

Trust Region Policy Optimisation in Multi-Agent Reinforcement Learning

*Jakub Grudzien Kuba, Ruiqing Chen,
Muning Wen, Ying Wen, Fanglei Sung, Jun Wang, Yaodong Yang*



Cooperative MARL: Problem Formulation

Problem Formulation

At time step t , n agents are at state s_t



state s_t

Problem Formulation

The agents take actions $\mathbf{a}_t^1 \sim \pi^1(\cdot^1 | s_t), \dots, \mathbf{a}_t^n \sim \pi^n(\cdot^n | s_t)$

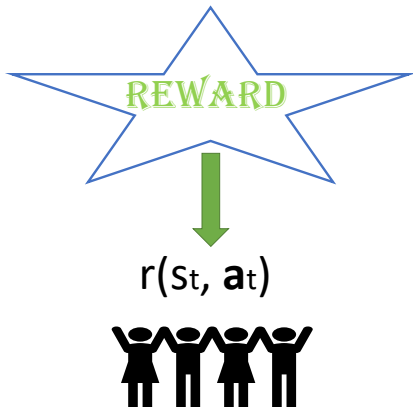


state s_t

Equivalently $\mathbf{a}_t \sim \pi(\cdot | s_t)$

Problem Formulation

The environment emits the joint reward $r(s_t, \mathbf{a}_t)$



Problem Formulation

The agents move to the next state

$$s_{t+1} \sim P(\cdot | s_t, \mathbf{a}_t)$$



state s_{t+1}

Problem Formulation

The agents want to maximise the joint return

$$J(\boldsymbol{\pi}) = \mathbb{E}_{s_0 \sim d^0, a_{0:\infty} \sim \boldsymbol{\pi}, s_{1:\infty} \sim P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, \mathbf{a}_t) \right]$$

Multi-Agent Trust Region Learning: A Summary of Existing Approaches

Existing Approaches

Methodology

Current Approaches

Methodology

- ▶ To endow all agents with a single policy (**homogeneity**),

$$\pi^i = \pi^j = \pi, \quad \forall i, j \in \mathcal{N}.$$

Existing Approaches

Methodology

- ▶ To endow all agents with a single policy (**homogeneity**),

$$\pi^i = \pi^j = \pi, \quad \forall i, j \in \mathcal{N}.$$

- ▶ To perform a trust-region update on π , for example, PPO

$$\mathbb{E}_{\mathbf{s} \sim \rho_\pi, \mathbf{a} \sim \pi} \left[\min \left(\frac{\bar{\pi}(\mathbf{a}|\mathbf{s})}{\pi(\mathbf{a}|\mathbf{s})} A_\pi(\mathbf{s}, \mathbf{a}), \text{clip} \left(\frac{\bar{\pi}(\mathbf{a}|\mathbf{s})}{\pi(\mathbf{a}|\mathbf{s})}, 1 \pm \epsilon \right) A_\pi(\mathbf{s}, \mathbf{a}) \right) \right].$$

Existing Approaches

Shortcomings

- ▶ Suboptimality (possibly exponential) of the solution in case of homogeneous policies

$$\frac{J_{\text{share}}^*}{J^*} = \frac{2}{2^n}.$$

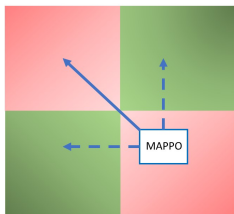
Current Approaches

Shortcomings

- ▶ Suboptimality (possibly exponential) of the solution in case of homogeneous policies

$$\frac{J_{\text{share}}^*}{J^*} = \frac{2}{2^n}.$$

- ▶ A loss of the monotonic improvement property for **heterogeneous** policies.



Heterogenous-Agent Trust Region Algorithms

Heterogeneous-Agent Trust Region Algorithms

Let $i_{1:m}$ and $j_{1:k}$ be disjoint, ordered subsets of agents.

Heterogeneous-Agent Trust Region Algorithms

Let $i_{1:m}$ and $j_{1:k}$ be disjoint, ordered subsets of agents.

- ▶ *The multi-agent state-action value function*

$$Q_{\pi}^{i_{1:m}}(s, \mathbf{a}^{i_{1:m}}) \triangleq \mathbb{E}_{\mathbf{a}^{-i_{1:m}} \sim \pi^{-i_{1:m}}} [Q_{\pi}(s, \mathbf{a}^{i_{1:m}}, \mathbf{a}^{-i_{1:m}})],$$

Heterogeneous-Agent Trust Region Algorithms

Let $i_{1:m}$ and $j_{1:k}$ be disjoint, ordered subsets of agents.

- ▶ *The multi-agent state-action value function*

$$Q_{\pi}^{i_{1:m}}(s, \mathbf{a}^{i_{1:m}}) \triangleq \mathbb{E}_{\mathbf{a}^{-i_{1:m}} \sim \pi^{-i_{1:m}}} [Q_{\pi}(s, \mathbf{a}^{i_{1:m}}, \mathbf{a}^{-i_{1:m}})],$$

- ▶ *The multi-agent advantage function*

$$A_{\pi}^{i_{1:m}}(s, \mathbf{a}^{j_{1:k}}, \mathbf{a}^{i_{1:m}}) \triangleq Q_{\pi}^{j_{1:k}, i_{1:m}}(s, \mathbf{a}^{j_{1:k}}, \mathbf{a}^{i_{1:m}}) - Q_{\pi}^{j_{1:k}}(s, \mathbf{a}^{j_{1:k}}).$$

Heterogeneous-Agent Trust Region Algorithms

Let $i_{1:m}$ and $j_{1:k}$ be disjoint, ordered subsets of agents.

- ▶ *The multi-agent state-action value function*

$$Q_{\pi}^{i_{1:m}}(s, \mathbf{a}^{i_{1:m}}) \triangleq \mathbb{E}_{\mathbf{a}^{-i_{1:m}} \sim \pi^{-i_{1:m}}} [Q_{\pi}(s, \mathbf{a}^{i_{1:m}}, \mathbf{a}^{-i_{1:m}})],$$

- ▶ *The multi-agent advantage function*

$$A_{\pi}^{i_{1:m}}(s, \mathbf{a}^{j_{1:k}}, \mathbf{a}^{i_{1:m}}) \triangleq Q_{\pi}^{j_{1:k}, i_{1:m}}(s, \mathbf{a}^{j_{1:k}}, \mathbf{a}^{i_{1:m}}) - Q_{\pi}^{j_{1:k}}(s, \mathbf{a}^{j_{1:k}}).$$

Lemma. For any state $s \in \mathcal{S}$, joint action $\mathbf{a}^{i_{1:n}} \in \mathcal{A}$, and permutation of agents $i_{1:n}$,

$$A_{\pi}^{i_{1:n}}(s, \mathbf{a}^{i_{1:n}}) = \sum_{m=1}^n A_{\pi}^{i_m}(s, \mathbf{a}^{i_{1:m-1}}, \mathbf{a}^{i_m}).$$

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.
- ▶ **for** $k = 0, 1, 2, \dots$

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.
- ▶ **for** $k = 0, 1, 2, \dots$
- ▶ Compute $C = \frac{4\gamma \max_{s, \mathbf{a}} |A_{\pi}(s, \mathbf{a})|}{(1-\gamma)^2}$.

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.
- ▶ **for** $k = 0, 1, 2, \dots$
- ▶ Compute $C = \frac{4\gamma \max_{s, \mathbf{a}} |A_{\pi}(s, \mathbf{a})|}{(1-\gamma)^2}$.
- ▶ Draw a random permutation of agents $i_{1:n}$.

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.
- ▶ **for** $k = 0, 1, 2, \dots$
- ▶ Compute $C = \frac{4\gamma \max_{s, \mathbf{a}} |A_{\pi}(s, \mathbf{a})|}{(1-\gamma)^2}$.
- ▶ Draw a random permutation of agents $i_{1:n}$.
- ▶ **for** $m = 1, \dots, n$

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.
- ▶ **for** $k = 0, 1, 2, \dots$
- ▶ Compute $C = \frac{4\gamma \max_{s, \mathbf{a}} |A_{\pi}(s, \mathbf{a})|}{(1-\gamma)^2}$.
- ▶ Draw a random permutation of agents $i_{1:n}$.
- ▶ **for** $m = 1, \dots, n$

$$\pi_{k+1}^{i_m} = \arg \max_{\pi^{i_m}} \mathbb{E}_{s \sim \rho_{\pi_k}, \mathbf{a}^{i_{1:m-1}} \sim \pi_{k+1}^{i_{1:m-1}}, \mathbf{a}^{i_m} \sim \pi^{i_m}} [A_{\pi_k}^{i_m}(s, \mathbf{a}^{i_{1:m-1}}, \mathbf{a}^{i_m})] \\ - \text{CD}_{\text{KL}}^{\max}(\pi_k^{i_m}, \pi^{i_m}).$$

Heterogeneous-Agent Trust Region Algorithms

Algorithm 1

- ▶ Initialise a joint policy π_0 at random.
- ▶ **for** $k = 0, 1, 2, \dots$
- ▶ Compute $C = \frac{4\gamma \max_{s, \mathbf{a}} |A_{\pi}(s, \mathbf{a})|}{(1-\gamma)^2}$.
- ▶ Draw a random permutation of agents $i_{1:n}$.
- ▶ **for** $m = 1, \dots, n$

$$\pi_{k+1}^{i_m} = \arg \max_{\pi^{i_m}} \mathbb{E}_{s \sim \rho_{\pi_k}, \mathbf{a}^{i_{1:m-1}} \sim \pi_{k+1}^{i_{1:m-1}}, \mathbf{a}^{i_m} \sim \pi^{i_m}} [A_{\pi_k}^{i_m}(s, \mathbf{a}^{i_{1:m-1}}, \mathbf{a}^{i_m})] - \text{CD}_{\text{KL}}^{\max}(\pi_k^{i_m}, \pi^{i_m}).$$

Theorem. Algorithm 1 achieves the monotonic improvement property, $J(\pi_{k+1}) \geq J(\pi_k)$, and converges to a set of Nash equilibria.

Deep Heterogeneous-Agent Trust Region Algorithms

Deep Heterogeneous-Agent Trust Region Algorithms

Update the parameters to θ_{k+1}^{im} by maximising

Deep Heterogeneous-Agent Trust Region Algorithms

Update the parameters to θ_{k+1}^{im} by maximising

HATRPO

$$\mathbb{E}_{s \sim \rho_{\pi_{\theta_k}}, \mathbf{a}^{1:m-1} \sim \pi_{\theta_{k+1}}^{1:m-1}, \mathbf{a}^{im} \sim \pi_{\theta_{k+1}^{im}}^{im}} [A_{\pi_{\theta_k}}^{im}(s, \mathbf{a}^{1:m-1}, \mathbf{a}^{im})],$$

$$\text{subject to } \mathbb{E}_{s \sim \rho_{\pi_{\theta_k}}} [\text{D}_{\text{KL}}(\pi_{\theta_k}^{im}(\cdot|s), \pi_{\theta_{k+1}^{im}}^{im}(\cdot|s))] \leq \delta.$$

Deep Heterogeneous-Agent Trust Region Algorithms

Update the parameters to θ_{k+1}^{im} by maximising

HATRPO

$$\mathbb{E}_{s \sim \rho_{\pi_{\theta_k}}, \mathbf{a}^{i_1:m-1} \sim \pi_{\theta_{k+1}^{i_1:m-1}}, \mathbf{a}^{im} \sim \pi_{\theta_{k+1}^{im}}} [A_{\pi_{\theta_k}^{im}}(s, \mathbf{a}^{i_1:m-1}, \mathbf{a}^{im})],$$

$$\text{subject to } \mathbb{E}_{s \sim \rho_{\pi_{\theta_k}}} [\text{D}_{\text{KL}}(\pi_{\theta_k^{im}}(\cdot|s), \pi_{\theta_{k+1}^{im}}(\cdot|s))] \leq \delta.$$

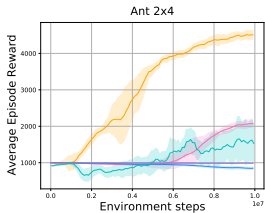
HAPPO

$$\mathbb{E}_{s \sim \rho_{\pi_{\theta_k}}, \mathbf{a} \sim \pi_{\theta_k}} \left[\min \left(\frac{\pi_{\theta_k^{im}}^{im}(\mathbf{a}^i|s)}{\pi_{\theta_k^{im}}^{im}(\mathbf{a}^i|s)} M^{i_1:m}(s, \mathbf{a}), \text{clip} \left(\frac{\pi_{\theta_k^{im}}^{im}(\mathbf{a}^i|s)}{\pi_{\theta_k^{im}}^{im}(\mathbf{a}^i|s)}, 1 \pm \epsilon \right) M^{i_1:m}(s, \mathbf{a}) \right) \right],$$

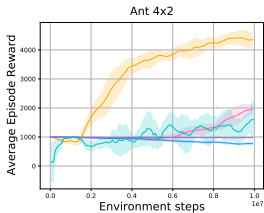
$$\text{where } M^{i_1:m}(s, \mathbf{a}) = \frac{\pi_{\theta_{k+1}^{i_1:m-1}}^{i_1:m-1}(\mathbf{a}^{i_1:m-1}|s)}{\pi_{\theta_k^{i_1:m-1}}^{i_1:m-1}(\mathbf{a}^{i_1:m-1}|s)} A_{\pi}(s, \mathbf{a}).$$

Empirical Results: the New State of the Art

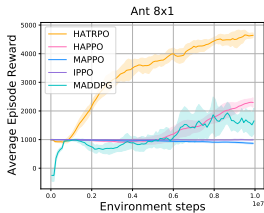
Multi-Agent MuJoCo



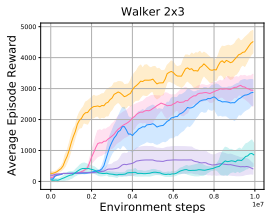
(a) 2x4-Agent Ant



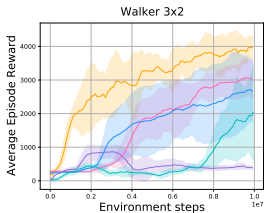
(b) 4x2-Agent Ant



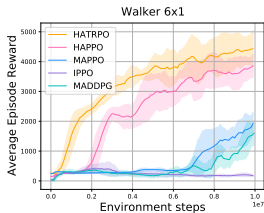
(c) 8x1-Agent Ant



(d) 2x3-Agent Walker



(e) 3x2-Agent Walker



(f) 6x1-Agent Walker

Thank you for your attention!

- ▶ *Jakub Grudzien Kuba* (University of Oxford, Huawei Technologies)
- ▶ *Ruiqing Chen* (ShanghaiTech University)
- ▶ *Muning Wen* (Shanghai Jiao Tong University)
- ▶ *Ying Wen* (Shanghai Jiao Tong University)
- ▶ *Fanglei Sun* (ShanghaiTech University)
- ▶ *Jun Wang* (University College London)
- ▶ *Yaodong Yang* (Peking University)