# DISTRIBUTED ALGORITHM FOR MULTI-OBJECTIVE MULTI-AGENT REINFORCEMENT LEARNING

**Anonymous authors**Paper under double-blind review

000

001

002003004

006

008

009 010

011

012

013

014

016

018

019

021

025

026

027 028 029

031

032

033

034

037

040

041 042

043

044

046

047

048

051

052

#### **ABSTRACT**

Multi-objective reinforcement learning (MORL) aims to optimize multiple conflicting objectives for a single agent, where finding Pareto-optimal solutions is NP-hard and existing algorithms are often centralized with high computational complexity, limiting their practical applicability. Multi-objective multi-agent reinforcement learning (MOMARL) extends MORL to multiple agents, which not only increases computational complexity exponentially due to the global stateaction space, but also introduces communication challenges, as agents cannot continuously communicate with a central coordinator in large-scale scenarios. This necessitates distributed algorithm, where each agent relies only on the information of its neighbors within a limited range rather than depending on the global scale. To address these challenges, we propose a distributed MOMARL algorithm in which each agent leverages only the state of its  $\kappa$ -hop neighbors and locally adjusts the weights of multiple objectives through a consensus protocol. We introduce an approximated policy gradient that reduces the dependency on global actions and a linear function approximation that limits the state space to local neighborhoods. Each agent i's computational complexity is thus reduced from  $\mathcal{O}(|\mathcal{S}||\mathcal{A}|)$ with global state-action space in centralized algorithms to  $\mathcal{O}(|\mathcal{S}_{\mathcal{N}_i^{\kappa}}||\mathcal{A}_i|)$  with  $\kappa$ neighborhood state and local action space. We prove that the algorithm converges to a Pareto-stationary solution at a rate of  $\mathcal{O}(1/T)$  and demonstrate in simulations for robot path planning that our approach achieves higher multi-objective values than state-of-the-art method.

## 1 Introduction

As real-world tasks grow increasingly complex, many scenarios naturally involve multiple conflicting objectives, motivating the study of multi-objective reinforcement learning (MORL). For instance, in robotic path planning (Zhang et al., 2016), sa robotic agsent may aim to simultaneously minimize path length, avoid collisions, and maximize information collection.

Different from the rapid development of traditional reinforcement learning (RL) (Grondman et al, 2012; Zhang et al, 2021), research on MORL (Ge et al., 2022; Stamenkovic et al, 2022) remains in its infancy due to the inherent conflicts among multiple objectives. Unlike scalar-reward RL, in MORL the improvement of one objective may degrade others, making standard policy optimization insufficient. A common approach to tackle MORL is to assign fixed weights to objectives and reduce the problem to a single-objective RL (Blondin & Hale, 2020); however, this requires prior knowledge of objective importance and may fail to explore the full Pareto front. To address this limitation, a more rigorous metric is Pareto optimality, where no objective can be improved without degrading others.

However, for non-convex MORL problems, finding exact Pareto-optimal solutions is NP-hard. Consequently, practical algorithms aim for  $\epsilon$ -Pareto stationary solutions (Sener & Koltun, 2018), which provide a necessary condition for approximate Pareto optimality. On the algorithmic side, for MORL problem with continuous action spaces, (Chen et al., 2021) proposed an actor-critic MORL algorithm based on deterministic policy gradients (Silver et al., 2014) to directly optimize multiple objectives. For MORL with discrete action spaces, (Zhou et al., 2024) introduced a unified multi-objective actor-critic framework applicable to both discounted and average-reward settings, where stochastic policy parameters are updated via a multi-gradient descent approach (Désidéri, 2012), ensuring convergence toward  $\epsilon$ -Pareto stationary solutions.

The aforementioned methods are all directed towards addressing the MORL problem in a centralized setting or for a single agent. However, practical applications of MORL problems often involve multiagents. For instance, teams of robots need to decide themselves how to explore distinct regions by simultaneously minimizing energy consumption and travel time. In comparison to the MORL problem with single-agent, the multi-objective multi-agent problem (MOMARL) (Rădulescu, 2020) is more intricate as it encompasses not only potential conflicts among different objectives but also

interactions between the distributed agents with limited communication. An intuitive approach to the MOMARL problem is to consider it as a MORL problem with a single agent, where the state and action are represented by the joint states and joint actions of all agents, respectively. However, as the number of agents increases, the size of their joint state-action space grows exponentially, and in large-scale scenarios, agents cannot communicate with a central coordinator continuously. These characteristic renders the current algorithms used for solving MORL problems with a single agent (Chen et al., 2021; Zhou et al., 2024) unsuitable for large-scale scenarios with multi-agents. Consequently, the MOMARL problem poses new challenges to the design of scalable algorithms and their theoretical analysis.

In this paper, we aims to address the following problem: *How to develop a fully distributed algorithm for the MOMARL problem and ensure its convergence to Pareto-stationary of the multi-objective function?* The contributions of this paper are described as follows.

- 1. We first propose a novel approximated policy gradient for each agent i, which reduces the global action a required by centralized algorithms to the agent's local action  $a_i$ . Furthermore, to reduce the dimensionality of state information, we employ a linear function approximation to restrict agent i's state to the neighborhood state  $s_{\mathcal{N}_i^{\kappa}}$ , encompassing only its  $\kappa$ -hop neighbors. The per-agent computational complexity is thus reduced from  $\mathcal{O}(|\mathcal{S}||\mathcal{A}|)$  in centralized algorithms to  $\mathcal{O}(|\mathcal{S}_{\mathcal{N}_i^{\kappa}}||\mathcal{A}_i|)$ .
- 2. We propose a novel distributed algorithm in which each agent only uses the policy gradient information from its immediate neighbors, while collaboratively adjusting the weights of multiple objectives via a consensus protocol. This design enables the agents to perform cooperative optimization toward a Pareto-stationary solution, without requiring access to the policy gradient information of agents beyond the direct neighbors.
- 3. We prove that the proposed distributed algorithm, despite relying only on local neighborhood information, achieves convergence to an  $\epsilon$ -Pareto-stationary solution at a rate of  $\mathcal{O}(1/T)$ , matching the convergence speed of centralized algorithms. Moreover, we run simulations in a robot path planning environment and show our algorithm converges to greater multi-objective values as compared to the extension of the latest MORL algorithm (Zhou et al., 2024), and performs close to the central optimum with much shorter running time.

For the sake of convenience, some key functions in this paper are presented in Table 1.

Symbols Annotation  $Q^m(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{\theta})$ Global Q-function in the m-th objective under joint policy  $\pi_{\theta}$  $\tilde{Q}_i^m(m{s},m{a};m{ heta})$ Local Q-function of agent i in the m-th objective under joint policy  $\pi_{\theta}$  $Q_{\mathrm{tru},i}^m(s_{\mathcal{N}_i^{\kappa}},a_{\mathcal{N}_i^{\kappa}};\boldsymbol{\theta})$ Graph-truncated Q-function of agent i in the m-th objective under joint policy  $\pi_{\theta}$  $\nabla_{ heta_i} J_{ ext{tru},i}^m(oldsymbol{ heta})$ Graph-truncated policy gradient of agent i in the m-th objective under joint policy  $\pi_{\theta}$  $\widehat{Q_i^m}(\boldsymbol{s},a_i;\boldsymbol{\theta})$ Action-averaged Q-function of agent i in the m-th objective under policy  $\pi_{\theta}$  $\nabla_{\theta_i} J^m_{\mathrm{app},i}(oldsymbol{ heta})$ Approximated policy gradient of agent i in the m-th objective under joint policy  $\pi_{\theta}$  $\hat{Q}_i^m(s_{\mathcal{N}_i^{\kappa}}, a_i; w_i^m)$ Linear approximation function of agent i in the m-th objective

Table 1: Symbols and functions.

## 2 THE NEW MOMARL PROBLEM FORMULATION AND PRELIMINARIES

#### 2.1 MODEL OF THE MOMARL PROBLEM

The MOMARL problem is described as  $(\mathcal{N}, \mathcal{M}, \mathcal{G}(\mathcal{N}, \mathcal{E}), \{\mathcal{S}_i\}_{i \in \mathcal{N}}, \{\mathcal{A}_i\}_{i \in \mathcal{N}}, \{\mathcal{P}_i\}_{i \in \mathcal{N}}, \rho, \{r_i^m\}_{i \in \mathcal{N}, m \in \mathcal{M}}, \gamma)$ , where  $\mathcal{N} = \{1, \cdots, N\}$  and  $\mathcal{M} = \{1, \cdots, M\}$  represent the agent set and the objective set, respectively.  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  represents the communication network among agents with  $\mathcal{E}$  being the set of edges<sup>1</sup>. For integer  $\kappa \geq 1$ , denote  $\mathcal{N}_i^{\kappa}$  as the  $\kappa$ -hop neighborhood of agent i and  $\mathcal{N}_{-i}^{\kappa} = \mathcal{N} \setminus \mathcal{N}_i^{\kappa}$ .

**State and action**:  $S_i$  and  $A_i$  represent the local state space and the local action space of agent i, respectively. Denote  $S = \prod_{i=1}^N S_i$  and  $A = \prod_{i=1}^N A_i$  as the global state space and the global action space, respectively. Denote  $s = (s_1, \cdots, s_N) \in S$  and  $s = (a_1, \cdots, a_N) \in A$  as the global state and the global action of agents, where  $s_i \in S_i$  and  $s_i \in A_i$  represent the local state and local action

<sup>&</sup>lt;sup>1</sup>For the case of time-varying neighbor agents, our algorithm is still applicable if the agent communicates intermittently (or delays communication) with its initial neighbor. In the process of convergence analysis of the algorithm, we just need to introduce an additional error to (22) caused by communication disconnection or delay.

of agent  $i \in \mathcal{N}$ , respectively. For integer  $\kappa \geq 1$ , denote  $s_{\mathcal{N}_i^{\kappa}}$  and  $a_{\mathcal{N}_i^{\kappa}}$  as the state and action of agent i's  $\kappa$ -hop neighbors, respectively. Denote  $\mathcal{S}_{\mathcal{N}_i^{\kappa}} = \prod_{j \in \mathcal{N}_i^{\kappa}} \mathcal{S}_j$  and  $\mathcal{A}_{\mathcal{N}_i^{\kappa}} = \prod_{j \in \mathcal{N}_i^{\kappa}} \mathcal{A}_j$  as the state space and the action space of agent i's  $\kappa$ -hop neighbors, respectively. Moreover, denote  $\mathcal{S}_{\mathcal{N}_{-i}^{\kappa}}$  and  $\mathcal{A}_{\mathcal{N}_{-i}^{\kappa}}$  as the state space and the action space of agents excluding agent i's  $\kappa$ -hop neighbors, respectively.  $a_{\mathcal{N}_{-i}^{\kappa}} \in \mathcal{S}_{\mathcal{N}_{-i}^{\kappa}}$ 

State transition probability function:  $\mathcal{P}_i(s_i'|s_{\mathcal{N}_i^1},a_i): \mathcal{S}_{\mathcal{N}_i^1} \times \mathcal{A}_i \times \mathcal{S}_i \to [0,1]$  is the state transition probability function of agent i, dependent of its 1-hop neighborhood state and its local action. Denote  $\mathcal{P}(s'|s,a) = \prod_{i=1}^N \mathcal{P}_i(s_i'|s_{\mathcal{N}_i^1},a_i): \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0,1]$  as the global state transition probability function. Note that the definition of the state transition probability function  $\prod_{i=1}^N \mathcal{P}_i(s_i'|s_{\mathcal{N}_i^1},a_i)$  is common in the literature. For example, it applies to the scenario of traffic signal control problem (Chu et al., 2020), where the traffic flow at each intersection is influenced by the traffic flow at its neighboring intersections and its own signal light.

**Initial state distribution**:  $\rho$  is the distribution of the initial state  $s_0$ .

**Reward function:**  $r_i^m(s_i, a_i) : \mathcal{S}_i \times \mathcal{A}_i \to \mathbb{R}$  is the reward function of agent  $i \in \mathcal{N}$  in the objective  $m \in \mathcal{M}$ . Denote  $s_t = (s_{1,t}, \cdots, s_{N,t})$  and  $a_t = (a_{1,t}, \cdots, a_{N,t})$  as the global state and the global action at time t, respectively. The reward of agent  $i \in \mathcal{N}$  in the objective  $m \in \mathcal{M}$  at time t can be represented as  $r_{i,t}^m = r_i^m(s_{i,t}, a_{i,t})$ .

**Discount factor**:  $\gamma = (\gamma^1, \cdots, \gamma^M)^\top \in \mathbb{R}^M$  with  $\gamma^m \in (0, 1)$  being the discount factor in the objective  $m \in \mathcal{M}$ .

**Softmax policy**: In this paper, we use the parameterized softmax policy  $\pi_{\theta_i}(a_i|s_i)$  with parameter  $\theta_i \in \mathbb{R}^{|\mathcal{S}_i||\mathcal{A}_i|}$ , which is described as

$$\pi_{\theta_i}(a_i|s_i) = \frac{\exp(\theta_{i,s_i,a_i})}{\sum_{a_i'} \exp(\theta_{i,s_i,a_i'})},$$
(1)

where  $\theta_{i,s_i,a_i}$  represents the element corresponding to  $(s_i,a_i)$  in  $\theta_i$ . Denote  $\boldsymbol{\theta} = (\theta_1^\top,\cdots,\theta_N^\top)^\top \in \mathbb{R}^{\sum_{i=1}^N |S_i||A_i|}$  as the joint policy parameter of agents and  $\boldsymbol{\pi}_{\boldsymbol{\theta}}(\boldsymbol{a}|\boldsymbol{s}) = \prod_{i=1}^N \pi_{\theta_i}(a_i|s_i)$  be the joint policy of all agents. Note that the softmax policy is used in RL to ensure the exploration of agents (Zhou et al., 2023; Zhang et al., 2022).

In the MOMARL problem, given a joint policy parameter  $\theta$ , the m-th objective of all agents is defined as

$$J^{m}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{s} \sim \boldsymbol{\rho}} \left[ \frac{1}{N} \sum_{t=0}^{\infty} \sum_{i=1}^{N} (\gamma^{m})^{t} r_{i,t}^{m} | \boldsymbol{s}_{0} = \boldsymbol{s}, \boldsymbol{a}_{t} \sim \boldsymbol{\pi}_{\boldsymbol{\theta}}(\cdot | \boldsymbol{s}_{t}) \right], \tag{2}$$

which is represents the average discounted reward of all agents over all time t. The goal of agents in the MOMARL problem is to find a joint policy parameter  $\theta$  to maximize the following composite objective, i.e.,

$$\max_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}) = [J^1(\boldsymbol{\theta}), \cdots, J^M(\boldsymbol{\theta})]^\top \in \mathbb{R}^M.$$
 (3)

In order to address the potential conflicts among the  $J(\theta)$  in (3), the notions of Pareto-optimality and  $\epsilon$ -Pareto-stationarity are introduced as follows.

**Definition 1** (Pareto-optimality) A solution  $\theta$  dominates solution  $\theta'$  if and only if  $J^m(\theta) \geq J^m(\theta')$ ,  $\forall m \in \mathcal{M}$  and  $\exists m' \in \mathcal{M}$ ,  $J^{m'}(\theta) > J^{m'}(\theta')$ . A solution  $\theta$  is Pareto-optimal if it is not dominated by any other solution.

Considering that finding Pareto-optimal solutions for non-convex MOMARL problems is NP-hard, it is generally more practical to seek the  $\epsilon$ -Pareto-stationary solution instead of the Pareto-optimal solution (Kumar et al., 2019).

**Definition 2** ( $\epsilon$ -Pareto-stationarity) Define  $\nabla_{\theta} J(\theta)$  as the gradient of  $J(\theta)$  respect to  $\theta$ . A solution  $\theta$  is  $\epsilon$ -Pareto stationary if there exists  $\lambda = (\lambda^1, \dots, \lambda^M)^{\top} \in \mathbb{R}^M$  such that  $\min_{\lambda \in \mathbb{R}^M} \|\nabla_{\theta} J(\theta)^{\top} \lambda\|_2^2 \leq \epsilon$  with  $\lambda \geq 0$ ,  $\|\lambda\|_1 = 1$ , and  $\epsilon > 0$ .

Based on Definitions 1-2, it is obvious that the Pareto-stationarity is a necessary condition for a solution to be Pareto-optimal. Specifically, in the context of convex MOMARL problems, the solutions that are Pareto-stationary also qualify as Pareto-optimal. Given the complexity associated with the MOMARL problem, this paper focuses on developing a distributed scalable algorithm to identify and achieve Pareto-stationarity.

#### 2.2 PRELIMINARIES IN THE MOMARL PROBLEM

In the MOMARL problem, for any joint policy parameter  $\theta$  and  $m \in \mathcal{M}$ , the global Q-function  $Q^m(s, a; \theta)$  in m-th objective is defined as

$$Q^{m}(s, \boldsymbol{a}; \boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\pi}\boldsymbol{\theta}} \left[ \frac{1}{N} \sum_{t=0}^{\infty} \sum_{i=1}^{N} (\gamma^{m})^{t} r_{i,t}^{m} | \boldsymbol{s}_{0} = \boldsymbol{s}, \boldsymbol{a}_{0} = \boldsymbol{a} \right], \tag{4}$$

which represents the value of the stat-action pair (s, a) in m-th objective under join policy  $\pi_{\theta}$ . In the MOMARL problem, given the joint policy parameter  $\theta$ , define  $d_{\rho}^{\theta,m}(s)$  as the discounted state visitation distribution, which is represented as

$$d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(\boldsymbol{s}) = (1 - \gamma^m) \sum_{t=0}^{\infty} (\gamma^m)^t \operatorname{Pr}^{\boldsymbol{\pi}_{\boldsymbol{\theta}}}(\boldsymbol{s}_t = \boldsymbol{s} | \boldsymbol{s}_0 \sim \boldsymbol{\rho}), \tag{5}$$

where  $\Pr^{\pi_{\theta}}(s_t = s | s_0 \sim \rho)$  represents the probability of  $s_t = s$  at time t under the initial state distribution  $\rho$  and the joint policy  $\pi_{\theta}$ .

Recall that the policy gradient theorem (Sutton et al., 2000) is the foundation of algorithm design in RL. Inspired by the theorem, in our MOMARL problem, we also have the following policy gradient lemma.

**Lemma 1** In the MOMARL problem, for any joint policy parameter  $\theta$ , the gradient of  $J^m(\theta)$  in m-the objective with respect to  $\theta$  is given by:

$$\nabla_{\boldsymbol{\theta}} J^{m}(\boldsymbol{\theta}) = \frac{1}{1 - \gamma^{m}} \mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}, m}, \boldsymbol{a} \sim \boldsymbol{\pi}_{\boldsymbol{\theta}}} [\nabla_{\boldsymbol{\theta}} \log \boldsymbol{\pi}_{\boldsymbol{\theta}}(\boldsymbol{a} | \boldsymbol{s}) Q^{m}(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{\theta})], \forall m \in \mathcal{M}.$$
 (6)

For the policy gradients involved in Definition 2, Lemma 1 shows that the calculation of the policy gradient  $\nabla_{\theta}J^{m}(\theta)$  depends on  $Q^{m}(s,a;\theta)$ , which involves global state-action (s,a). Consequently, there are two challenges in applying (6): (i) the computational complexity of handling the global state-action (s,a) in a centralized setting is high; (ii) achieving distributed decision making among multi-agents with limited communication.

# 3 DISTRIBUTED SCALABLE ACTOR-CRITIC ALGORITHM FOR MOMARL PROBLEM

Before presenting our proposed method, we first revisit the centralized approach for solving MO-MARL problems. In the centralized setting, the policy update requires the global state-action pair (s,a). This immediately leads to an exponential growth of the joint state-action space with the number of agents, i.e., of order  $\mathcal{O}(|\mathcal{S}||\mathcal{A}|)$ . Such computational complexity makes centralized algorithms prohibitive for large-scale systems. Moreover, centralized training implicitly assumes that agents can constantly communicate with a central controller, which is unrealistic in many real-world scenarios (e.g., swarm robotics, sensor networks). To overcome these limitations, the natural way forward is to make the algorithm scalable, meaning that the per-agent computational cost should remain polynomial in the state-action dimension of its neighbors within a limited range rather than depending on the global scale. The most effective way to achieve scalability is to design a distributed algorithm, where (i) each agent i makes decisions independently using only  $(s_{\mathcal{N}_i^{\kappa}}, a_i)$ , where  $s_{\mathcal{N}_i^{\kappa}}$  is the state from its  $\kappa$ -hop neighbors and (2) exchanges its local Lagrangian multiplier  $\lambda_i$  with its directly neighbors for balancing multiple objectives. This shift from centralized to distributed design is the key idea underlying our algorithm: by restricting each agent's decision to local stateaction information and coordinating objective trade-offs through a consensus protocol, we preserve scalability while maintaining rigorous convergence guarantees.

To reduce the reliance of algorithm on global information, we design a distributed algorithm, where each agent (i) estimates a local policy gradient based only on its own action, (ii) leverages a linear approximation restricted to  $\kappa$ -hop neighborhood states, and (iii) updates multi-objective weights through a consensus protocol to cooperatively approach Pareto-stationary solutions.

# 3.1 A NOVEL APPROXIMATED POLICY

In contrast to the global Q-function utilized in (6), which relies on global state-action information, we introduce a novel concept the "action-averaged Q-function" for each agent i. This formulation leverages rewards from agent i's  $\kappa$ -hop neighbors to effectively reduce the dependence on the full joint action a by focusing on the local action  $a_i$ , as defined below:

$$\widehat{Q_i^m}(\boldsymbol{s}, a_i; \boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\pi}_{\boldsymbol{\theta}}} \left[ \frac{1}{N} \sum_{t=0}^{\infty} (\gamma^m)^t \sum_{j \in \mathcal{N}_i^n} r_j^m(s_{j,t}, a_{j,t}) | \boldsymbol{s}_0 = \boldsymbol{s}, a_{i,0} = a_i \right].$$
 (7)

To approximate the true policy gradient  $\nabla_{\theta} J^m(\theta)$  in (6), we define  $\nabla_{\theta_i} J^m_{app}(\theta)$  as an approximated policy gradient for agent i, derived using the action-averaged Q-function presented in (7), as follows:

$$\nabla_{\theta_i} J_{app,i}^m(\boldsymbol{\theta}) = \frac{1}{1 - \gamma^m} \mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}, a_i \sim \pi_{\theta_i}} \left[ \widehat{Q_i^m}(\boldsymbol{s}, a_i; \boldsymbol{\theta}) \nabla_{\theta_i} \log \pi_{\theta_i}(a_i | s_i) \right]. \tag{8}$$

Unlike the policy gradient in centralized algorithm that requires global action a, (8) only requires the local action  $a_i$  of agent i. The approximation error between  $\nabla_{\theta_i} J^m_{app,i}(\theta)$  and original  $\nabla_{\theta_i} J^m(\theta)$  in (6) can be well bounded for the MOMARL problem in the following theorem.

**Theorem 1** *In the MOMARL problem, given a joint policy*  $\pi_{\theta}$ *, for any agent*  $i \in \mathcal{N}$  *and objective*  $m \in \mathcal{M}$ *, it holds that* 

$$\|\nabla_{\theta_i} J_{app,i}^m(\boldsymbol{\theta}) - \nabla_{\theta_i} J^m(\boldsymbol{\theta})\|_2 \le \frac{\sqrt{2}R}{(1 - \gamma^m)^2} (\gamma^m)^{\kappa + 1}. \tag{9}$$

The proof of Theorem 1 is provided in Appendix A.2. The policy gradient has been approximated so far by constructing  $\widehat{Q_i^m}(s,a_i;\boldsymbol{\theta})$  in (7) and  $\nabla_{\theta_i}J_{app,i}^m(\boldsymbol{\theta})$  in (8), which reduces the action dimension of each agent i to its local action  $a_i$ . However, the expression of  $\widehat{Q_i^m}(s,a_i;\boldsymbol{\theta})$  still requires the global state. Therefore, in the following, we will focus on reducing the dimensionality of agents' state information.

# 3.2 CRITIC STEP: LINEAR FUNCTION APPROXIMATION

# Algorithm 1: Linear function approximation

```
1 Require: The number of samples K, the learning-rate \eta_w^m and \varepsilon > 0;
```

- 2 Initialization: Initialize the  $\varepsilon$ -exploration policy  $\pi_{\theta}^{\varepsilon} = \prod_{i=1}^{N} \pi_{\theta_i}^{\varepsilon}$ , where  $\pi_{\theta_i}^{\varepsilon}(a_i|s_i) = (1-\varepsilon)\pi_{\theta_i}(a_i|s_i) + \frac{\varepsilon}{|\mathcal{A}_i|}$  for all  $i \in \mathcal{N}$ . The initial values of the parameters  $w_{i,0}^m$  is set as  $w_{i,0}^m = \mathbf{0}_{d_i}$  for all  $i \in \{1, 2, \cdots, N\}$ ;
- 3 The agents execute the  $\varepsilon$ -exploration policy  $\pi^{\varepsilon}_{\theta}$  and each agent  $i \in \mathcal{N}$  collects a sequence of samples  $\{(s_{i,k}, a_{i,k}, r^{i,k}_{i,k})\}_{0 \leq k \leq K}$  in m-the objective;

```
4 for i=1,2,\cdots,N do

5 For each objective m\in\mathcal{M}, agent i\in\mathcal{N} collects the state information \{s_j\}_{j\in\mathcal{N}_i^\kappa} of its \kappa-hop neighbors and reward \{r_j^m\}_{j\in\mathcal{N}_i^\kappa} from its \kappa-hop neighbors to form a sample set \{s_{\mathcal{N}_i^\kappa,k},a_{i,k},r_{\mathcal{N}_i^\kappa,k}^m\}_{0\leq k\leq K};

for k=0,1,2,\cdots,K-1 do

Each agent i\in\mathcal{N} estimates its local temporal difference error:

\delta_{i,k}^m=\phi_i(s_{\mathcal{N}_i^\kappa,k},a_{i,k})^\top w_{i,k}^m-\frac{1}{N}\sum_{j\in\mathcal{N}_i^\kappa}r_{j,k}^m-\gamma^m\phi_i(s_{\mathcal{N}_i^\kappa,k+1},a_{i,k+1})^\top w_{i,k}^m;
w_{i,k+1}^m=w_{i,k}^m-\eta_w^m\delta_{i,k}^m\phi_i(s_{\mathcal{N}_i^\kappa,k+1},a_{i,k+1});
end
```

10 end

11 Output:  $\{w_{i,K}^m\}_{i\in\mathcal{N},m\in\mathcal{M}}$ 

In this subsection, we use the localized stochastic approximation and propose a linear function in (10) to reduce the dimension of the state-action required by agent  $i \in \mathcal{N}$  to  $(s_{\mathcal{N}_i^\kappa}, a_i)$ . Specially, the linear function  $\hat{Q}_i^m(s_{\mathcal{N}_i^\kappa}, a_i; w_i^m)$  of agent i to approximate  $\widehat{Q}_i^m(s, a_i; \theta)$  in (7) is given as

$$\hat{Q}_i^m(s_{\mathcal{N}_i^\kappa}, a_i; w_i^m) = \phi_i(s_{\mathcal{N}_i^\kappa}, a_i)^\top w_i^m, \tag{10}$$

where  $\phi_i(s_{\mathcal{N}_i^{\kappa}}, a_i): \mathcal{S}_{\mathcal{N}_i^{\kappa}} \times \mathcal{A}_i \to \mathbb{R}^{d_i}$  is the feature vector mapping and  $w_i^m \in \mathbb{R}^{d_i}$  is the parameter of agent i in m-th objective. By the definition of  $\widehat{Q_i^m}(s, a_i; \theta)$  in (7), the parameter with initial value  $w_{i,0}^m$  can be updated by sample sequence  $\{s_{\mathcal{N}_i^{\kappa},k}, a_{i,k}, r_{\mathcal{N}_i^{\kappa},k}^m \}_{0 \le k \le K}$  as

$$w_{i,k+1}^{m} = w_{i,k}^{m} - \eta_{w}^{m} \delta_{i,k}^{m} \phi_{i}(s_{\mathcal{N}_{i}^{\kappa},k+1}, a_{i,k+1}), \tag{11}$$

where K is the number of sample for linear parameter training and  $\delta^m_{i,k}$  is the local temporal difference error, which is discribed as

$$\delta_{i,k}^{m} = \phi_{i}(s_{\mathcal{N}_{i}^{\kappa},k}, a_{i,k})^{\top} w_{i,k}^{m} - \frac{1}{N} \sum_{i \in \mathcal{N}^{\kappa}} r_{j,k}^{m} - \gamma^{m} \phi_{i}(s_{\mathcal{N}_{i}^{\kappa},k+1}, a_{i,k+1})^{\top} w_{i,k}^{m}, \tag{12}$$

and  $\eta_w^m$  is the fixed learning rate of parameters  $w_i^m$ . The detailed description of linear function approximation is given in Algorithm 1.

#### ACTOR STEP: POLICY PARAMETER UPDATE

To estimate the policy gradient and ensure convergence toward a Pareto-stationary solution, we need to dynamically adjust the weights associated with multiple objectives. In centralized algorithms, this requires access to all agents' policy gradients, which is often infeasible and unscalable (Zhou et al., 2024). To overcome this limitation, we design a distributed algorithm in which agents collaboratively adjust the objective weights through a consensus protocol. This design removes the dependence on global policy gradient while still enabling cooperative optimization toward Pareto-stationary so-

To estimate the approximated policy gradient  $\nabla_{\theta_i} J^m_{app,i}(\boldsymbol{\theta})$  in (8), for joint policy  $\boldsymbol{\pi}_{\boldsymbol{\theta}_t}$ , our estimate  $g^m_{i,t}(B)$  is calculated iteratively based on the sample sequence  $\{(s^b_{\mathcal{N}^\kappa_i,h},a^b_{i,h})\}_{0\leq b\leq B-1,0\leq h\leq H-1}$ :

$$g_{i,t}^{m}(b+1) = \frac{b}{b+1}g_{i,t}^{m}(b) + \frac{1}{b+1}\widehat{\nabla}_{\theta_{i}}J_{app,i}^{m,b}(\boldsymbol{\theta}_{t}), \tag{13}$$

where  $g_{i,t}^m(0) = \mathbf{0}_{|\mathcal{S}_i||\mathcal{A}_i|}$  and

$$\widehat{\nabla}_{\theta_i} J_{app,i}^{m,b}(\boldsymbol{\theta}_t) = \sum_{h=0}^{H-1} (\gamma^m)^h \phi_i(s_{\mathcal{N}_i^{\kappa},h}^{\kappa}, a_{i,h}^b)^{\top} w_i^m(t) \nabla_{\theta_i} \log \pi_{\theta_{i,t}}(a_{i,h}^b | s_{i,h}^b)$$
(14)

with  $w_i^m(t)$  being the output of Algorithm 1 in t-th iteration of policy parameters. Let  $g_{i,t}^m =$  $g^m_{i,t}(B)^{ op}$  and  $\boldsymbol{g}^m_t = \left((g^m_{1,t})^{ op},\cdots,(g^m_{N,t})^{ op}
ight)^{ op} \in \mathbb{R}^{\sum_{i=1}^N |\mathcal{S}_i||\mathcal{A}_i|}$ . Following Pareto-stationarity in Definition 1, we denote  $\boldsymbol{\lambda}^*_t = (\lambda^{*1}_t,\cdots,\lambda^{*M}_t)^{ op} \in \mathbb{R}^M$  as solution of the following quadratic programming problem:

$$\min_{\boldsymbol{\lambda}_{t} = (\lambda_{t}^{1}, \dots, \lambda_{t}^{M})^{\top} \in \mathbb{R}^{M}} J_{t}^{g}(\boldsymbol{\lambda}_{t}) = \left\| \sum_{m=1}^{M} \lambda_{t}^{m} \boldsymbol{g}_{t}^{m} \right\|_{2}^{2}$$
s.t.  $\boldsymbol{\lambda}_{t} \geq 0, \|\boldsymbol{\lambda}_{t}\|_{1} = 1.$  (15)

For the network  $\mathcal{G}(\mathcal{N}, \mathcal{E})$  among agents, we define its weight matrix as  $W^{\mathcal{G}} = [w_{ij}^{\mathcal{G}}]_{N \times N}$ , where each element  $w_{ij}^{\mathcal{G}}$  represents the weight of the edge from agent j to agent i, which is defined as

$$w_{ij}^{\mathcal{G}} = \begin{cases} \frac{1}{1 + \max(|\mathcal{N}_i|, |\mathcal{N}_j|)}, & j \in \mathcal{N}_i, \\ 1 - \sum_{l \in \mathcal{N}_i} w_{il}^{\mathcal{G}}, & j = i, \\ 0, & \text{otherwise.} \end{cases}$$
 (16)

By using the definition of  $W^{\mathcal{G}}$  in (16), we solve the problem (15) in a distributed way, which is presented in the following Algorithm 2.

# Algorithm 2: Distributed computation to solve problem (15) of objective weight adjustment

```
305
          1 Initialization: Each agent i \in \mathcal{N} sets \lambda_i(0) = \frac{1}{M} \mathbf{1}_M and chooses step sizes \alpha_k = \frac{2}{k+2} for all k \geq 0;
306
          2 for k = 0, 1, 2, \dots, K_{\lambda} - 1 do
307
                    Each agent i initializes x_{i,t}(k) = \sum_{m=1}^{M} \lambda_{i,m}(k)g_{i,t}^{m};
308
                    for m=1,2,\cdots,M do
                           Each agent i initializes y_{i,t}^m(k,0) = \langle x_{i,t}(k), g_{i,t}^m \rangle;
          5
                          while \exists i \in \mathcal{N}, y_{i,t}^m(l_1) \neq \frac{1}{N} \sum_{i=1}^N y_{i,t}^m(k,0) \text{ do}
y_{i,t}^m(k,l_1+1) = \sum_{j \in \mathcal{N}_i} w_{ij}^g y_{j,t}^m(k,l_1);
l_1 \leftarrow l_1 + 1;
310
311
          7
312
           8
313
314
                           Each agent i obtains local output y_{i,t}^m(k);
         10
315
         11
                     Each agent obtains u_{i,t}^*(k) = \arg\min_m y_{i,t}^m(k);
         12
316
                     Each agent updates \lambda_i(k+1) = (1-\alpha_k)\lambda_i(k) + \alpha_k e_{u_{i,t}^*(k)}, where e_{u_{i,t}^*(k)} is an M-dimensional unit
         13
317
                     vector with the u_{i,t}^*(k)-th element being 1, and the other elements being 0;
318
319
```

15 Output:  $\widehat{\lambda}_t = \lambda_i(K_{\lambda})$  for all  $i \in \mathcal{N}$ ;

In Algorithm 2 each agent iteratively computes policy gradients, engages in consensus steps to evaluate the quadratic objective (ref. Lines 5-9), and updates its local weight vector via a Frank-Wolfe update rule (ref. Lines 12-13). Through repeated consensus and update steps, all agents asymptot-

ically agree on the optimal weight vector  $\hat{\lambda}_t$ , thereby achieving a distributed solution that approximates the centralized Pareto-stationary weighting without requiring global information exchange.

After computing  $\lambda_t$  by Algorithm 2, we update the weight  $\lambda_t$  as

$$\lambda_t = (1 - \eta_{\lambda,t})\lambda_{t-1} + \eta_{\lambda,t}\widehat{\lambda}_t, \tag{17}$$

where  $\eta_{\lambda,t}$  is the learning rate of  $\lambda_t$ . Denote  $g_t = \sum_{m=1}^M \lambda_t^m g_t^m$ , the update of  $\theta_{t+1}$  is presented as  $\theta_{t+1} = \theta_t + \eta_{\theta,t} g_t$ , (18)

where  $\eta_{\theta,t}$  is the learning rate of policy parameter. In the NMARL problem, the agents can use  $\theta_t$  to execute the actions based on (1).

### 3.4 OVERALL DISTRIBUTED ALGORITHM FOR MOMARL PROBLEM

Based on the distributed designs in the previous three subsections, we propose a distributed MO-MARL algorithm, which is given in Algorithm 3.

# **Algorithm 3:** Distributed algorithm for MOMARL problem

```
Require: The non-negative integers T, B, \overline{H}, the learning-rates \{\eta_{\lambda,t}\}_{t\in\{1,\cdots,T\}} and \{\eta_{\theta,t}\}_{t\in\{1,\cdots,T\}};
339
          2 Initialization: Initialize \lambda_0 = \frac{1}{M} \mathbf{1}_M \in \mathbb{R}^M, the policy parameter \theta_{i,1} \in \mathbb{R}^{|\mathcal{S}_i| \times |\mathcal{A}_i|} to follow Gaussian
340
341
             distribution for all i \in \{1, 2, \dots, N\};
         3 for t = 1, 2, \dots, T do
342
                  Initial policy gradient estimation g_{i,t}^m(0) = \mathbf{0}_{|\mathcal{S}_i||\mathcal{A}_i|} for all i \in \mathcal{N};
343
         5
                   Critic step: All agents use Algorithm 1 and output the weight vectors \{w_i^m(t)\}_{i\in\mathcal{N}};
344
                  for b = 0, 1, 2, \dots, B - 1 do
345
                        All agents execute the joint policy \pi_{\theta_t} in H-1 horizon;
         8
                        Each agent i \in \mathcal{N} collects a sequence of samples, which includes the state information \{s_j\}_{j \in \mathcal{N}_i^e}
         9
347
                         from its \kappa-hop neighbors and its local action information a_i, i.e., \{(s_{\mathcal{N}_i^{\kappa},h}^b, a_{i,h}^b)\}_{0 \leq h \leq H-1};
348
                        Each agent i estimates the local policy gradient in m-th objective according to (13);
         10
349
                  end
         11
350
                   All agents calculate g_{i,t}^m = g_{i,t}^m(B) by (13) and achieve \boldsymbol{g}_t^m = \left( (g_{1,t}^m)^\top, \cdots, (g_{N,t}^m)^\top \right)^\top for all m \in [M];
         12
351
                   Compute \hat{\lambda}_t by Algorithm 2 as the approximation solution to problem (15);
         13
352
                   Update the weight \lambda_t according to (17);
         14
353
         15
                   Update the policy parameter \theta_{t+1} according to (18);
354
         16 end
        17 Output: \pi_{\theta_{\hat{T}}} with \hat{T} chosen uniformly from \{1, \cdots, T\}
355
```

Algorithm 3 incorporates linear function approximation in Algorithm 1 and distributed consensusbased adjustment of multiplier objective weights in Algorithm 2 into a unified framework. The primary advantages are as follows: (1) each agent relies solely on state information from its  $\kappa$ -hop neighbors, thereby avoiding the exponential expansion of the centralized joint state-action space (ref. Line 9); (2) Each agent in Line 13 only uses the policy gradient estimations of its direct neighbors, which eliminates the requirement for a central coordinator while supporting collaborative multiobjective optimization.

## 4 PARETO-STATIONARY CONVERGENCE OF ALGORITHM 3

Before the convergence analysis of the algorithm, some assumptions are introduced in the following.

**Assumption 1** In the MOMARL problem, for any joint policy parameter  $\theta$  and objective  $m \in \mathcal{M}$ ,  $\xi_{\rho}^{\theta,m}(s,a)$  satisfies that

$$\inf_{\boldsymbol{\theta}} \min_{(\boldsymbol{s}, \boldsymbol{a}) \in \boldsymbol{S} \times \boldsymbol{\mathcal{A}}} \xi_{\boldsymbol{\rho}}^{\boldsymbol{\theta}, m}(\boldsymbol{s}, \boldsymbol{a}) > 0.$$
 (19)

**Assumption 2** In the MOMARL problem, for any agent  $i \in \mathcal{N}$  and objective  $m \in \mathcal{M}$ , there exists constant R > 1 such that the instantaneous reward  $r_{i,t}^m$  at time  $t \geq 0$  satisfies  $|r_{i,t}^m| \leq R$ .

**Assumption 3** *In the MOMARL problem, the network*  $\mathcal{G}(\mathcal{N}, \mathcal{E})$  *among agents is connected graph.* 

Assumption 1 ensures that for any joint policy  $\pi_{\theta}$ ,  $(s, a) \in \mathcal{S} \times \mathcal{A}$  is visited with a non-zero probability, Assumption 2 provides an upper bound on the reward, and Assumption 3 is provided for designing distributed algorithm. Assumptions 1-2 are standard prerequisite for the convergence

analysis of RL algorithms (e.g. (Zhou et al., 2023; Zhang et al., 2022)) and Assumption 3 is commonly in (Olfati-Saber & Murray, 2004).

Our process to prove the Pareto-stationary convergence of Algorithm 3 is as follows: (i) We start establish the smoothness of objective function (i.e., 2); (ii) We from the definition of Pareto-stationarity in Definition 2 and analyze the error between the true gradient  $\nabla_{\theta_i} J^m(\theta_t)$  and the calculated gradient  $g^m_{i,t}$  in (13) (i.e., Lemma 3); (iii) We control  $\lambda_t$  by setting the step size  $\eta_{\theta,t}$  to ensure that Algorithm 3 converges to Pareto-stationary solution in Theorem 2.

**Lemma 2** In the MOMARL problem, define  $L_J = \max_{m \in \mathcal{M}} \frac{6N}{(1-\gamma^m)^3}$ . For any objective  $m \in \mathcal{M}$ , the objective  $J^m(\theta)$  is  $L_J$ -smooth.

The detailed proof of Lemma 2 can be found in Appendix A.3. Lemma 2 establishes that each individual objective in the MOMARL problem is  $L_J$ -smooth with respect to the joint policy parameters. This smoothness property implies that the gradient of  $J^m(\theta)$  does not change abruptly when the policy parameters are slightly perturbed, which is crucial for analyzing the stability and convergence of gradient-based algorithms. In particular,  $L_J$ -smoothness allows us to control the error propagation when using approximate gradients or local updates in a distributed setting, and it forms a key technical result for deriving convergence rates toward Pareto-stationary solutions.

Define  $\varepsilon_{critic}$  below as the linear approximation error in Algorithm 1:

$$\varepsilon_{critic} = \sup_{m \in \mathcal{M}} \sup_{\boldsymbol{\theta}} \sup_{i \in \mathcal{N}} \mathbb{E} \left[ \sup_{\boldsymbol{s}, a_i} \left| \hat{Q}_i(s_{\mathcal{N}_i^{\kappa}}, a_i; w_{i,K}^m) - \widehat{Q}_i^m(\boldsymbol{s}, a_i; \boldsymbol{\theta}) \right|^2 \right]. \tag{20}$$

Based on  $\varepsilon_{critic}$ , we further define the gradient approximation error in Section 3.3 as

$$\varepsilon_{actor}^{m} = \underbrace{\frac{8R^{2}}{(1 - \gamma^{m})^{4}} (\gamma^{m})^{2\kappa + 2}}_{\text{Truncation error}} + \underbrace{\frac{32}{(1 - \gamma^{m})^{2}B} + \frac{8(\gamma^{m})^{2H}}{(1 - \gamma^{m})^{4}}}_{\text{Sampling error}} + \underbrace{\frac{8\varepsilon_{critic}}{(1 - \gamma^{m})^{2}}}_{\text{Linear approximation error}}.$$
 (21)

**Lemma 3** In Algorithm 3, for joint policy parameter  $\theta_t$ , any agent  $i \in \mathcal{N}$ , and objective  $m \in \mathcal{M}$ , we have

$$\mathbb{E}[\|\nabla_{\theta_i} J^m(\boldsymbol{\theta}_t) - g_{i,t}^m\|_2^2] \le \varepsilon_{actor}^m.$$

The proof of the Lemma 3 is given in Appendix A.4. Lemma 3 provides an upper bound on the error between the true policy gradient and the estimated gradient  $g^m_{i,t}$  used by each agent in Algorithm 3. This result quantifies the approximation error induced by the use of a critic or a finite-horizon estimator in computing policy gradients. Importantly, the bound  $\varepsilon^m_{actor}$  captures the combined effects of truncation, sampling, and linear function approximation errors. Based on Lemma 3, the convergence of Algorithm 2 is established in the following theorem.

**Proposition 1** In Algorithm 2, for any iteration t, we define  $G_t = (\boldsymbol{g}_t^1, \dots, \boldsymbol{g}_t^m) \in \mathbb{R}^{M \times \sum_{i=1}^N |\mathcal{S}_i||\mathcal{A}_i|}$  and have the following result:

- (i) In Line 10 of Algorithm 2,  $y_{i,t}^m(k) = \frac{1}{N} \sum_{i=1}^N \langle \sum_{m=1}^M \lambda_{i,m}(k) g_{i,t}^m, g_{i,t}^m \rangle$ ;
- (ii) In Problem 15,  $J_t^g(\cdot)$  is  $L_t^g$ -Lipschitz continuous with  $L_t^g=2\sigma_{max}^2(G_t)$ , where  $\sigma_{max}(G_t)$  is the largest singular value of  $G_t$ ;

$$(iii) |J_t^g(\widehat{\boldsymbol{\lambda}}_t) - J_t^g(\boldsymbol{\lambda}_t^*)| \le \frac{4L_t^g}{K_{\lambda}+1}.$$

The proof of Proposition 1 is presented in Appendix A.5. Proposition 1 guarantees that, despite the distributed nature of the update and the use of only local gradient information from neighbors, the agents can collectively achieve a close approximation to the globally optimal weight allocation for combining multiple objectives. This property is fundamental for ensuring the cooperative convergence of the distributed algorithm toward a Pareto-stationary solution. Based on Proposition 1, the Pareto-stationary convergence of Algorithm 3 is presented in the following theorem.

**Theorem 2** In Algorithm 3, let  $\eta_{\theta,t} = \frac{1}{3L_J}$ , and  $\eta_{\lambda,t} = \frac{1}{(t+1)^2}$ . Our policy parameter sequences  $\{\theta_t\}_{t=1}^T$  generated by Algorithm 3 satisfies:

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_{t})^{\top} \widehat{\boldsymbol{\lambda}}_{t}\|_{2}^{2}] \leq \frac{216N}{(1-\|\boldsymbol{\gamma}\|_{\infty})^{3}T} \left(1 + \sum_{t=1}^{T} \eta_{\boldsymbol{\lambda},t}\right) + 5 \max_{m \in \mathcal{M}} \varepsilon_{actor}^{m} + \frac{8}{K_{\lambda} + 1} \left(\max_{m \in \mathcal{M}} (\varepsilon_{actor}^{m})^{2} + \max_{m \in \mathcal{M}} \frac{2R^{2}}{(1-\gamma^{m})^{4}}\right). \tag{22}$$

The proof of Theorem 2 can be found in Appendix A.6. Theorem 2 shows that Algorithm 3 can converge to an approximate Pareto-stationary solution at a rate of  $\mathcal{O}(1/T)$  with some approximate error terms. These errors are not significant, as we can control the upper bound of their upper bounds by setting the feature vector in the linear approximation, sample batch B, sample size H for policy gradient approximation, and the  $K_{\lambda}$  for the calculation of objective weights.

#### 5 SIMULATION EXPERIMENTS

In this section, we employ a path planning for multiple robotics to travel environment analogous to the one described in (Zhou et al., 2023). While (Zhou et al., 2023) examines a road network comprising 8 nodes, our simulation focuses on a larger network containing 18 nodes. Specifically, the path planning problems of N robots (i.e., agents) on a typical acyclic path network in Fig. 1, where the "blue" nodes in  $\{b_1, b_2, \cdots, b_5\}$  represent the set of starting nodes for agents. "purple" node, "orange" node, and "green" node in sthe right-hand-side represent the destination nodes of objective 1, objective 2, and objective 3, respectively. In the path planning problem, the local state space  $\mathcal{S}_i$  of agent  $i \in \mathcal{N}$  is defined as  $\mathcal{S}_i = \{b_1, b_2, b_3, b_4, b_5, c_1, c_2, c_3, c_4, c_5, d_1, d_2, d_3, d_4, d_5, e_1, e_2, e_3\}$ , and the local action space of agent  $i \in \mathcal{N}$  is defined as  $\mathcal{A}_i = \{0, 1, 2, 3\}$  with 3 being the maximum out degree of nodes in path network

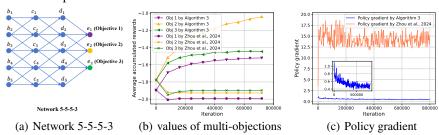


Figure 1: (a) Acyclic network. (b,c) The evolution of the objective performance  $J(\theta_t)$  and the norm of policy gradient  $||g_t||_2$  of the policy sequence generated by Algorithm 3, respectively.

In Network 5-5-5-3, consider agent i at node  $b_2$  for illustration, the action "0" means that the agent i remains stationary at the current node for one time step, while "1", "2", and "3" indicate that the agent i follows the edge  $(b_2, c_1)$ , edge  $(b_2, c_2)$ , and edge  $(b_2, c_3)$ , respectively. It should be noted that the agent i remains at the current node even when the action selected by it exceeds the out degree of the current node. For example, if agent  $j \in \mathcal{N}$  selects the action "3" at node  $b_1$ , then it will remain stationary at  $b_1$  for one time step.

The reward settings of each agent include: (i) the time run cost -0.5 at each step; (ii) the collision penalty when the agents share a path to move; (iii) the additional rewards for achieving different objectives. Specifically, when an agent reaches objective 1, objective 2, and objective 3, it receives the additional rewards of 0.5, 1.5, and 1, respectively. In this path planning problem, the agents want to efficiently reach the destinations while avoiding collision. The objective of agents is to find a joint policy parameter  $\theta$  to maximizes (3).

Our robot path planning problem includes 10 agents, whose initial positions are set to  $\{b_1,b_2,b_3,b_4,b_5,b_1,b_2,b_3,b_4,b_5\}$ . In this simulation, both the proposed Algorithm 3 and the latest MORL centralized algorithm from (Zhou et al., 2024) are tested under different random seeds. The discounted average cumulative reward  $\{J^m(\theta_t)\}_{m\in\{1,2,3\}}$  of the policy sequence generated by Algorithm 3 and the centralized algorithm are depicted in Fig. 1(b), where Algorithm 3 outperforms the centralized algorithm on performance of each objective. The main reason is that the centralized algorithm applies linear approximation to the global state-action space, causing large errors, whereas our distributed algorithm only approximates local Q-functions, leading to higher accuracy.

Moreover, the norm of policy gradients (i.e.,  $\|g_t\|_2$ ) generated by Algorithm 3 and the centralized algorithm are showed in Fig. 1(c). The norm of the policy gradient in Algorithm 3 exhibits a fast convergence trend towards to 0. However, the policy gradient in the centralized algorithm does not converge but fluctuates over many iterations from 0 due to the excessively large state-action dimension, resulting in a substantial approximation error in the linear approximation.

# 6 CONCLUSIONS

We propose a distributed algorithm for MOMARL and prove its convergence to a close-to-Pareto-stationary point. Each agent only requires state-action information  $(s_{\mathcal{N}_i^\kappa}, a_i)$ , ensuring scalability. This framework itself is a significant contribution and may inspire other scalable RL methods in networked systems

#### REFERENCES

- Rick Zhang and Marco Pavone. Control of robotic mobility-on-demand systems: a queueing-theoretical perspective. *The International Journal of Robotics Research*, 35(1-3): 186-203, 2016.
- Marina Weck, Eric Blake Jackson, Markus Sihvonen, and Ingrid Pappel. Building smart living environments for ageing societies: Decision support for cross-border e-services between Estonia and Finland. *Technology in Society*, 71: 102066, 2022.
  - Jiayi Xu, Mario Di Nardo, and Shi Yin. Improved swarm intelligence-based logistics distribution optimizer: decision support for multimodal transportation of cross-border e-commerce. *Mathe-matics*, 12(5): 763, 2024.
  - Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 1998.
  - Ivo Grondman, Lucian Buşoniu, Gabriel A. D. Lopes, and Robert Babuška. A survey of actor-critic reinforcement learning: Standard and natural policy gradients. *IEEE Transactions on Systems, Man, and Cybernetics, part C (applications and reviews)*, 42(6): 1291-1307, 2012.
  - Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, 321-384, 2021.
  - Yingqiang Ge, Xiaoting Zhao, Lucia Yu, Saurabh Paul, Diane Hu, Chu-Cheng Hsieh, and Yongfeng Zhang. Toward pareto efficient fairness-utility tradeoff in recommendation through reinforcement learning. In *Proceedings of the ACM international conference on web search and data mining*, ppages 316-324, 2022.
  - Dusan Stamenkovic, Alexandros Karatzoglou, Ioannis Arapakis, Xin Xin, and Kleomenis Katevas. Choosing the best of both worlds: Diverse and novel recommendations through multi-objective reinforcement learning. In *Proceedings of the ACM International Conference on Web Search and Data Mining*, pages 957-965, 2022
  - Maude J. Blondin and Matthew Hale. An algorithm for multi-objective multi-agent optimization. In *Proceedings of the IEEE American Control Conference*, pages 1489-1494, 2020.
  - Haibo Yang, Zhuqing Liu, Jia Liu, Chaosheng Dong, Michinari Momma. Federated multi-objective learning. In *Advances in Neural Information Processing Systems*, 36, 2024.
  - Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In *Advances in Neural Information Processing Systems*, 31, 2018.
  - Xu Chen, Yali Du, Long Xia, and Jun Wang. Reinforcement recommendation with user multi-aspect preference. In *Proceedings of the Web Conference*, pages 425-435, 2021.
  - David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *Proceedings of the International Conference on Machine Learning*, pages 387-395, 2014.
  - Tianchen Zhou, FNU Hairi, Haibo Yang, Jia Liu, Tian Tong, Fan Yang, Michinari Momma, and Yan Gao. Finite-time convergence and sample complexity of actor-critic multi-objective reinforcement learning. *arXiv preprint*, arXiv: 2405.03082, 2024.
  - Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization. *Comptes Rendus Mathematique*, 350(5-6): 313-318, 2012.
  - Roxana Rădulescu, Patrick Mannion, Diederik M. Roijers, and Ann Nowé. Multi-objective multi-agent decision making: a utility-based analysis and survey. *Autonomous Agents and Multi-Agent Systems*, 34(1): 10, 2020.
  - Yiheng Lin, Guannan Qu, Longbo Huang, and Adam Wierman. Multi-agent reinforcement learning in stochastic networked systems. In *Advances in Neural Information Processing Systems*, pages 7825-7837, 2021.
  - Zhaoyi Zhou, Zaiwei Chen, Yiheng Lin, and Adam Wierman. Convergence rates for localized actorcritic in networked markov potential games. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, pages 2563-2573, 2023.

- Tianshu Chu, Jie Wang, Lara Codecà, and Zhaojian Li. Multi-agent deep reinforcement learning for large-scale traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 21(3): 1086-1095, 2020.
- Pengcheng Dai, Wenwu Yu, He Wang, and Jiahui Jiang. Applications in traffic signal control: a distributed policy gradient decomposition algorithm. *IEEE Transactions on Industrial Informatics*, 20(2): 2762-2775, 2024.
- Runyu Zhang, Jincheng Mei, Bo Dai, Dale Schuurmans, and Na Li. On the global convergence rates of decentralized softmax gradient play in markov potential games. In *Advances in Neural Information Processing Systems*, pages 1923-1935, 2022.
- Harshat Kumar, Alec Koppel, and Alejandro Ribeiro. On the sample complexity of actor-critic for reinforcement learning. In *Advances in Neural Information Processing Systems*, 2019.
- Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, pages 1057-1063, 2000.
- Guannan Qu, Adam Wierman, and Na Li. Scalable reinforcement learning of localized policies for multi-agent networked systems. In *Proceedings of the Conference on Learning for Dynamics and Control*, pages 256-266, 2020.
- Guannan Qu, Yiheng Lin, Adam Wierman, and Na Li. Scalable Multi-Agent Reinforcement Learning for Networked Systems with Average Reward. In Advances in Neural Information Processing Systems, pages 2074-2086, 2020.
- John Tsitsiklis and Benjamin Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5): 674-690, 1997.
- R. Olfati-Saber and R. M. Murray, "Consensus problems in networks of agents with switching topology and time-delays," *IEEE Trans. Autom. Control*, vol. 49, no. 9, pp. 1520-1533, Sep. 2004.

# A APPENDIX

#### A.1 Preliminary definitions for Theorem 1

In this subsection, we introduce the formal definition of the exponential decay property in the MO-MARL problem.

Different from the definition of the global Q-function in (4), for each agent  $i \in \mathcal{N}$ , its local Q-function  $Q_i^m(s, a; \theta)$  in m-th objective is defined as

$$Q_i^m(s, \boldsymbol{a}; \boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\pi}_{\boldsymbol{\theta}}} \Big[ \sum_{t=0}^{\infty} (\gamma^m)^t r_{i,t}^m | \boldsymbol{s}_0 = \boldsymbol{s}, \boldsymbol{a}_0 = \boldsymbol{a} \Big].$$
 (23)

Based on the definitions of the global Q-function (4) and the local Q-function (23), we have

$$Q^{m}(s, \boldsymbol{a}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} Q_{i}^{m}(s, \boldsymbol{a}; \boldsymbol{\theta}), \tag{24}$$

which shows the global Q-function can be decomposed into the sum of the local Q-functions of all agents.

**Definition 3** The MOMARL satisfies the  $(\vartheta, \varrho)$ -exponential decay property with  $\vartheta = (\vartheta^1, \dots, \vartheta^M)^\top \in \mathbb{R}^M, \varrho = (\varrho^1, \dots, \varrho^M)^\top \in \mathbb{R}^M$ , if for any joint policy  $\pi_{\theta}$ , agent  $i \in \mathcal{N}$ , objective  $m \in \mathcal{M}$ ,  $s_{\mathcal{N}_i^{\kappa}} \in \mathcal{S}_{\mathcal{N}_i^{\kappa}}$ ,  $a_{\mathcal{N}_i^{\kappa}} \in \mathcal{A}_{\mathcal{N}_i^{\kappa}}$ ,  $s_{\mathcal{N}_{-i}^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}} \in \mathcal{S}_{\mathcal{N}_{-i}^{\kappa}}$ , and  $a_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}} \in \mathcal{A}_{\mathcal{N}_{-i}^{\kappa}}$ ,  $Q_i^{m}(s, \boldsymbol{a}; \boldsymbol{\theta})$  satisfies

$$\left| Q_i^m(s_{\mathcal{N}_i^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_i^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}}; \boldsymbol{\theta}) - Q_i^m(s_{\mathcal{N}_i^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}}', a_{\mathcal{N}_i^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}}'; \boldsymbol{\theta}) \right| \le \vartheta^m(\varrho^m)^{\kappa + 1}. \tag{25}$$

The exponential decay property of the MOMARL problem indicates that the dependence of agent i's local Q-function  $Q_i^m(s,a;\theta)$  on other agents shrinks rapidly as the distance between them increases. By Assumption 2, we can directly obtain the following lemma.

**Lemma 4** The MOMARL problem satisfies  $\left(\left(\frac{R}{1-\gamma^1},\cdots,\frac{R}{1-\gamma^M}\right)^\top,\gamma\right)$ -exponential decay property.

**Proof.** For any objective  $m \in \mathcal{M}$  and agent  $i \in \mathcal{N}$ , by using Lemma 3 in (Qu et al., 2020a), we have that

$$\left|Q_i^m(s_{\mathcal{N}_i^{\kappa}},s_{\mathcal{N}_{-i}^{\kappa}},a_{\mathcal{N}_i^{\kappa}},a_{\mathcal{N}_{-i}^{\kappa}};\boldsymbol{\theta}) - Q_i^m(s_{\mathcal{N}_i^{\kappa}},s_{\mathcal{N}_{-i}^{\kappa}}',a_{\mathcal{N}_i^{\kappa}},a_{\mathcal{N}_{-i}^{\kappa}}';\boldsymbol{\theta})\right| \leq \frac{R}{1-\gamma^m}(\gamma^m)^{\kappa+1},$$

which can further deduce that the MOMARL problem satisfies the  $\left(\left(\frac{R}{1-\gamma^1},\cdots,\frac{R}{1-\gamma^M}\right)^\top,\gamma\right)$ -exponential decay property.

Lemma 4 provides a possibility for agents to approximate  $Q_i^m(s, a; \theta)$  by only using its  $\kappa$ -hop neighbors' information. Inspired by exponential decay property in Lemma 4, we design a proper class of graph-truncated Q-functions:

$$Q_{tru,i}^{m}(s_{\mathcal{N}_{i}^{\kappa}}, a_{\mathcal{N}_{i}^{\kappa}}; \boldsymbol{\theta}) = \sum_{s_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}}} \xi_{\boldsymbol{\rho}}^{\boldsymbol{\theta}, m}(s_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}} | s_{\mathcal{N}_{i}^{\kappa}}, a_{\mathcal{N}_{i}^{\kappa}})$$

$$Q_{i}^{m}(s_{\mathcal{N}_{i}^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}}; \boldsymbol{\theta}), \tag{26}$$

where  $\xi^{m{ heta},m}_{m{
ho}}(s_{\mathcal{N}^{\kappa}_{-i}},a_{\mathcal{N}^{\kappa}_{-i}}|s_{\mathcal{N}^{\kappa}_{i}},a_{\mathcal{N}^{\kappa}_{i}})$  is the weight coefficient and satisfies

$$\xi_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}}|s_{\mathcal{N}_{i}^{\kappa}}, a_{\mathcal{N}_{i}^{\kappa}}) \\
= \frac{\xi_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s_{\mathcal{N}_{i}^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}}, a_{\mathcal{N}_{i}^{\kappa}}, a_{\mathcal{N}_{-i}^{\kappa}})}{\sum_{s_{\mathcal{N}_{-i}^{\kappa}}', a_{\mathcal{N}_{-i}^{\kappa}}'} \xi_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s_{\mathcal{N}_{i}^{\kappa}}, s_{\mathcal{N}_{-i}^{\kappa}}', a_{\mathcal{N}_{i}^{\kappa}}, a_{\mathcal{N}_{i}^{\kappa}}', a_{\mathcal{N}_{-i}^{\kappa}}')}.$$
(27)

Using (26), we define the graph-truncated policy gradient  $\nabla_{\theta_i} J_{tru,i}^m(\theta)$  as

$$\nabla_{\theta_{i}} J_{tru,i}^{m}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}, \boldsymbol{a} \sim \boldsymbol{\pi}_{\boldsymbol{\theta}}} \left[ \sum_{j \in \mathcal{N}_{i}^{\kappa}} Q_{tru,j}^{m}(s_{\mathcal{N}_{j}^{\kappa}}, a_{\mathcal{N}_{j}^{\kappa}}; \boldsymbol{\theta}) \right]$$

$$\nabla_{\theta_{i}} \log \pi_{\theta_{i}}(a_{i}|s_{i}) \frac{1}{(1-\gamma)N}. \tag{28}$$

The graph-truncated policy gradient approximation error is presented in the following.

**Lemma 5** In the MOMARL problem, for any agent  $i \in \mathcal{N}$  and objective  $m \in \mathcal{M}$ , we have

$$\left\| \nabla_{\theta_i} J_{tru,i}^m(\boldsymbol{\theta}) - \nabla_{\theta_i} J^m(\boldsymbol{\theta}) \right\|_2 \le \frac{\sqrt{2}R}{(1 - \gamma^m)^2} (\gamma^m)^{\kappa + 1}. \tag{29}$$

The proof of Lemma 5 is similar to Lemma 4 in (Qu et al., 2020a), therefore omitted here. Lemma 5 shows that the graph-truncated Q-functions  $\{Q^m_{tru,j}(s_{\mathcal{N}^\kappa_j},a_{\mathcal{N}^\kappa_j};\boldsymbol{\theta})\}_{j\in\mathcal{N}^\kappa_i}$  can effectively approximate the policy gradient  $\nabla_{\theta_i}J^m(\boldsymbol{\theta})$  through the state-action  $(s_{\mathcal{N}^\kappa_i},a_{\mathcal{N}^\kappa_i})$ .

Unlike the graph-truncated policy gradient  $\nabla_{\theta_i} J^m_{tru,i}(\theta)$  in (28) that requires  $a_{\mathcal{N}^\kappa_i}$ , (8) only requires the local action  $a_i$ . Next, we establish the equivalence between graph-truncated policy gradient  $\nabla_{\theta_i} J^m_{tru,i}(\theta)$  and approximated policy gradient  $\nabla_{\theta_i} J^m_{app}(\theta)$  in the following proposition.

**Proposition 2** *In the MOMARL problem, given a joint policy*  $\pi_{\theta}$ *, for any agent*  $i \in \mathcal{N}$  *and objective*  $m \in \mathcal{M}$ *, it holds* 

$$\nabla_{\theta_i} J_{tru,i}^m(\boldsymbol{\theta}) = \nabla_{\theta_i} J_{app,i}^m(\boldsymbol{\theta}). \tag{30}$$

**Proof.** By the definition of  $Q_{tru}^m(s_{\mathcal{N}_i^{\kappa}}, a_{\mathcal{N}_i^{\kappa}}; \boldsymbol{\theta})$  in (26), we have

$$\mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}, m}, \boldsymbol{a} \sim \boldsymbol{\pi}_{\boldsymbol{\theta}}} \left[ \frac{1}{N} \sum_{j \in \mathcal{N}^{\kappa}} Q_{tru, j}^{m}(s_{\mathcal{N}_{j}^{\kappa}}, a_{\mathcal{N}_{j}^{\kappa}}; \boldsymbol{\theta}) \nabla_{\theta_{i}} \log \pi_{\theta_{i}}(a_{i} | s_{i}) \right]$$

$$= \mathbb{E}_{\boldsymbol{s} \sim d^{\boldsymbol{\theta},m}_{\boldsymbol{\rho}},\boldsymbol{a} \sim \boldsymbol{\pi}_{\boldsymbol{\theta}}} \Big[ \frac{1}{N} \sum_{j \in \mathcal{N}^{\kappa}_{i}} \sum_{\tilde{\boldsymbol{s}}_{-\mathcal{N}^{\kappa}_{i}}, \tilde{\boldsymbol{a}}_{-\mathcal{N}^{\kappa}_{j}}} \xi^{\boldsymbol{\theta},m}_{\boldsymbol{\rho}}(\tilde{\boldsymbol{s}}_{-\mathcal{N}^{\kappa}_{j}}, \tilde{\boldsymbol{a}}_{-\mathcal{N}^{\kappa}_{j}} | s_{\mathcal{N}^{\kappa}_{j}}, a_{i}, a_{\mathcal{U}^{\kappa}_{j,-i}}) Q^{m}_{j}(s_{\mathcal{N}^{\kappa}_{j}}, \tilde{\boldsymbol{s}}_{-\mathcal{N}^{\kappa}_{j}}, a_{i}, a_{\mathcal{U}^{\kappa}_{j,-i}}, \tilde{\boldsymbol{a}}_{-\mathcal{N}^{\kappa}_{j}}; \boldsymbol{\theta}) \\$$

$$\nabla_{\theta_i} \log \pi_{\theta_i}(a_i|s_i) \Big]$$

$$= \mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}, m}, \boldsymbol{a} \sim \boldsymbol{\pi}_{\boldsymbol{\theta}}} \left[ \frac{1}{N} \sum_{j \in \mathcal{N}_{\kappa}^{\kappa}} Q_{j}^{m}(s_{\mathcal{N}_{j}^{\kappa}}, s_{-\mathcal{N}_{j}^{\kappa}}, a_{i}, a_{\mathcal{U}_{j,-i}^{\kappa}}, a_{-\mathcal{N}_{j}^{\kappa}}; \boldsymbol{\theta}) \nabla_{\theta_{i}} \log \pi_{\theta_{i}}(a_{i}|s_{i}) \right]$$
(31)

$$= \mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}, a_{i} \sim \pi_{\theta_{i}}} \left[ \frac{1}{N} \mathbb{E}_{\boldsymbol{\pi}\boldsymbol{\theta}} \left[ \sum_{t=0}^{\infty} (\gamma^{m})^{t} \sum_{j \in \mathcal{N}_{i}^{\kappa}} r_{j}^{m}(s_{j,t}, a_{j,t}) | \boldsymbol{s}_{0} = \boldsymbol{s}, a_{i,0} = a_{i} \right] \nabla_{\theta_{i}} \log \pi_{\theta_{i}}(a_{i}|s_{i}) \right]$$
(32)

$$= \mathbb{E}_{\boldsymbol{s} \sim d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}, m}, a_{i} \sim \pi_{\theta_{i}}} \left[ \widehat{Q_{i}^{m}}(\boldsymbol{s}, a_{i}; \boldsymbol{\theta}) \nabla_{\theta_{i}} \log \pi_{\theta_{i}}(a_{i} | s_{i}) \right], \tag{33}$$

where the second equality (31) is obtained from the definition of  $\xi_{\rho}^{\theta,m}(s_{\mathcal{N}_{-i}^{\kappa}},a_{\mathcal{N}_{-i}^{\kappa}}|s_{\mathcal{N}_{i}^{\kappa}},a_{\mathcal{N}_{i}^{\kappa}})$  in (27), the third equality (32) comes from the definition of the local Q-function in (23), and the last equality (33) can be achieved by the definition of  $\widehat{Q_{i}^{m}}(s,a_{i};\theta)$  in (7).

Proposition 2 provides an equivalence between  $Q^m_{tru,i}(s_{\mathcal{N}_i^{\kappa}},a_{\mathcal{N}_j^{\kappa}};\boldsymbol{\theta})$  and  $\widehat{Q^m_i}(s,a_i;\boldsymbol{\theta})$  in policy gradient approximation.

#### A.2 The proof of Theorem 1

**Proof.** By the definition of  $\nabla_{\theta_i} J^m_{app,i}(\boldsymbol{\theta})$  in (8), we have

$$\|\nabla_{\theta_{i}}J_{app,i}^{m}(\boldsymbol{\theta}) - \nabla_{\theta_{i}}J^{m}(\boldsymbol{\theta})\|_{2} = \|\nabla_{\theta_{i}}J_{app,i}^{m}(\boldsymbol{\theta}) - \nabla_{\theta_{i}}J_{tru,i}^{m}(\boldsymbol{\theta}) + \nabla_{\theta_{i}}J_{tru,i}^{m}(\boldsymbol{\theta}) - \nabla_{\theta_{i}}J^{m}(\boldsymbol{\theta})\|_{2}$$

$$= \|\nabla_{\theta_{i}}J_{tru,i}^{m}(\boldsymbol{\theta}) - \nabla_{\theta_{i}}J^{m}(\boldsymbol{\theta})\|_{2}$$
(34)

$$\leq \frac{\sqrt{2}R}{(1-\gamma^m)^2} (\gamma^m)^{\kappa+1},\tag{35}$$

where the second equality comes from Proposition 2 and last inequality achieved by Lemma 5. Hence, the proof is completed.  $\Box$ 

## A.3 PROOF OF LEMMA 2

In the MOMARL problem, for any  $i \in \mathcal{N}$ , denote  $s_{-i} = s \setminus s_i$  as the state of agents other than agent i and  $a_{-i} = a \setminus a_i$  as the action of agents other than agent i. For any joint policy  $\pi_{\theta}$ , denote  $d_{\rho,i}^{\theta,m}(s_i)$  and  $d_{\rho,-i}^{\theta,m}(s_{-i})$  as the discounted state visitation distribution of  $s_i$  and  $s_{-i}$  in m-objective, respectively. For each agent  $i \in \mathcal{N}$ , define its local value function  $V_i^m(s;\theta)$  in m-th objective as

$$V_i^m(s;\theta) = \sum_{a} \pi_{\theta}(a|s) Q_i^m(s,a;\theta).$$
 (36)

Define the averaged value function, the averaged Q-function, and the averaged advantage function of agent  $i \in \mathcal{N}$  in the objective  $m \in \mathcal{M}$  as

$$\overline{V_i^m}(s_i; \boldsymbol{\theta}) = \frac{1}{N} \sum_{s'_{-i}} d_{\boldsymbol{\rho}, -i}^{\boldsymbol{\theta}, m}(s'_{-i}) \sum_{j \in \mathcal{N}} V_j^m(s_i, s'_{-i}; \boldsymbol{\theta}), \tag{37}$$

$$\overline{Q_i^m}(s_i, a_i; \boldsymbol{\theta}) = \frac{1}{N} \sum_{s'_{-i}, a'_{-i}} d_{\boldsymbol{\rho}, -i}^{\boldsymbol{\theta}, m}(s'_{-i}) \pi_{\boldsymbol{\theta}_{-i}}(a'_{-i}|s'_{-i})$$

$$\sum_{j \in \mathcal{N}} Q_j^m(s_i, s'_{-i}, a_i, a'_{-i}; \boldsymbol{\theta}), \tag{38}$$

$$\overline{A_i^m}(s_i, a_i; \boldsymbol{\theta}) = \overline{Q_i^m}(s_i, a_i; \boldsymbol{\theta}) - \overline{V_i^m}(s_i; \boldsymbol{\theta}). \tag{39}$$

**Lemma 6** (Softmax policy gradient) In the MOMARL problem, for any joint policy  $\pi_{\theta}$ , agent  $i \in \mathcal{N}$ , and objective  $m \in \mathcal{M}$ , the gradient of  $J^m(\theta)$  with respect to  $\theta_{i,s_i,a_i}$  is represented as

$$\frac{\partial J^m(\boldsymbol{\theta})}{\partial \theta_{i,s_i,a_i}} = \frac{1}{1 - \gamma^m} d_{\boldsymbol{\rho},i}^{\boldsymbol{\theta},m}(s_i) \pi_{\theta_i}(a_i|s_i) \overline{A_i^m}(s_i, a_i; \boldsymbol{\theta}). \tag{40}$$

**Proof.** According to the policy gradient lemma 1 and (24), we have

$$\frac{\partial J^{m}(\boldsymbol{\theta})}{\partial \theta_{i,s_{i},a_{i}}} = \frac{1}{1 - \gamma^{m}} \sum_{s',a'} d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}}(s') \pi_{\boldsymbol{\theta}}(a'|s') \frac{\partial \log \pi_{\boldsymbol{\theta}}(a'|s')}{\partial \theta_{i,s_{i},a_{i}}} \\
\left(\frac{1}{N} \sum_{j=1}^{N} Q_{j}^{m}(s', a'; \boldsymbol{\theta})\right) \\
= \frac{1}{1 - \gamma^{m}} \sum_{s'_{i},a'_{i}} d_{\boldsymbol{\rho},n}^{\boldsymbol{\theta},m}(s'_{i}) \pi_{\boldsymbol{\theta}_{i}}(a'_{i}|s'_{i}) \\
\sum_{s'_{-i},a'_{-i}} d_{\boldsymbol{\rho},-i}^{\boldsymbol{\theta},m}(s'_{-i}) \pi_{\boldsymbol{\theta}_{-i}}(a'_{-i}|s'_{-i}) \left(\mathbf{1}\{s'_{i} = s_{i}, a'_{i} = a_{i}\}\right) \\
-\mathbf{1}\{s'_{i} = s_{i}\} \pi_{\boldsymbol{\theta}_{i}}(a_{i}|s_{i})\right) \frac{1}{N} \left(\sum_{j \in \mathcal{N}} Q_{j}^{m}(s_{i}, s'_{-i}, a_{i}, a'_{-i}; \boldsymbol{\theta})\right) \\
= \frac{1}{1 - \gamma^{m}} d_{\boldsymbol{\rho},n}^{\boldsymbol{\theta},m}(s_{i}) \pi_{\boldsymbol{\theta}_{i}}(a_{i}|s_{i}) \overline{Q_{i}^{m}}(s_{i}, a_{i}; \boldsymbol{\theta}) \\
-\frac{1}{1 - \gamma^{m}} d_{\boldsymbol{\rho},n}^{\boldsymbol{\theta},m}(s_{i}) \pi_{\boldsymbol{\theta}_{i}}(a_{i}|s_{i}) \overline{Q_{i}^{m}}(s_{i}, a_{i}; \boldsymbol{\theta}) - \overline{V_{i}^{m}}(s_{i}; a_{i}; \boldsymbol{\theta}) \\
= \frac{1}{1 - \gamma^{m}} d_{\boldsymbol{\rho},n}^{\boldsymbol{\theta},m}(s_{i}) \pi_{\boldsymbol{\theta}_{i}}(a_{i}|s_{i}) \overline{Q_{i}^{m}}(s_{i}, a_{i}; \boldsymbol{\theta}) - \overline{V_{i}^{m}}(s_{i}; \boldsymbol{\theta}) \\
= \frac{1}{1 - \gamma^{m}} d_{\boldsymbol{\rho},n}^{\boldsymbol{\theta},m}(s_{i}) \pi_{\boldsymbol{\theta}_{i}}(a_{i}|s_{i}) \overline{A_{i}^{m}}(s_{i}, a_{i}; \boldsymbol{\theta}), \tag{42}$$

where the second equality (41) comes from the fact that

$$\frac{\partial \log \boldsymbol{\pi}_{\boldsymbol{\theta}}(\boldsymbol{a}'|\boldsymbol{s}')}{\partial \theta_{i,s_{i},a_{i}}} = \frac{\partial \log \boldsymbol{\pi}_{\theta_{i}}(a'_{i}|s'_{i})}{\partial \theta_{i,s_{i},a_{i}}}$$

$$= \mathbf{1}\{s'_{i} = s_{i}, a'_{i} = a_{i}\} - \mathbf{1}\{s'_{i} = s_{i}\}\boldsymbol{\pi}_{\theta_{i}}(a_{i}|s_{i}), \tag{43}$$

and the last equality (42) can be obtained from the definition of the averaged advantage function  $\overline{A_i^m}(s_i,a_i;\theta)$ .  $\square$  Based on Lemma 6, the proof of Lemma 2 are represented as follows.

Consider that for any different policies  $\pi_{\theta}$  and  $\pi_{\theta'}$ , we have

$$\|\nabla_{\boldsymbol{\theta}} J^m(\boldsymbol{\theta}') - \nabla_{\boldsymbol{\theta}} J^m(\boldsymbol{\theta})\|_2^2$$

$$\leq \sum_{i=1}^{N} \|\nabla_{\theta_i} J^m(\boldsymbol{\theta}') - \nabla_{\theta_i} J^m(\boldsymbol{\theta})\|_2^2$$

$$\leq \sum_{i=1}^{N} \|\nabla_{\theta_i} J^m(\boldsymbol{\theta}') - \nabla_{\theta_i} J^m(\boldsymbol{\theta})\|_1^2. \tag{44}$$

By Lemmas 6, we have

$$\|\nabla_{\theta_{i}}J^{m}(\boldsymbol{\theta}') - \nabla_{\theta_{i}}J^{m}(\boldsymbol{\theta})\|_{1}$$

$$= \frac{1}{1 - \gamma^{m}} \sum_{s_{i},a_{i}} \left| d_{\boldsymbol{\rho},i}^{\boldsymbol{\theta}',m}(s_{i})\pi_{\theta_{i}'}(a_{i}|s_{i})\overline{A_{i}^{m}}(s_{i},a_{i};\boldsymbol{\theta}') - d_{\boldsymbol{\rho},i}^{\boldsymbol{\theta},m}(s_{i})\pi_{\theta_{i}}(a_{i}|s_{i})\overline{A_{i}^{m}}(s_{i},a_{i};\boldsymbol{\theta}) \right|$$

$$\leq \frac{1}{1 - \gamma^{m}} \sum_{s,a} \left| d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}',m}(s)\pi_{\boldsymbol{\theta}'}(\boldsymbol{a}|s)A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}') - d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s)\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|s)A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}) \right|$$

$$\leq \frac{1}{1 - \gamma^{m}} \sum_{s,a} \left( \left| d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}',m}(s)\pi_{\boldsymbol{\theta}'}(\boldsymbol{a}|s) - d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s)\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|s) \right| \right) A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}') + d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s)\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|s) |A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}') - A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}) |$$

$$\leq \frac{1}{1 - \gamma^{m}} \sum_{s,a} \frac{1}{1 - \gamma^{m}} \left( \left| d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}',m}(s)\pi_{\boldsymbol{\theta}'}(\boldsymbol{a}|s) - d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s)\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|s) \right| \right) + \max_{s,a} |A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}') - A^{m}(s,\boldsymbol{a};\boldsymbol{\theta})|,$$

$$\leq \frac{1}{1 - \gamma^{m}} \sum_{s,a} \frac{1}{1 - \gamma^{m}} \left( \left| d_{\boldsymbol{\rho}}^{\boldsymbol{\theta}',m}(s)\pi_{\boldsymbol{\theta}'}(\boldsymbol{a}|s) - d_{\boldsymbol{\rho}}^{\boldsymbol{\theta},m}(s)\pi_{\boldsymbol{\theta}}(\boldsymbol{a}|s) \right| \right) + \max_{s,a} |A^{m}(s,\boldsymbol{a};\boldsymbol{\theta}') - A^{m}(s,\boldsymbol{a};\boldsymbol{\theta})|,$$

$$(46)$$

where the first inequality (45) can be obtained from the definition of  $\overline{A_i^m}(s_i,a_i;\theta)$  in (39) and the fact that  $|\sum_{i=1}^N x_i - \sum_{i=1}^N y_i| \le \sum_{i=1}^N |x_i - y_i|, \forall x_i,y_i \in \mathbb{R}$ , and the last inequality (46) comes from the fact that  $A^m(s,a;\theta) \le 1/(1-\gamma^m)$ .

For the right side of (46), we can use Corollary 35 and Lemma 32 in Zhang et al. (2022) to further obtain

$$\sum_{s} |d_{\rho}^{\theta',m}(s)\pi_{\theta'}(a|s) - d_{\rho}^{\theta,m}(s)\pi_{\theta}(a|s)|$$

$$\leq \frac{1}{1 - \gamma^{m}} \max_{s} \|\pi_{\theta'}(\cdot|s) - \pi_{\theta}(\cdot|s)\|_{1}$$
(47)

and

$$|A^{m}(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{\theta}') - A^{m}(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{\theta})|$$

$$\leq \frac{2}{(1 - \gamma^{m})^{2}} \max_{\boldsymbol{s}} \|\boldsymbol{\pi}_{\boldsymbol{\theta}'}(\cdot|\boldsymbol{s}) - \boldsymbol{\pi}_{\boldsymbol{\theta}}(\cdot|\boldsymbol{s})\|_{1}.$$
(48)

Substituting (47) and (48) into (46), we have

$$\|\nabla_{\theta_{i}} J^{m}(\boldsymbol{\theta}') - \nabla_{\theta_{i}} J^{m}(\boldsymbol{\theta})\|_{1}$$

$$\leq \frac{3}{(1 - \gamma^{m})^{3}} \max_{\boldsymbol{s}} \|\boldsymbol{\pi}_{\boldsymbol{\theta}'}(\cdot|\boldsymbol{s}) - \boldsymbol{\pi}_{\boldsymbol{\theta}}(\cdot|\boldsymbol{s})\|_{1}$$

$$= \frac{3}{(1 - \gamma^{m})^{3}} \sum_{i=1}^{N} \max_{s_{i}} \|\boldsymbol{\pi}_{\theta'_{i}}(\cdot|s_{i}) - \boldsymbol{\pi}_{\theta_{i}}(\cdot|s_{i})\|_{1}$$

$$\leq \frac{6}{(1 - \gamma^{m})^{3}} \sum_{i=1}^{N} \|\theta'_{i} - \theta_{i}\|_{2},$$
(50)

where the last inequality (50) is obtained from Corollary 37 in Zhang et al. (2022) that for any two difference softmax policies  $\pi_{\theta_i}$  and  $\pi_{\theta_i'}$ , and  $s_i \in \mathcal{S}_i$ ,  $\|\pi_{\theta_i}(\cdot|s_i) - \pi_{\theta_i'}(\cdot|s_i)\|_1 \le 2\|\theta_i - \theta_i'\|_2$ .

Combining (44) and (50), we further have

$$\|\nabla_{\boldsymbol{\theta}} J^{m}(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} J^{m}(\boldsymbol{\theta}')\|_{2}^{2}$$

$$\leq \sum_{i=1}^{N} \left(\frac{6}{(1-\gamma^{m})^{3}} \sum_{i=1}^{N} \|\theta'_{i} - \theta_{i}\|_{2}\right)^{2}$$

$$= \frac{36N}{(1-\gamma^{m})^{6}} \left(\sum_{i=1}^{N} \|\theta'_{i} - \theta_{i}\|_{2}\right)^{2}$$

$$\leq \frac{36N^{2}}{(1-\gamma^{m})^{6}} \sum_{i=1}^{N} \|\theta'_{i} - \theta_{i}\|_{2}^{2}$$

$$\leq \frac{36N^{2}}{(1-\gamma^{m})^{6}} \|\theta' - \theta\|_{2}^{2}, \tag{51}$$

where the second inequality (51) is obtained from that  $\left(\sum_{i=1}^N x_i\right)^2 \leq N(\sum_{i=1}^N x_i^2), \forall x_i \in \mathbb{R}$ .

According to  $L_J = \max_{m \in \mathcal{M}} \frac{6N}{(1-\gamma^m)^3}$  and (52), we further have that  $J^m(\boldsymbol{\theta})$  is  $L_J$ -smooth for all  $m \in \mathcal{M}$ .

#### A.4 PROOF OF LEMMA 3

The detailed proof of Lemma 3 is provided in the following.

**Proof.** By the update of  $g_{i,t}^m(b+1)$  in (13) and  $g_{i,t}^m=g_{i,t}^m(B)$ , we have

$$\nabla_{\theta_{i}} J^{m}(\boldsymbol{\theta}_{t}) - g_{i,t}^{m}$$

$$= \nabla_{\theta_{i}} J^{m}(\boldsymbol{\theta}_{t}) - \nabla_{\theta_{i}} J_{app,i}^{m}(\boldsymbol{\theta}_{t}) + \nabla_{\theta_{i}} J_{app,i}^{m}(\boldsymbol{\theta}_{t}) - g_{i,t}^{m}$$

$$= \nabla_{\theta_{i}} J^{m}(\boldsymbol{\theta}_{t}) - \nabla_{\theta_{i}} J_{app,i}^{m}(\boldsymbol{\theta}_{t})$$

$$+ \sum_{h=0}^{\infty} (\gamma^{m})^{h} \mathbb{E} \left[ \nabla_{\theta_{i}} \log \pi_{\theta_{i,t}}(a_{i,h}|s_{i,h}) \widehat{Q}_{i}^{m}(s_{h}, a_{i,h}; \boldsymbol{\theta}_{t}) \right]$$

$$- \frac{1}{B} \sum_{b=0}^{B-1} \sum_{h=0}^{H-1} (\gamma^{m})^{h} \phi_{i}(s_{\mathcal{N}_{i}^{\kappa},h}^{b}, a_{i,h}^{b})^{\top} w_{i,t}$$

$$\nabla_{\theta_{i}} \log \pi_{\theta_{i,t}}(a_{i,h}^{b}|s_{i,h}^{b})$$

$$= \underbrace{\nabla_{\theta_{i}} J^{m}(\boldsymbol{\theta}_{t}) - \nabla_{\theta_{i}} J_{app,i}^{m}(\boldsymbol{\theta}_{t})}_{T_{1}}$$

$$+ \underbrace{\sum_{h=0}^{H-1} (\gamma^{m})^{h} \mathbb{E} \left[ \nabla_{\theta_{i}} \log \pi_{\theta_{i,t}}(a_{i,h}|s_{i,h}) \widehat{Q}_{i}^{m}(s_{h}, a_{i,h}; \boldsymbol{\theta}_{t}) \right]}_{T_{2}}$$

$$+ \underbrace{\sum_{h=H}^{\infty} (\gamma^{m})^{h} \mathbb{E} \left[ \nabla_{\theta_{i}} \log \pi_{\theta_{i,t}}(a_{i,h}|s_{i,h}) \widehat{Q}_{i}^{m}(s_{h}, a_{i,h}; \boldsymbol{\theta}_{t}) \right]}_{T_{3}}$$

$$- \underbrace{\frac{1}{B} \sum_{b=0}^{B-1} \sum_{h=0}^{H-1} (\gamma^{m})^{h} \nabla_{\theta_{i}} \log \pi_{\theta_{i,t}}(a_{i,h}^{b}|s_{i,h}) \widehat{Q}_{i}^{m}(s_{h}^{b}, a_{i,h}^{b}; \boldsymbol{\theta}_{t})}_{T_{4}}$$

$$+ \underbrace{\frac{1}{B} \sum_{b=0}^{B-1} \sum_{h=0}^{H-1} (\gamma^{m})^{h} \nabla_{\theta_{i}} \log \pi_{\theta_{i,t}}(a_{i,h}^{b}|s_{i,h}^{b})}_{T_{4}}$$

$$\underbrace{\left(\widehat{Q_i^m}(s_h^b, a_{i,h}^b; \boldsymbol{\theta}_t) - \phi_i(s_{\mathcal{N}_i^{\kappa}, h}^b, a_{i,h}^b)^{\top} w_{i,t}\right)}_{\mathcal{T}_{\epsilon}},$$
(54)

where the equality (53) can be obtained by the policy gradient theorem variant (i.e., Lemma F.1 in Zhou et al. (2023)). Based on (54), we have

$$\mathbb{E}[\|\nabla_{\theta_{i}}J^{m}(\theta_{t}) - g_{i,t}^{m}\|_{2}^{2}] \\
= \mathbb{E}[\|\mathcal{T}_{1} + \mathcal{T}_{2} + \mathcal{T}_{3} - \mathcal{T}_{4} + \mathcal{T}_{5}\|_{2}^{2}] \\
\leq 4\mathbb{E}[\|\mathcal{T}_{1}\|_{2}^{2} + \|\mathcal{T}_{2} - \mathcal{T}_{4}\|_{2}^{2} + \|\mathcal{T}_{3}\|_{2}^{2} + \|\mathcal{T}_{5}\|_{2}^{2}]. \tag{55}$$

$$\leq \frac{8R^{2}}{(1 - \gamma^{m})^{4}} (\gamma^{m})^{2\kappa + 2} + \frac{32}{(1 - \gamma^{m})^{2}B} + \frac{8(\gamma^{m})^{2H}}{(1 - \gamma^{m})^{4}} \\
+ \frac{8\varepsilon_{critic}^{\theta_{t}}}{(1 - \gamma^{m})^{2}}. \tag{56}$$

where (56) can be obtained by (35) and the definition of  $\varepsilon_{critic}^{\theta_t}$  in (20).

## A.5 PROOF OF 1

**Proof.** Based on (Olfati-Saber & Murray, 2004), the results of (i) and (ii) are straightforward; therefore, we will focus on proving (iii). Since  $J_t^g$  is  $L_t^g$ -Lipschitz continuous and  $\lambda_i(k) = \lambda_j(k)$  for all  $i, j \in \mathcal{N}$ , we define  $\lambda(k) = \lambda_i(k)$  and have

$$J_{t}^{g}(\boldsymbol{\lambda}(k+1)) \leq J_{t}^{g}(\boldsymbol{\lambda}(k)) + \alpha_{k} \langle \nabla J_{t}^{g}(\boldsymbol{\lambda}(k)), \boldsymbol{e}_{u_{i,t}^{*}(k)} - \boldsymbol{\lambda}(k) \rangle$$

$$+ \frac{L_{t}^{g}}{2} \alpha_{k}^{2} \|\boldsymbol{e}_{u_{i,t}^{*}(k)} - \boldsymbol{\lambda}(k)\|^{2}$$

$$\leq J_{t}^{g}(\boldsymbol{\lambda}(k)) + \alpha_{k} \langle \nabla J_{t}^{g}(\boldsymbol{\lambda}(k)), \boldsymbol{e}_{u_{i,t}^{*}(k)} - \boldsymbol{\lambda}(k) \rangle$$

$$+ \alpha_{k}^{2} L_{t}^{g}, \tag{57}$$

where the last inequality can be obtained by the fact that  $\|e_{u_{i,t}^*(k)} - \lambda(k)\| \le \sqrt{2}$ . Since  $u_{i,t}^*(k) = \arg\min_{m} y_{i,t}^m(k)$ , we can obtain that

$$\langle \nabla J_t^g (\boldsymbol{\lambda}(k)), \boldsymbol{e}_{u_{i,t}^*(k)} - \boldsymbol{\lambda}(k) \rangle \leq \langle \nabla J_t^g (\boldsymbol{\lambda}(k)), \boldsymbol{\lambda}_t^* - \boldsymbol{\lambda}(k) \rangle$$
$$\leq J_t^g (\boldsymbol{\lambda}_t^*) - f(\boldsymbol{\lambda}(k)). \tag{58}$$

Combining (58) with (57), we have

$$J_t^g(\boldsymbol{\lambda}(k+1)) - J_t^g(\boldsymbol{\lambda}_t^*) \le (1 - \alpha_k) \left( J_t^g(\boldsymbol{\lambda}(k)) - J_t^g(\boldsymbol{\lambda}_t^*) \right) + \alpha_k^2 L_t^g.$$
(59)

By induction, (59) implies that  $J_t^g(\lambda(k+1)) - J_t^g(\lambda_t^*) \leq \frac{4L_t^g}{k+1}$ . By the fact that  $\widehat{\lambda}_t = \lambda(K_\lambda)$ , we can obtain  $J_t^g(\widehat{\lambda}_t) - J_t^g(\lambda_t^*) \leq \frac{4L_t^g}{K_\lambda + 1}$ .

# A.6 PROOF OF THEOREM 2

**Proof.** According to the smoothness of  $J^m(\theta)$  in Lemma 2 and  $L_J = \max_{m \in \mathcal{M}} \frac{6N}{(1-\gamma^m)^3}$ , we can have

$$J^{m}(\boldsymbol{\theta}_{t+1}) \geq J^{m}(\boldsymbol{\theta}_{t}) + \langle \nabla_{\boldsymbol{\theta}} J^{m}(\boldsymbol{\theta}_{t}), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \rangle - \frac{L_{J}}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t}\|_{2}^{2}, \forall m \in \mathcal{M}.$$

$$(60)$$

Taking  $\lambda_t$  weighted summation over (60), we have

$$\begin{aligned} & \boldsymbol{\lambda}_t^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t+1}) \\ \geq & \boldsymbol{\lambda}_t^{\top} \boldsymbol{J}(\boldsymbol{\theta}_t) + \langle \nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_t)^{\top} \boldsymbol{\lambda}_t, \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle \\ & - \frac{L_J}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2^2 \\ = & \boldsymbol{\lambda}_t^{\top} \boldsymbol{J}(\boldsymbol{\theta}_t) + \eta_{\boldsymbol{\theta},t} \Big\langle \nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_t)^{\top} \boldsymbol{\lambda}_t, \sum_{t=1}^{M} \lambda_t^m \boldsymbol{g}_t^m \Big\rangle \end{aligned}$$

$$-\frac{L_{J}\eta_{\theta,t}^{2}}{2} \left\| \sum_{m=1}^{M} \lambda_{t}^{m} g_{t}^{m} \right\|_{2}^{2}$$

$$= \lambda_{t}^{\top} J(\theta_{t}) + \eta_{\theta,t} \left\langle \nabla_{\theta} J(\theta_{t})^{\top} \lambda_{t}, \right.$$

$$\left. \sum_{m=1}^{M} \lambda_{t}^{m} \left( g_{t}^{m} - \nabla_{\theta} J^{m}(\theta_{t}) + \nabla_{\theta} J^{m}(\theta_{t}) \right) \right\rangle$$

$$-\frac{L_{J}\eta_{\theta,t}^{2}}{2} \left\| \sum_{m=1}^{M} \lambda_{t}^{m} g_{t}^{m} \right\|_{2}^{2}$$

$$= \lambda_{t}^{\top} J(\theta_{t}) + \eta_{\theta,t} \left\langle \nabla_{\theta} J(\theta_{t})^{\top} \lambda_{t}, \sum_{m=1}^{M} \lambda_{t}^{m} \nabla_{\theta} J^{m}(\theta_{t}) \right\rangle$$

$$+ \eta_{\theta,t} \left\langle \nabla_{\theta} J(\theta_{t})^{\top} \lambda_{t}, \sum_{m=1}^{M} \lambda_{t}^{m} \left( g_{t}^{m} \right)$$

$$- \nabla_{\theta} J^{m}(\theta_{t}) \right) - \frac{L_{J}\eta_{\theta,t}^{2}}{2} \left\| \sum_{m=1}^{M} \lambda_{t}^{m} g_{t}^{m} \right\|_{2}^{2}$$

$$\geq \lambda_{t}^{\top} J(\theta_{t}) + \frac{\eta_{\theta,t}}{2} \left\| \nabla_{\theta} J(\theta_{t})^{\top} \lambda_{t} \right\|_{2}^{2}$$

$$- \frac{\eta_{\theta,t}}{2} \left\| \sum_{m=1}^{M} \lambda_{t}^{m} \left( g_{t}^{m} - \nabla_{\theta} J^{m}(\theta_{t}) \right) \right\|_{2}^{2}$$

$$- \frac{L_{J}\eta_{\theta,t}^{2}}{2} \left\| \sum_{m=1}^{M} \lambda_{t}^{m} \left( g_{t}^{m} - \nabla_{\theta} J^{m}(\theta_{t}) + \nabla_{\theta} J^{m}(\theta_{t}) \right) \right\|_{2}^{2}$$

$$\geq \lambda_{t}^{\top} J(\theta_{t}) + \left( \frac{\eta_{\theta,t}}{2} - L_{J}\eta_{\theta,t}^{2} \right) \left\| \nabla_{\theta} J(\theta_{t})^{\top} \lambda_{t} \right\|_{2}^{2}$$

$$- \left( \frac{\eta_{\theta,t}}{2} + L_{J}\eta_{\theta,t}^{2} \right) \left\| \sum_{m=1}^{M} \lambda_{t}^{m} \left( g_{t}^{m} - \nabla_{\theta} J^{m}(\theta_{t}) \right) \right\|_{2}^{2} ,$$

$$= \left( \frac{\eta_{\theta,t}}{2} + L_{J}\eta_{\theta,t}^{2} \right) \left\| \sum_{m=1}^{M} \lambda_{t}^{m} \left( g_{t}^{m} - \nabla_{\theta} J^{m}(\theta_{t}) \right) \right\|_{2}^{2} ,$$

$$(62)$$

where the equality (61) comes from (18), the inequality (62) can be obtained by the fact that  $\langle x, y \rangle \ge -\frac{1}{2}(\|x\|^2 + \|y\|^2), \forall x, y \in \mathbb{R}^{\sum_{i=1}^N |\mathcal{S}_i||\mathcal{A}_i|}$ , and the inequality (63) can be get by the fact that  $\|x + y\|_2^2 \le 2(\|x\|_2^2 + \|y\|_2^2), \forall x, y \in \mathbb{R}^{\sum_{i=1}^N |\mathcal{S}_i||\mathcal{A}_i|}$ . By (63), we have

$$\|\nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_{t})^{\top} \boldsymbol{\lambda}_{t}\|_{2}^{2} \leq \frac{2\left(\boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t+1}) - \boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t})\right)}{\eta_{\boldsymbol{\theta},t} - 2\eta_{\boldsymbol{\theta},t}^{2} L_{J}} + \frac{\eta_{\boldsymbol{\theta},t} + 2\eta_{\boldsymbol{\theta},t}^{2} L_{J}}{\eta_{\boldsymbol{\theta},t} - 2\eta_{\boldsymbol{\theta},t}^{2} L_{J}} \left\| \sum_{m=1}^{M} \lambda_{t}^{m} \left(\nabla_{\boldsymbol{\theta}} J^{m}(\boldsymbol{\theta}_{t}) - \boldsymbol{g}_{t}^{m}\right) \right\|_{2}^{2}.$$

$$(64)$$

Consider that  $\hat{\lambda}_t$  is the optimal of problem (15), we have

$$\|\nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_t)^{\top} \widehat{\boldsymbol{\lambda}}_t \|_2^2 \le \|\nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_t)^{\top} \boldsymbol{\lambda}_t \|_2^2 + \frac{4L_t^g}{K_{\lambda} + 1}.$$
 (65)

Using the setting of the learning rate as  $\eta_{\theta,t} = \frac{1}{3L_J}$  and taking expectation on both side of (64), we further have

$$\mathbb{E}[\|\nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_{t})^{\top} \widehat{\boldsymbol{\lambda}}_{t} \|_{2}^{2}]$$

$$\leq 18L_{J} \mathbb{E}[\boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t+1}) - \boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t})] + \mathbb{E}\left[\frac{4L_{t}^{g}}{K_{\lambda} + 1}\right]$$

$$+ 5\mathbb{E}\left[\sum_{m=1}^{M} \lambda_{t}^{m} \|\nabla_{\boldsymbol{\theta}} J^{m}(\boldsymbol{\theta}_{t}) - \boldsymbol{g}_{t}^{m} \|_{2}^{2}\right]$$

$$\leq 18L_{J}\mathbb{E}[\boldsymbol{\lambda}_{t}^{\top}\boldsymbol{J}(\boldsymbol{\theta}_{t+1}) - \boldsymbol{\lambda}_{t}^{\top}\boldsymbol{J}(\boldsymbol{\theta}_{t})] + 5 \max_{m \in \mathcal{M}} \varepsilon_{actor}^{m} + \frac{8}{K_{\lambda} + 1} \Big( \max_{m \in \mathcal{M}} (\varepsilon_{actor}^{m})^{2} + \max_{m \in \mathcal{M}} \frac{2R^{2}}{(1 - \gamma^{m})^{4}} \Big),$$
(66)

where the last inequality comes from Lemma 3. Taking average of (66) over T, we have

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla_{\boldsymbol{\theta}} \boldsymbol{J}(\boldsymbol{\theta}_{t})^{\top} \widehat{\boldsymbol{\lambda}}_{t}\|_{2}^{2}]$$

$$\leq \frac{1}{T} \sum_{t=1}^{T} 18L_{J} \mathbb{E}[\boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t+1}) - \boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t})]$$

$$+ 5 \max_{m \in \mathcal{M}} \varepsilon_{actor}^{m} + \frac{8}{K_{\lambda} + 1} \Big( \max_{m \in \mathcal{M}} (\varepsilon_{actor}^{m})^{2} + \max_{m \in \mathcal{M}} \frac{2R^{2}}{(1 - \gamma^{m})^{4}} \Big). \tag{67}$$

Considering that

$$\sum_{t=1}^{T} \mathbb{E}[\boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t+1}) - \boldsymbol{\lambda}_{t}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t})]$$

$$= \mathbb{E}\Big[\sum_{t=1}^{T-1} (-\boldsymbol{\lambda}_{t+1} + \boldsymbol{\lambda}_{t})^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{t+1}) - \boldsymbol{\lambda}_{1}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{1})$$

$$+ \boldsymbol{\lambda}_{T}^{\top} \boldsymbol{J}(\boldsymbol{\theta}_{T+1})\Big]$$

$$\leq \mathbb{E}\Big[\sum_{t=1}^{T-1} \|-\boldsymbol{\lambda}_{t+1} + \boldsymbol{\lambda}_{t}\|_{1} \|\boldsymbol{J}(\boldsymbol{\theta}_{t+1})\|_{\infty} + |\boldsymbol{\lambda}_{1}|_{1} \|\boldsymbol{J}(\boldsymbol{\theta}_{1})\|_{\infty}$$

$$+ |\boldsymbol{\lambda}_{T}|_{\infty} \|\boldsymbol{J}(\boldsymbol{\theta}_{T+1})\|_{\infty}\Big]$$
(68)

$$\leq \sum_{t=1}^{T-1} \left[ \eta_{\boldsymbol{\lambda}, t} \mathbb{E}[\|\boldsymbol{\lambda}_t - \widehat{\boldsymbol{\lambda}}_t\|_1] \frac{1}{1 - \|\boldsymbol{\gamma}\|_{\infty}} \right] + \frac{2}{1 - \|\boldsymbol{\gamma}\|_{\infty}}$$
 (69)

$$\leq \frac{2}{1 - \|\gamma\|_{\infty}} \left( 1 + \sum_{t=1}^{T} \eta_{\lambda, t} \right), \tag{70}$$

where the inequality (68) comes from the fact that  $x^\top y \leq \|x\|_1 \|y\|_\infty, \forall x, y \in \mathbb{R}^M$  and the inequality (69) is obtained from the update of  $\lambda_t$  in (17). Taking (70) into (67), we can prove this theorem.