policy Optimization(TRPO)---Review: REINFORCE

$$\begin{aligned}\nabla \eta(\theta) &= E_{\pi}[\gamma^t Q_{\pi}(S_t, A_t) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)}] \\ &= E_{\pi}[\gamma^t G_t \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)}]\end{aligned} \quad \longrightarrow \quad \begin{aligned}\theta' &= \theta + \alpha \gamma^t G_t \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \\ \theta' &= \theta + \alpha \gamma^t G_t \nabla \log \pi(A_t | S_t, \theta)\end{aligned}$$

With Baseline:

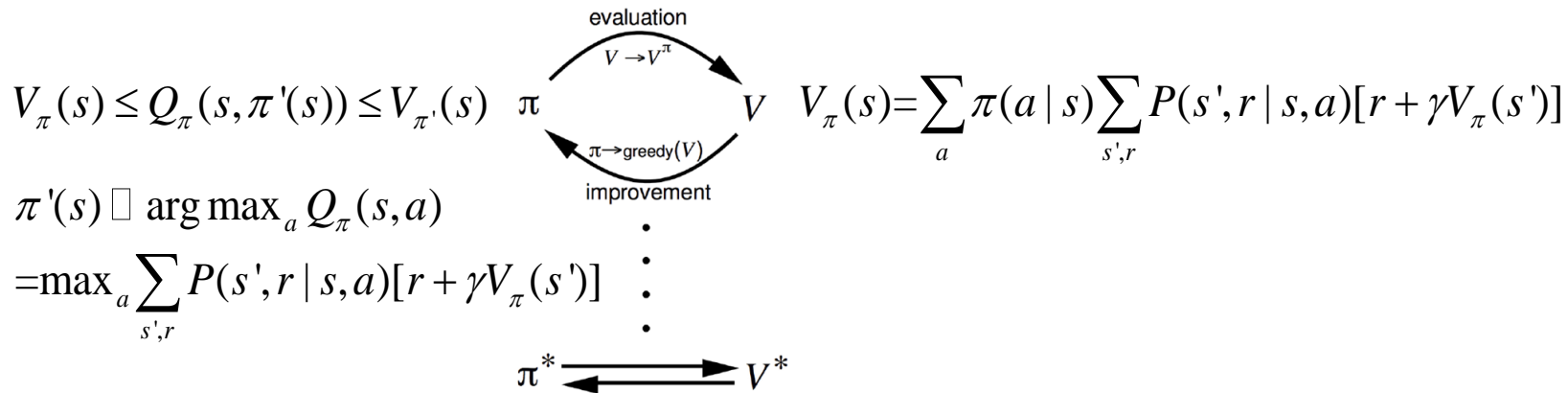
$$\begin{aligned}\sum_a b(s) \nabla_{\theta} \pi(a | s, \theta) &= 0 \\ \nabla \eta(\theta) &\square \sum_s d_{\pi}(s) \sum_a (Q_{\pi}(s, a) - b(s)) \nabla \pi(a | s, \theta) \\ \theta' &= \theta + \alpha \gamma^t (G_t - b(S_t)) \nabla \log \pi(A_t | S_t, \theta)\end{aligned}$$

$b(S_t) \square V(S_t)$ can reduce the variance.



Trust Region Policy Optimization(TRPO)---Review: Dynamic Programing

➤ Policy Iteration:



➤ Value Iteration:

$$V_*(s) \square \max_\pi V_\pi(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_*(s')]$$

$$\pi'(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_*(s')]$$

$$\begin{aligned}
 V_\pi(s) &\leq Q_\pi(s, \pi'(s)) \\
 &= E_{\pi'}[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s] \\
 &\leq E_{\pi'}[R_{t+1} + \gamma Q_\pi(S_{t+1}, \pi'(s)) | S_t = s] \\
 &= E_{\pi'}[R_{t+1} + \gamma E_{\pi'}[R_{t+2} + \gamma V_\pi(S_{t+2})] | S_t = s] \\
 &= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 V_\pi(S_{t+2}) | S_t = s] \\
 &\dots \\
 &\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\
 &= V_{\pi'}(s)
 \end{aligned}$$



Trust Region Policy Optimization(TRPO)

Policy Improvement :

$$\eta(\pi) \geq \eta(\pi)$$

So how to find policy π ?

$$\begin{aligned} \eta(\pi) &= \eta(\pi) + \sum_{t=0}^{\infty} \sum_s P(s_t = s | \pi) \sum_a \pi(a | s) \gamma^t A_{\pi}(s, a) \\ &= \eta(\pi) + \sum_s \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi) \sum_a \pi(a | s) A_{\pi}(s, a) \\ &= \eta(\pi) + \sum_s \rho_{\pi}(s) \sum_a \pi(a | s) A_{\pi}(s, a) \end{aligned}$$



Trust Region Policy Optimization(TRPO)

Policy Improvement :

$$\eta(\pi) \geq \eta(\pi)$$

Define:

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

We have:

$$L_{\pi_{\theta_0}}(\pi_{\theta_0}) = \eta(\pi_{\theta_0})$$
$$\nabla_{\theta} L_{\pi_{\theta_0}}(\pi_{\theta}) |_{\theta=\theta_0} = \nabla_{\theta} \eta(\pi_{\theta}) |_{\theta=\theta_0}$$

A sufficiently small step that improves $L_{\pi_{\theta_0}}(\pi_{\theta_0})$ will also improve η



Trust Region Policy Optimization (TRPO)

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

Notice:

$$\pi' = \arg \max_{\pi'} L_{\pi_{old}}(\pi')$$
$$\pi_{new}(a|s) = (1 - \alpha)\pi_{old}(a|s) + \alpha\pi'$$
$$\varepsilon = \max_s |E_{a \sim \pi'(a|s)}[A_\pi(s, a)]|$$
$$\eta(\pi_{new}) \geq L_{\pi_{old}}(\pi_{new}) - \frac{2\varepsilon\gamma}{(1-\gamma)^2} \alpha^2$$

Then:

$$D_{TV}(p \parallel q) = \frac{1}{2} \sum_i |p_i - q_i|$$
$$D_{TV}^{\max}(\pi, \pi) = \max_s D_{TV}(\pi(\square|s) \parallel \pi(\square|s))$$
$$\alpha = D_{TV}^{\max}(\pi_{old}, \pi_{new})$$
$$\varepsilon = \max_{s,a} |A_\pi(s, a)|$$
$$\eta(\pi_{new}) \geq L_{\pi_{old}}(\pi_{new}) - \frac{4\varepsilon\gamma}{(1-\gamma)^2} \alpha^2$$



Trust Region Policy Optimization (TRPO)

$$\eta(\pi) = \eta(\pi) + \sum_s \rho_\pi(s) \sum_a \pi(a|s) A_\pi(s, a)$$

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

Notice:

$$\pi' = \arg \max_{\pi'} L_{\pi_{old}}(\pi')$$

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

$$\varepsilon = \max_s |E_{a \sim \pi'(a|s)}[A_\pi(s, a)]|$$

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

Then:

$$D_{TV}(p \parallel q) = \frac{1}{2} \sum_i |p_i - q_i|$$

$$D_{TV}^{\max}(\pi, \pi) = \max_s D_{TV}(\pi(\square|s) \parallel \pi(\square|s))$$

$$\alpha = D_{TV}^{\max}(\pi_{old}, \pi_{new})$$

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

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



Trust Region Policy Optimization (TRPO)

$$D_{TV}(p \parallel q)^2 \leq D_{KL}(p \parallel q)$$

Define:

$$D_{KL}^{\max}(\pi, \pi) = \max_s D_{KL}(\pi(\square s) \parallel \pi(\square s))$$

$$\eta(\pi_{new}) \geq L_{\pi_{old}}(\pi_{new}) - \frac{4\varepsilon\gamma}{(1-\gamma)^2} \alpha^2 \quad \eta(\pi_{new}) \geq L_{\pi_{old}}(\pi_{new}) - CD_{KL}^{\max}(\pi_{old}, \pi_{new})$$

Define:

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

We have:

$$\eta(\pi_{i+1}) \geq M_i(\pi_{i+1})$$

$$\eta(\pi_i) = M_i(\pi_i)$$

$$\eta(\pi_{i+1}) - \eta(\pi_i) \geq M_i(\pi_{i+1}) - M_i(\pi_i)$$

So when:

$$\pi_{i+1} = \arg \max_{\pi} M_i(\pi)$$

$$\eta(\pi_{i+1}) - \eta(\pi_i) \geq M_i(\pi_{i+1}) - M_i(\pi_i) \geq 0$$



Trust Region Policy Optimization(TRPO)

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

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

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

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

end for



Trust Region Policy Optimization (TRPO)

$$\begin{aligned}
 \pi_{i+1} &= \arg \max_{\pi} M_i(\pi) \\
 &= \arg \max_{\pi} [L_{\pi_i}(\pi) - CD_{KL}^{\max}(\pi_i, \pi)] \\
 &= \arg \max_{\pi} [\eta(\pi_i) + \sum_s \rho_{\pi_i}(s) \sum_a \pi(a|s) A_{\pi_i}(s, a) - CD_{KL}^{\max}(\pi_i, \pi)] \\
 &= \arg \max_{\pi} [\sum_s \rho_{\pi_i}(s) \sum_a \pi(a|s) A_{\pi_i}(s, a) - CD_{KL}^{\max}(\pi_i, \pi)]
 \end{aligned}$$

Trust Region Constraint:

$$\begin{aligned}
 &\text{maximize}_{\pi} \sum_s \rho_{\pi_i}(s) \sum_a \pi(a|s) A_{\pi_i}(s, a) \\
 &\text{Subject to} \quad D_{KL}^{\max}(\pi_i, \pi) \leq \delta
 \end{aligned}$$

Approximation:

$$\begin{aligned}
 &\text{maximize}_{\theta} \sum_s \rho_{\theta_{old}}(s) \sum_a \pi_{\theta}(a|s) A_{\theta_{old}}(s, a) \\
 &\text{Subject to} \quad \bar{D}_{KL}^{\rho_{\theta_{old}}}(\theta_{old}, \theta) \leq \delta \\
 &\bar{D}_{KL}^{\rho}(\theta_1, \theta_2) \square E_{s \square \rho} [D_{KL}(\pi_{\theta_1}(\square s) \parallel \pi_{\theta_2}(\square s))]
 \end{aligned}$$



Trust Region Policy Optimization (TRPO)

$$\sum_a \pi_\theta(a | s) A_{\theta_{old}}(s, a) = E_{a \sim q} \left[\frac{\pi_\theta(a | s_n)}{q(a | s_n)} A_{\theta_{old}}(s_n, a) \right]$$

$$\text{maximize}_\theta E_{s \sim \rho_{\theta_{old}}, a \sim q} \left[\frac{\pi_\theta(a | s_n)}{q(a | s_n)} Q_{\theta_{old}}(s_n, a) \right]$$

$$\text{Subject to } E_{s \sim \rho_{\theta_{old}}} [D_{KL}(\pi_{\theta_{old}}(\cdot | s) \| \pi_\theta(\cdot | s))] \leq \delta$$

Look Back: **optimal step size, trust region**



Proximal Policy Optimization Algorithms(PPO)

Policy Gradient Methods:

Gradient Estimator:

$$g = E_t[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_t]$$

Objective:

$$L^{PG} = E_t[\log \pi_{\theta}(a_t | s_t) A_t]$$

Trust Region Methods:

$$\max_{\theta} E_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} A_t \right]$$

$$\text{Subject to } E_t [D_{KL}(\pi_{\theta_{old}}(\cdot | s) \| \pi_{\theta}(\cdot | s))] \leq \delta$$

Objective:

$$\max_{\theta} E_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} A_t \right] - \beta D_{KL}(\pi_{\theta_{old}}(\cdot | s) \| \pi_{\theta}(\cdot | s))$$



Proximal Policy Optimization Algorithms(PPO)

Define:
$$r_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$$

TRPO Objective:
$$L^{CPI} = E_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} A_t \right] = E_t [r_t(\theta) A_t]$$

PPO Objective:
$$L^{CLIP} = E_t [\min(r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \varepsilon, 1 + \varepsilon) A_t)]$$

谢谢!

撒花~

撒花~

