

# Optimality and Approximation with Policy Gradient Methods in Markov Decision Processes

Alekh Agarwal\* Sham M. Kakade† Jason D. Lee‡ Gaurav Mahajan§

## 【强化学习 98】PG Theory Summary



张楚琦

清华大学 交叉信息院博士在读

38 人赞同了该文章

这篇文章太经典了，今天在想某个问题的时候需要用到，由推导了一遍，做了个总结。

### 原文传送门

Agarwal, Alekh, et al. "Optimality and Approximation with Policy Gradient Methods in Markov Decision Processes." arXiv preprint arXiv:1908.00261 (2019).

### 正文

Algorithm	Error	Iterations
Direct parametrization w/ projected gradient ascent $\theta^{(t+1)} \leftarrow P_{\mathcal{A}(\theta)}(\theta^{(t)} + \eta \nabla_{\theta} V^{(t)}(\rho))$	$\min_{\theta \in \mathcal{A}} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon$ (Theorem 4.2)	$\frac{64\gamma S  A }{(1-\gamma)^3\epsilon^2} \left\  \frac{d_{\theta^*}^{\pi}}{\mu} \right\ _{\infty}^2$
Softmax parametrization (w/o regularization) $\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \nabla_{\theta} V^{(t)}(\rho)$	Only asymptotic convergence (Theorem 5.1)	May need exponential number of iterations
Softmax parametrization w/ relative entropy regularization $L_{\theta}(\theta) = V^{\pi_{\theta}}(\rho) + \frac{\lambda}{ S  A } \sum_{s,a} \log \pi_{\theta}(a s)$ $\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \nabla_{\theta} L_{\theta}(\theta^{(t)})$	$\min_{\theta \in \mathcal{A}} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon$ (Corollary 5.4)	$\frac{320 S ^2 A ^2}{(1-\gamma)^3\epsilon^2} \left\  \frac{d_{\theta^*}^{\pi}}{\mu} \right\ _{\infty}^2$
NPG with softmax parametrization $\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \nabla_{\theta} \left( \theta^{(t)} \right)^{-1} \nabla_{\theta} V^{(t)}(\rho)$ $\Leftrightarrow \theta_{s,a}^{(t+1)} \leftarrow \theta_{s,a}^{(t)} + \frac{\eta}{1-\gamma} A^{(t)}(s,a)$	$V^*(\rho) - V^{(t)}(\rho) \leq \epsilon$ (Theorem 5.7)	$\frac{2}{(1-\gamma)^2\epsilon}$
Unconstrained approximate NPG $\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \omega^{(t)}$ ( $\omega^{(t)}$ is sample based NPG) $\Leftrightarrow \omega^{(t)} \in \arg \min_{\omega} L_{\omega}(\theta; \theta)$ $L_{\omega}(\omega; \theta) = \mathbb{E}_{s,a \sim \pi_{\theta}} [A^{\pi_{\theta}}(s,a) - \omega \cdot \nabla_{\theta} \log \pi_{\theta}(a s)]^2$ $v(s,a) \leftarrow d_{\theta}^{\pi}(s) \pi_{\theta}(a s)$ $v(s,a) \leftarrow d_{\theta}^{\pi}(s) \pi_{\theta}(a s)$	$\min_{\theta \in \mathcal{A}} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon + \frac{\sqrt{\epsilon_{\text{approx}}}}{1-\gamma}$ where $L_{\omega}(\omega^{(t)}, \theta^{(t)}) \leq \epsilon_{\text{approx}}$ (Theorem 6.4)	$\frac{2\beta W^2 \log  A }{(1-\gamma)^2\epsilon^2}$ where $\beta$ smooth, $\ \omega^{(t)}\ _2 \leq W$
Unconstrained approximate NPG $v(s,a) \leftarrow d_{\theta}^{\pi}(s) \pi_{\theta}(a s)$ $\omega^{(t)}$ can achieve $\arg \min_{\omega} L_{\omega}(\omega; \theta)$	$\min_{\theta \in \mathcal{A}} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon + \frac{1}{\sqrt{(1-\gamma)^3}} \left\  \frac{v^*}{\mu} \right\ _{\infty} \epsilon_{\text{approx}}$ where $L_{\omega}^*(\theta^{(t)}) = \min_{\omega} L_{\omega}(\omega; \theta^{(t)}) \leq \epsilon_{\text{approx}}$ (Corollary 6.5)	$\frac{2\beta W^2 \log  A }{(1-\gamma)^2\epsilon^2}$
Unconstrained approximate NPG $v(s,a) \leftarrow d_{\theta}^{\pi}(s) \pi_{\theta}(a s)$ Consider estimation error for $\min_{\omega} L_{\omega}(\omega; \theta)$	$\min_{\theta \in \mathcal{A}} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon + \frac{1}{\sqrt{(1-\gamma)^3}} \left\  \frac{v^*}{\mu} \right\ _{\infty} (\sqrt{\epsilon_{\text{approx}}} + \frac{\epsilon_1}{\sqrt{N}})$ (Corollary 6.5)	$\frac{2\beta W^2 \log  A }{(1-\gamma)^2\epsilon^2}$ where $\ \omega^{(t)}\ _2 \leq W$
Projected Policy Gradient for constrained policy class $\theta^{(t+1)} \leftarrow P_{\mathcal{A}}(\theta^{(t)} + \eta \nabla_{\theta} V^{(t)}(\rho))$	$\min_{\theta \in \mathcal{A}} \{V^*(\rho) - V^{(t)}(\rho)\} \leq (W^* + 1)\epsilon + \frac{1}{\sqrt{(1-\gamma)^3}} \left\  \frac{d_{\theta^*}^{\pi}}{\mu} \right\ _{\infty} \epsilon_{\text{approx}}$ where $L_{\theta^*}^*(\theta^{(t)}) \leq \epsilon_{\text{approx}}$ and $\ \omega^{(t)}\ _2 \leq W^*$ (Corollary 6.14)	$\frac{8\beta}{(1-\gamma)^3\epsilon^2} \left\  \frac{d_{\theta^*}^{\pi}}{\mu} \right\ _{\infty}$

知乎 @张楚琦

## Direct parametrization

Direct parametrization with projected gradient ascent

$$\pi^{(t+1)} = P_{\Delta(A)^{|S|}}(\pi^{(t)} + \eta \nabla_{\pi} V^{(t)}(\mu))$$

**Theorem 4.2.** The projected gradient ascent algorithm (9) on  $V^{\pi}(\mu)$  with stepsize  $\eta = \frac{(1-\gamma)^3}{2\gamma|A|}$  satisfies for all distributions  $\rho \in \Delta(S)$ ,

$$\min_{t \leq T} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon \quad \text{whenever} \quad T > \frac{64\gamma|S||A|}{(1-\gamma)^6\epsilon^2} \left\| \frac{d_{\rho}^{\pi^*}}{\mu} \right\|_{\infty}^2.$$

Theorem 4.2 proof sketch

1. Smoothness of  $V^{\pi}(\mu)$  with  $\beta = \frac{2\gamma|A|}{(1-\gamma)^3}$  (Lemma E3)
2. #samples to magnitude of the update (projected gradient ascent, Bech, 2017)
3. Small magnitude of the update to small gradient (projection, Ghadimi and Lan, 2016)
4. Small gradient to optimality (Lemma 4.1)

## Softmax parametrization w/ relative entropy regularization

$$L_\lambda(\theta) = V^{\pi_\theta}(\mu) + \frac{\lambda}{|S||A|} \sum_{s,a} \log \pi_\theta(a|s), \quad \theta^{(t+1)} = \theta^{(t)} + \eta \nabla_\theta L_\lambda(\theta^{(t)})$$

**Corollary 5.4.** (Iteration complexity with relative entropy regularization) Let  $\beta_\lambda := \frac{8\gamma}{(1-\gamma)^3} + \frac{2\lambda}{|S|}$ . Starting from any initial  $\theta^{(0)}$ , consider the updates (13) with  $\lambda = \frac{\epsilon(1-\gamma)}{2 \left\| \frac{d_\mu^{\pi^*}}{\mu} \right\|_\infty}$  and  $\eta = 1/\beta_\lambda$ . Then for all starting state distributions  $\rho$ , we have

$$\min_{t \leq T} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon \quad \text{whenever} \quad T \geq \frac{320|S|^2|A|^2}{(1-\gamma)^6 \epsilon^2} \left\| \frac{d_\mu^{\pi^*}}{\mu} \right\|_\infty^2.$$

Corollary 5.4 proof sketch

1. Smoothness of  $L_\lambda(\theta)$  with  $\beta_\lambda = \frac{8\gamma}{(1-\gamma)^3} + \frac{2\lambda}{|S|}$  (Lemma E4)
2. #samples to small gradient (unconstrained gradient ascent, Ghadimi and Lan, 2013)
3. Small gradient to optimality (Theorem 5.3)

知乎 @张楚珩

## Natural policy gradient with softmax parametrization and exact gradient

$$F_\rho(\theta) = \mathbb{E}_{s \sim \rho} \mathbb{E}_{a \sim \pi_\theta} \mathbb{E}_{a' \sim \pi_\theta} [(\nabla_\theta \log \pi_\theta(a|s)) (\nabla_\theta \log \pi_\theta(a'|s))^T] \quad \theta^{(t+1)} = \theta^{(t)} + \eta F_\rho(\theta^{(t)})^{-1} \nabla_\theta V^{(t)}(\rho)$$

**Theorem 5.7** (Global convergence for Natural Policy Gradient Ascent). Suppose we run the NPG updates (15) using  $\rho \in \Delta(S)$  and with  $\theta^{(0)} = 0$ . Fix  $\eta > 0$ . For all  $T > 0$ , we have:

$$V^{(T)}(\rho) \geq V^*(\rho) - \frac{\log |A|}{\eta T} - \frac{1}{(1-\gamma)^2 T}.$$

In particular, setting  $\eta \geq (1-\gamma)^2 \log |A|$ , we see that NPG finds an  $\epsilon$ -optimal policy in a number of iterations that is at most:

$$T \leq \frac{2}{(1-\gamma)^2 \epsilon}.$$

Theorem 5.7 proof sketch

1. NPG update has a simple form under softmax parametrization (Lemma 5.6)
2. Improvement lower bound exists in terms of the partition function (Lemma 5.8)

NPG gives large weights to rare situations, and therefore the algorithm does not need  $\mu$  and the sample complexity does not involve distribution mismatch coefficient using exact gradients.

知乎 @张楚珩

## Unconstrained approximate NPG

$$\theta^{(t+1)} = \theta^{(t)} + \eta \omega^{(t)} \quad (\omega^{(t)} \text{ is sample based NPG}) \Leftrightarrow \omega^{(t)} \in \arg \min_{\omega} L_{\nu}(\omega; \theta)$$

**Theorem 6.4.** (NPG approximation) Fix a comparison policy  $\pi$  and a state distribution  $\rho$ . Define  $\nu$  as the induced state-action measure under  $\pi$ , i.e.

$$\nu(s, a) = d_{\pi}^{\pi}(s) \pi(a|s).$$

Suppose that the update rule (18) starts with  $\theta^{(0)} = 0$  and uses the (arbitrary) sequence of weights  $w^{(0)}, \dots, w^{(T)}$ ; that Assumption 6.2 holds; and that for all  $t < T$ ,

$$\frac{1}{T} \sum_{t=0}^{T-1} L_{\nu}(w^{(t)}; \theta^{(t)}) \leq \bar{\epsilon}_{\text{approx}}, \quad \|w^{(t)}\|_2 \leq W.$$

We have that:

$$\min_{t \leq T} \{V^{\pi}(\rho) - V^{(t)}(\rho)\} \leq \frac{1}{1-\gamma} \left( \sqrt{\epsilon_{\text{approx}}} + \frac{\log |A|}{\eta T} + \frac{\eta \beta W^2}{2} \right).$$

Theorem 6.4 proof sketch

1. NPG update is equivalent to minimize  $\arg \min_{\omega} L_{\nu}(\omega; \theta)$  (Kakade, 2001)
2. Assume that the parametrized policies are  $\beta$ -smooth (Assumption 6.2)
3. Performance difference lemma (Lemma 3.2) and Jensen's inequality

知乎 @张楚琦

发布于 2019-10-31

强化学习 (Reinforcement Learning)

深度学习 (Deep Learning)

赞同 38



1 条评论

分享

喜欢

收藏



文章被以下专栏收录



强化学习前沿  
读呀读paper

进入专栏