

Local Optimization Achieves Global Optimality in Multi-Agent Reinforcement Learning

Yulai Zhao

Zhuoran Yang

Zhaoran Wang

Jason Lee

Princeton University

Yale University

Northwestern University

Princeton University

ABSTRACT

Policy optimization methods with function approximation are widely used in multi-agent reinforcement learning. However, it remains elusive how to design such algorithms with statistical guarantees. Leveraging a multi-agent performance difference lemma that characterizes the landscape of multi-agent policy optimization, we find that the localized action value function serves as an ideal descent direction for each local policy. Motivated by the observation, we present a multi-agent PPO algorithm in which the local policy of each agent is updated similarly to vanilla PPO. We prove that with standard regularity conditions on the Markov game and problem-dependent quantities, our algorithm converges to the globally optimal policy at a sublinear rate. We extend our algorithm to the off-policy setting and introduce pessimism to policy evaluation, which aligns with experiments. To our knowledge, this is the first provably convergent multi-agent PPO algorithm in cooperative Markov games.

MOTIVATION

Multi-agent reinforcement learning (MARL) has demonstrated many empirical successes, e.g. strategic games (Go, StarCraft II..)

• Main Challenges in MARL (Zhang 2021)

1. non-stationarity: each action taken by one agent affects the total reward and the transition of state.
2. scalability: taking other agents into consideration, each individual agent would face the joint action space, whose dimension increases exponentially with the number of agents
3. function approximation: closely related to the scalability issue, the state space and joint action space are often immense in MARL

Despite the empirical successes, theoretical studies of policy optimization in MARL are very limited. Even for the cooperative setting where the agents share a common goal: maximizing the total reward function. In this paper, we aim to answer the following fundamental question:

Can we design a provably convergent multi-agent policy optimization algorithm in the cooperative setting with function approximation?

We answer this question affirmatively!

PROBLEM SETUP

• Fully-cooperative Markov Games

- a tuple $M = (N, \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$: A party of participants N , a set of states \mathcal{S} , a set of actions \mathcal{A} , a transition probability $\mathcal{P}: \mathcal{S} \times \mathcal{A} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, a reward function: $r: \mathcal{S} \times \mathcal{A} \times \mathcal{A} \rightarrow [0, 1]$, a discounted factor $\gamma \in [0, 1]$.

- define policies as probability distributions over action space: $\pi \in \mathcal{S} \rightarrow \Delta(\mathcal{A})$. We shall use bold π to represent the joint policy.

Value function: $V^\pi(s) = E_{a \sim \pi} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t) | s_0 = s]$.

• Multi-agent Notations

We write index k on superscript when we refer to the specific k -th agent. When bold symbols are used without any superscript (e.g., \mathbf{a}), they consider all agents. For simplicity, let $(m; m')$ be shorthand for set: $\{i | m \leq i \leq m', i \in N\}$.

Definition 3.1. Let P be a subset in N . The multi-agent action value function associated with agents in P is: $Q_P^k(s, \mathbf{a}^P) = E_{\tilde{a} \sim \tilde{\pi}} [Q_\pi(s, \mathbf{a}^P, \tilde{\mathbf{a}})]$, we use a tilde over symbols to refer to the complement agents, namely $\tilde{\mathbf{a}} = \{\mathbf{a}^i | i \notin P, i \in N\}$.

Let $P, P' \subset N$ be two disjoint subsets of agents. The multiagent advantage function is defined below. Essentially, it accounts for the improvements of setting agents $\mathbf{a}^{P'}$ upon setting agents \mathbf{a}^P , while all other agents follow π : $A_\pi^P(s, \mathbf{a}^P, \mathbf{a}^{P'}) = Q_{\pi^{P \cup P'}}(s, \mathbf{a}^P, \mathbf{a}^{P'}) - Q_\pi^P(s, \mathbf{a}^P)$.

RESULTS

• Multi-agent PPO for online setting

For this setting, we adopt a useful log-linear parametrization to build conditional dependency structure. Such a conditional structure enables us to implement local optimization for every agent while still aiming at global optimality, because action value functions are ideal descent directions for local policies, a result from the multi-agent advantage function decomposition (Kuba 2022).

At the k -th iteration, we have the current policy π_{θ_k} , and we need to:

- (1) perform **policy evaluation** to obtain the action value function estimates $\hat{Q}_{\pi_{\theta_k}}$ for determining the quality of π_{θ_k} .
- (2) perform **policy improvement** to update policy to $\pi_{\theta_{k+1}}$. In practice, we approximate such policy within a certain parametrization class, which is achieved by implementing SGD updates.

Parametrization For the m -th agent ($m \in \mathcal{A}$), its conditional policy depends on all prior ordered agents $\mathbf{a}^{1:m-1}$. Given a coefficient vector $\theta^m \in \Theta$, where $\Theta = \{\theta | \|\theta\|_2 \leq R\}$ is a convex, norm-constrained set. The probability of choosing action \mathbf{a}^m under state s is

$$\pi_m(\mathbf{a}^m | s, \mathbf{a}^{1:m-1}) = \frac{\exp(\phi^\top(s, \mathbf{a}^{1:m-1}, \mathbf{a}^m) \theta^m)}{\sum_{\mathbf{a}^m \in \mathcal{A}} \exp(\phi^\top(s, \mathbf{a}^{1:m-1}, \mathbf{a}^m) \theta^m)} \quad (2)$$

where ϕ is a set of feature vector representations. Without loss of generality, we impose a regularity condition such that every $\|\phi\|_2 \leq 1$. This parametrization has been widely used in RL literature [Bravanan et al., 2009, Gimpel and Smith, 2010, Hoss et al., 2013, Agarwal et al., 2020, Zhao et al., 2022].

Algorithm 3 Policy Improvement Solver for MA-PPO

Input: MG $(\mathcal{N}, \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, iterations T , stepsize η , samples $\{s_t, \mathbf{a}_t^{1:m-1}, \mathbf{a}_t^m\}_{t=1}^T$.

Output: Policy update θ .

- 1: Initialize $\theta_0 = 0$.
- 2: **for** $t = 0, 1, \dots, T-1$ **do**
- 3: Let $(s_t, \mathbf{a}_t^{1:m-1}, \mathbf{a}_t^m) \leftarrow (s_t, \mathbf{a}_t^{1:m-1}, \mathbf{a}_t^m)$.
- 4: $\theta(t + \frac{1}{2}) \leftarrow \theta(t) - 2\eta\phi(s_t, \mathbf{a}_t^{1:m-1}, \mathbf{a}_t^m) \left((\theta(t) - \theta^m)^\top \phi(s_t, \mathbf{a}_t^{1:m-1}, \mathbf{a}_t^m) - \beta_t^{-1} Q_{\pi_{\theta(t)}}^m(s_t, \mathbf{a}_t^{1:m-1}, \mathbf{a}_t^m) \right)$
- 5: $\theta(t+1) \leftarrow \Pi_{\Theta} \theta(t + \frac{1}{2})$
- 6: **end for**
- 7: Calculate average: $\theta \leftarrow \frac{1}{T} \sum_{t=1}^T \theta_t$.

RESULTS (CONTINUE)

Theorem 1 (informal): For this setting, after K iterations, we have $J(\pi^*) - J(\bar{\pi})$ upper bounded by $O\left(\frac{N}{1-\gamma} \sqrt{\frac{\log(LA)}{K}}\right)$.

• Pessimistic MA-PPO with Linear Function Approximation

We perform pessimistic policy evaluation via regularization to reduce such overestimation aligning with experimental works.

Theorem 2 (informal): For this setting, after K iterations, we have $J(\pi^*) - J(\bar{\pi})$ upper bounded by $O\left(\frac{N}{1-\gamma} \sqrt{\frac{\log(LA)}{K}} + \frac{c}{(1-\gamma)} \sqrt{\frac{\log(LA)}{K}}\right)$.

REFERENCES

- Kaifu Zhang, Zhuoran Yang, and Tamer Basar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. Handbook of Reinforcement Learning and Control, pp. 321–384, 2021.
- Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong Yang. Trust region policy optimisation in multiagent reinforcement learning. In International Conference on Learning Representations, 2022.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. Optimality and approximation with policy gradient methods in Markov decision processes. In Conference on Learning Theory, pages 64–66. PMLR, 2020.

CONCLUSIONS

1. We propose a multi-agent PPO algorithm in which the local policy of each agent is updated sequentially in a similar fashion as vanilla PPO algorithm (Schulman et al., 2017).
2. We adopt the log-linear function approximation for the policies. We prove that multi-agent PPO converges to a sublinear $O\left(\frac{N}{1-\gamma} \sqrt{\frac{\log(LA)}{K}}\right)$ rate up to some statistical errors incurred in evaluating/improving policies.
3. Moreover, we propose an off-policy variant of the multi-agent PPO algorithm and introduce pessimism into policy evaluation.

ACKNOWLEDGEMENTS

JDL acknowledges support of the ARO under MURI Award W911NF-11-1-0304, the Sloan Research Fellowship, NSF CCF 2002272, NSF IIS 2107304, NSF CIF 2212262, ONR Young Investigator Award, and NSF CAREER Award 2144994