

Decentralized Policy Gradient Method for Mean-Field Linear Quadratic Regulator with Global Convergence

Abstract

The scalability of multi-agent reinforcement learning methods to a large number of population is drawing more and more attention in both practice and theory. We consider the basic yet important model, *i.e.*, linear quadratic regulator (LQR), in a mean-field approximation scheme against the curse of the action space dimensions and the exponential growth of agent interactions. Several methods proposed in the mean-field setting require a centralized controller, which is unrealistic in practice. In this paper, we present the first decentralized policy gradient method (MF-DPGM) for mean-field multi-agent reinforcement learning, where exchangeable agents of a large team communicate via a connected network. After a linear transformation of states and policies, we update the new local and mean-field policies by a decentralized gradient primal-dual algorithm respectively in a decoupled way, in order to achieve a global policy consensus. We also give a rigorous proof of the global convergence rate of MF-DPGM by studying the geometry of the problem and estimating one-step progress under a decentralized scheme. In addition, extensive experiments are conducted to support our theoretical findings.

1. Introduction

Recent years have witnessed a promising resurgence of multi-agent reinforcement learning (MARL) in data-driven settings and large population applications. Motivating applications span over multi-robotics systems [11], autonomous driving [38, 48], and sensor networks [44, 12]. MARL involves a set of agents learning to make decisions that minimize their accumulative cost by iterative interactions with a shared environment [50, 51, 7]. As a result, a fundamental difficulty in MARL is that changes in the policy of one agent will affect those of the others, and vice versa [41]. Furthermore, large modern multi-agent systems result in an exponential growth of the dimension of the joint action space with the number of agents. Hence the classical MARL methods [5, 35, 32] via either equilibrium-solving or a few controllers stagger in large-scale applications. Additionally, although a central controller receiving costs and determining actions reduces MARL to a classical MDP which can be solved by existing single-agent RL approaches, the central controller is usually costly to install and the communication overhead degrades the scalability and robustness.

Motivation of mean-field settings. In this paper, we consider homogeneous large-scale MARL systems with symmetry, where each agent has the same cost function and state transition rule. To address the complicated correlations in multi-agent systems, [14] and [42] consider accounting for the extra information from conjecturing the policies of other agents, while [30] and [61] study the decentralized actor-critic algorithm. On the other hand, the mean field approximation [28, 22] serves as an effective alternative to modelling strategic interactions for large populations with symmetry. To characterize the mean-field effect in finite-agent systems, [1] shows that any exchangeable system, where exchanging any two agents does not affect the dynamics or costs, is equivalent to one where the dynamics and costs are coupled across agents through the mean field (empirical mean state).

More importantly, compared to other formulations, the mean-field formulation greatly alleviates the curse of action space dimensionality by a symmetric global optimal policy and a cost function for all agents while decoupling complex correlations in large interactive systems.

However, neither decentralized algorithms nor accompanied theoretical guarantees are studied for the well-behaved mean-field setting. To fill this void, we study the decentralized exchangeable multi-agent systems in the collaborative setting where each agent seeks the optimal policy that minimizes the accumulative global cost over all the agents, via neighborhood communications by a connected network. We propose the first decentralized policy learning scheme with smaller exploration space, less communication, and more robustness. Moreover, we study the nonconvex landscape of the cost functions, which are combined with one-step iterate progress to establish a sublinear global convergence rate for LQR. Our contributions are concluded as follows: (1) First, we formulate the policy gradients for MARL under the mean-field setting. (2) Then we proposed the first decentralized algorithm (MF-DPGM) to effectively learn the optimal policy for mean-field MARL. (3) We present a novel global convergence guarantee for MF-DPGM under mild assumptions and simulation results to justify the practical performance of our algorithm.

Related Work. There has been a line of work on solving normal MARL problems. Based on the seminal work for the framework of Markov games [36], follow-up works, such as [29, 37, 20], studied both collaborative and competitive relationships among agents. Recently, MARL with a large population [47, 38] becomes increasingly popular, such as urban transportation [48, 38], social dilemmas [31, 23], multi-robotics systems [11], and power grids [8], wherein the curse of dimensionality for learning and control [41] arises. Our work is in the line of collaborative settings, where a central controller can help solve MARL by existing single-agent algorithms [6, 55, 40]. Nonetheless, due to high cost to set central controllers in large-scale applications [26], a series of decentralized methods [58, 61, 30] are developed following [62] to learn optimal policies with local rewards and actions.

Mean-field approximation of the system stems from [52], which is generalized to multi-agent scenarios by [21, 27, 28, 22]. Our setting is also closely related to mean-field control for exchangeable agents [39, 1] but from a model-free reinforcement learning perspective. Although another independent simultaneous work [9] investigate MARL for a mean-field case, only the centralized algorithm with infinite players is considered and directly reduced to single-agent LQR when doing variable transformation.

Moreover, Our method is also related to the line of work addressing the challenges of non-convexity in distributed / decentralized optimization; see recent developments in [3, 63, 19, 57, 16]. [18] develops a non-convex ADMM based methods with distributed consensus, which enjoys a similar first global sublinear convergence rate as ours. Nevertheless, the network considered therein is a star network with a central controller. A primal-dual method for unconstrained problem over a connected network and a global convergence rate is derived in [19]. References [16, 25, 34] developed methods for distributed stochastic zeroth and first-order non-convex optimization.

Notations. We introduce some frequently used notations here, while others will be defined later. We denote by $[n]$ the set of integers $\{1, 2, \dots, n\}$, by $\bar{\mathcal{N}}$ all of the nonnegative integers. With slight abuse of notation, we use $\langle \cdot, \cdot \rangle$ to denote the inner-product for vectors, matrices, tensors, and block-wise cases according to the context. For a matrix $X \in \mathbb{R}^{d_1 \times d_2}$, we denote by $\text{vec}(X) \in \mathbb{R}^{d_1 d_2}$ the vectorization of X . For a vector v , we define $\|v\|_X^2 := v^\top X v$ as the squared norm of v under metric matrix X . When v is a matrix or tensor, $\|v\|$ refers to the norm of its vectorization. The

mode- n matrix product of tensor $\mathbb{X} \in \mathbb{R}^{d_1 \times \dots \times d_p}$ and a matrix $X \in \mathbb{R}^{e \times d_n}$ is a tensor $X\mathbb{X}_{(n)} \in \mathbb{R}^{d_1 \times \dots \times d_{n-1} \times e \times \dots \times d_p}$. Without specifying the mode, $X\mathbb{X}$ denotes mode-1 matrix product. We also denote by $\mathbf{1}$ and $\mathbf{0}$ the all-one and all-zero vector (resp. tensor) respectively by the context, by $\sigma_{\min}(X)$ the second smallest eigenvalue of X .

2. Preliminaries

We study discrete-time linear-quadratic MARL under mean-field settings with exchangeable finite n agents. The states of the system at time-step t is given by $\{x_t^{(1)}, x_t^{(2)}, \dots, x_t^{(n)}\}$, where $x_t^{(i)} \in \mathbb{R}^d$ denotes the state vector of the i^{th} agent ($i \in [n]$). Then the system dynamics is described as follows,

$$x_{t+1}^{(i)} = Ax_t^{(i)} + Bu_t^{(i)} + \bar{A}\bar{x}_t + w_t^{(i)}, \quad \forall i \in [n], \quad (2.1)$$

where $u_t^{(i)} \in \mathbb{R}^m$ denotes the control (the action), a Gaussian noise $w_t^{(i)} \sim N(0, I_d)$ independent of each agent and time-step is added to the succeeding state, $\bar{x}_t = 1/n \sum_{i=1}^n x_t^{(i)}$ denotes the mean-field state of the system at time-step t . For such dynamics, we have a collective cost function of the distributed system at time-step t ,

$$c_t = \sum_{i=1}^n x_t^{(i)\top} Q x_t^{(i)} + u_t^{(i)\top} R u_t^{(i)} + \bar{x}_t^\top \bar{Q} \bar{x}_t. \quad (2.2)$$

We also define the cost for agent i as $c_t^{(i)} = x_t^{(i)\top} Q x_t^{(i)} + u_t^{(i)\top} R u_t^{(i)} + \bar{x}_t^\top \bar{Q} \bar{x}_t$ at time-step t , where Q, R, \bar{Q} are positive semidefinite matrices. Furthermore, it is shown that optimal control for the i^{th} agent can be written as a linear combination of $x_t^{(i)}$ and \bar{x}_t , $u_t^{(i)} = Kx_t^{(i)} + L\bar{x}_t$, where the same matrices K, L apply to all agents by the symmetry results in optimal control [1]. By plugging $u_t^{(i)}$ into (2.1), we can rewrite $x_{t+1}^{(i)}$ as $x_{t+1}^{(i)} = (A + BK)x_t^{(i)} + (\bar{A} + BL)\bar{x}_t + w_t^{(i)}$. Let $\Theta = (K; L) \in \mathbb{R}^{2 \times m \times d}$, beginning with initial states $\{x_0^{(i)}\}_{i=1}^n$, our goal is to find the controls $\{u_t^{(i)}\}_{i=1}^n$ ($t \geq 0$) minimizing the long-term collective cost,

$$C(\Theta) = \mathbb{E}_{\mathbf{x}_0, \mathbf{w}} \left[\sum_{t=0}^{\infty} \gamma^t c_t \right] = \mathbb{E}_{\mathbf{x}_0, \mathbf{w}} \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^n (x_t^{(i)\top} Q x_t^{(i)} + u_t^{(i)\top} R u_t^{(i)} + \bar{x}_t^\top \bar{Q} \bar{x}_t). \quad (2.3)$$

We denote by $\gamma \in (0, 1)$ the discounted factor in the infinite horizon. The expectation is taken *w.r.t.* all the initial states $\mathbf{x}_0 = (x_0^{(1)\top}, \dots, x_0^{(n)\top})^\top$ and all independent noise terms $\mathbf{w} = \{w_t^{(i)}\}_{i \in [n], t \in \mathbb{N}}$. Therefore, our goal is to solve the following infinite horizon mean-field LQR problem,

$$\text{minimize } C(\Theta) \text{ s.t. } x_{t+1}^{(i)} = Ax_t^{(i)} + Bu_t^{(i)} + \bar{A}\bar{x}_t + w_t^{(i)}, \quad x_0^{(i)} \sim \mathcal{D} \text{ for each } i \in [n]. \quad (2.4)$$

3. Algorithm

Since direct policy gradients over the matrix parameters K (resp. L) lead to gradients of K (resp. L) containing the other parameter L (resp. K) due to the correlation mean-field term \bar{x}_t across different agents, it is hard to apply analysis of standard policy gradients. To deal with the correlated states in the mean-field setting, we adopt a reparametrization trick to obtain the dynamics below,

$$u_t^{(i)} = Kx_t^{(i)} + L\bar{x}_t = M(x_t^{(i)} - \bar{x}_t) + N\bar{x}_t \triangleq My_t^{(i)} + N\bar{y}_t. \quad (3.1)$$

Consequently, we have the total cost in new states as below.

$$C(\Theta) = \sum_{i=1}^n \text{Tr}[(Q + M^\top RM)\Sigma_M^{(i)}] + n \text{Tr}[(\bar{Q} + N^\top RN)\Sigma_N],$$

where $\Sigma_M^{(i)} = \mathbb{E}_{y_0, w} [\sum_{t=0}^{\infty} \gamma^t y_t^{(i)} y_t^{(i)\top}]$, $\Sigma_N = \mathbb{E}_{\bar{y}_0, w} [\sum_{t=0}^{\infty} \gamma^t \bar{y}_t \bar{y}_t^\top]$, and $\tilde{Q} = Q + \bar{Q}$. See Section A for details. However, without a central controller, agent i possesses $M^{(i)}$ and $N^{(i)}$, defined as the policy iterates of agent i in decentralized scenarios, as its local policy before convergence and is restricted to communicating policies with neighbor agents over a network, although each agent has access to the global mean state. Similar problems emerge when updating $N^{(i)}$'s, which is the mean-field policy of each agent before finding the optimum. Hence, we resort to performing a decentralized optimization scheme with global consensus. As M and N are decoupled into two similar update processes, below we denote by $\tilde{\Theta}^{(i)}$ either $M^{(i)}$ or $N^{(i)}$. We concatenate policy parameters together as a higher dimensional tensor, *i.e.*, $\tilde{\Theta} = [\tilde{\Theta}^{(1)}; \tilde{\Theta}^{(2)}; \dots; \tilde{\Theta}^{(n)}] \in \mathbb{R}^{n \times m \times d}$ to provide a compact formulation of the decentralized optimization problem,

$$\min_{\tilde{\Theta}} \tilde{C}(\tilde{\Theta}) = \frac{1}{n} \sum_{i=1}^n C(\tilde{\Theta}^{(i)}), \quad \text{s.t.} \quad \tilde{\Theta}^{(i)} = \tilde{\Theta}^{(j)} \quad \text{for } (i, j) \in \mathcal{E}, \quad (3.2)$$

where the communication network is considered as an *unweighted* and *undirected* connected graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ with vertex set \mathcal{V} and edge set \mathcal{E} . Such a description gives a separable global objective function $\tilde{C}(\tilde{\Theta})$ and linear constraints indicating the connectivity property of the communication network. In next section, we will introduce a more tractable formulation for algorithms by graph theory. Note that the summands in (3.2) are disjoint components of $\tilde{\Theta}$. Therefore, the gradient of the global objective should be expanded as: $\nabla \tilde{C}(\tilde{\Theta}) = 1/n \left(\nabla C(\tilde{\Theta}^{(1)}); \dots; \nabla C(\tilde{\Theta}^{(n)}) \right)$. For simplicity in some context, we also define the vectorization of $\tilde{\Theta}$ as $\Theta = [\text{vec}(\tilde{\Theta}^{(1)})^\top, \dots, \text{vec}(\tilde{\Theta}^{(n)})^\top]^\top$.

3.1 Communication Structure and Algebraic Representations

In this section, we capture the structure of \mathcal{G} with $|\mathcal{V}| = n$ and $|\mathcal{E}| = e$ by some tools from spectral graph theory. Define the edge-associated (Ω_θ) and agent-associated (Γ_θ) parameters:

$$\Omega_\theta := \text{diag}(\sigma_1^\theta, \dots, \sigma_e^\theta) \succ 0, \quad \Gamma_\theta := \text{diag}(\gamma_1^\theta, \dots, \gamma_n^\theta) \succ 0, \quad (3.3)$$

where θ can be replaced by M or N to represent two sets of configurations since the policy parameters are decoupled into two optimization processes. Hereafter if we use symbols without specifying M or N , the statement will hold for both M and N , respectively. Further, when $\tilde{\Theta}^{(i)} = M^{(i)}$ (resp. $N^{(i)}$), we set $\theta = M$ (resp. N) in the rest of the paper. We will see how Ω_θ and Γ_θ serve as step-size for MF-DPGM in Section 3.2. We also use σ_{ij}^θ , where $(i, j) \in \mathcal{E}$ as an alternative to σ_k^θ ($k \in [e]$) to emphasize σ_k^θ is assigned to the edge $(i, j) \in \mathcal{E}$. The quantity assigned to agent i is used for $M^{(i)}$ or $N^{(i)}$, or both by the context. We define $D = \text{diag}(d_1, \dots, d_n)$ with d_i as the degree of vertex (agent) i . Denote the graph Laplacian matrix and its scaled version [10] as \mathcal{L} and $\tilde{\mathcal{L}}$ respectively, where $\mathcal{L} = I_A D^{1/2}$ with I_A being the incidence matrix [10]. We further define $\mathcal{L}_a := |\tilde{\mathcal{L}}| \in \mathbb{R}^{e \times n}$ by taking element-wise absolute value on $\tilde{\mathcal{L}}$. By definitions we can find the following relations: $D = 1/2(\tilde{\mathcal{L}}^\top \tilde{\mathcal{L}} + \mathcal{L}_a^\top \mathcal{L}_a)$, $\tilde{\mathcal{L}} \mathbf{1} = 0$. For notational simplicity, we define

$$\tilde{\Omega}_\theta = \text{diag} \left(\left\{ \sum_{j:j \sim i} (\sigma_{ij}^\theta)^2 \right\}_{j \in \mathcal{G}} \right) = 1/2(\tilde{\mathcal{L}}^\top \Omega_\theta^2 \tilde{\mathcal{L}} + \mathcal{L}_a \Omega_\theta^2 \mathcal{L}_a), \quad (3.4)$$

$$H_\theta := \mathcal{L}_a^\top \Omega_\theta^2 \mathcal{L}_a + \Gamma_\theta^2, \quad (\eta_i^\theta)^2 = 2 \sum_{j:j \sim i} (\sigma_{ij}^\theta)^2 + (\gamma_i^\theta)^2. \quad (3.5)$$

With these parameters, we can rewrite (3.2) in a more tractable form for numerical algorithms:

$$\min_{\tilde{\Theta}} \tilde{C}(\tilde{\Theta}) = 1/n \sum_{i=1}^n C(\tilde{\Theta}^{(i)}), \quad \text{s.t.} \quad \tilde{\mathcal{L}} \tilde{\Theta} = \mathbf{0}.$$

The idea of MF-DPGM relies on this formulation from a primal-dual view.

3.2 The MF-DPGM Algorithm

In this section, we introduce the first decentralized algorithm to learn the optimal policy in mean-field LQR setting (2.4). We relax the linear constraints in (3.2), and gradually enforces them with the development of the algorithm. Our method mainly involves a policy evaluation step and a communication and update step.

Policy evaluation: At the beginning of the k -th iteration, we sample n_p trajectories of each agent i according to its own current policy $\tilde{\Theta}_k^{(i)}$, to estimate the action-value function by $1/n_p \sum_{p=1}^{n_p} \hat{Q}_{i,p}^\pi(x_t^{(i)}, u_t^{(i)})$ for policy evaluation in model-free settings¹. Furthermore, to retain accuracy for large populations, we adopt gradient estimator in REINFORCE [54] for the cost function instead of the biased zeroth order method [13]. In practice, the deterministic policy $u_t^{(i)}$ can be realized by a Gaussian policy $\pi_{\Theta^{(i)}}$ with a nearly zero variance.

Communication and update: After running policies at the k -th iteration, each agent i collects policies from neighbors and its own gradient estimators of the cost function at the k -th and $(k-1)$ -th iterations, and combines them linearly by Ω_θ and Γ_θ as a decentralized version of policy gradient method [2] to update the local policy $\tilde{\Theta}_k^{(i)}$. As a result, we can view tunable parameters Ω_θ and Γ_θ as step-sizes in MF-DPGM. See (3.11) for a detailed update. A practical choice [15, 49] is $\Omega_\theta^2 := \alpha^2 I$ as a multiple of identity, which involves the sum of neighbor iterates.

Considering the decentralized nature of the problem, our convergence measure ε not only takes cost error into account, but also consider consensus error to guarantee that the estimated policies for agents converge to the identical optimal policy in Section 2. To see this, we demonstrate the connection between our algorithm and the primal-dual paradigm commonly used for solving constrained optimization. Write $\tilde{\Theta}_k = [\tilde{\Theta}_k^{(1)}; \tilde{\Theta}_k^{(2)}; \dots; \tilde{\Theta}_k^{(n)}] \in \mathbb{R}^{n \times m \times d}$ for $k \geq -1$. We introduce the augmented Lagrangian function (\mathcal{A}) in tensor variables:

$$\mathcal{A}(\tilde{\Theta}_k, \Lambda_k) = \tilde{C}(\tilde{\Theta}_k) + \langle \Lambda_k, \tilde{\mathcal{L}}\tilde{\Theta}_k \rangle + \frac{1}{2} \|\Omega_\theta \tilde{\mathcal{L}}\tilde{\Theta}_k\|^2, \quad (3.6)$$

where $\Lambda_k \in \mathbb{R}^{e \times m \times d}$ is the dual variable at iteration k , which is updated as

$$\Lambda_{k+1} = \Lambda_k + \Omega_\theta \tilde{\mathcal{L}}\tilde{\Theta}_{k+1}. \quad (3.7)$$

We define $\mathcal{A}_k := \mathcal{A}(\tilde{\Theta}_k, \Lambda_k)$. By plugging (3.7) into (3.11) and using the definition of H_θ in (3.5) we have

$$\nabla \tilde{C}(\tilde{\Theta}_k) + H_\theta(\tilde{\Theta}_{k+1} - \tilde{\Theta}_k) + \tilde{\mathcal{L}}\tilde{\Theta}_k + \tilde{\mathcal{L}}^\top \Omega^2 \tilde{\mathcal{L}}\tilde{\Theta}_{k+1} = \mathbf{0}. \quad (3.8)$$

Note that $\mathbf{0} \in \mathbb{R}^{n \times m \times d}$ here. Viewing (3.8) as the optimal (first-order) condition of some function, we observe that update (3.11) is an implementation of

$$\begin{aligned} \tilde{\Theta}_{k+1} = \arg \min_{\tilde{\Theta}} & \left\langle \nabla \tilde{C}(\tilde{\Theta}_k) + \tilde{\mathcal{L}}^\top \Lambda_k, \tilde{\Theta} - \tilde{\Theta}_k \right\rangle \\ & + \frac{1}{2} \|\Omega_\theta \tilde{\mathcal{L}}\tilde{\Theta}\|^2 + \frac{1}{2} \|\Omega_\theta \mathcal{L}_a(\tilde{\Theta} - \tilde{\Theta}_k)\|^2 + \frac{1}{2} \|\Gamma_\theta(\tilde{\Theta} - \tilde{\Theta}_k)\|^2 \end{aligned} \quad (3.9)$$

together with (3.7), where the third term of (3.9) encodes the network structure that we utilize for neighborhood averaging. By definitions in § 3.1, the terms $\|\Omega_\theta \tilde{\mathcal{L}}\tilde{\Theta}\|^2$ shows how close policies of agents are to each other with policy $\tilde{\Theta}$. Such a primal-dual interpretation is closely related to some classical constrained optimization methods, such as the Uzawa method [56] and the prox-

1. Also, we can share part of noise with term $u_t^{(i)}$ in (2.1) to make stochastic policies. To highlight our main idea, we focus on theoretical results with exact gradients.

Algorithm 1: Mean-Field Decentralized Policy Gradient Method (MF-DPGM)

Data: Agent dynamics n, A, B, \bar{A} ; Cost parameters Q, R, \bar{Q} ; Initial state distribution \mathcal{D} .

Input: Network \mathcal{G} ; Ω, Γ ; Length of horizon T ; Number of sample paths n_p ; Estimation error ϵ .

Output: Estimation of optimal control policies $\hat{\Theta}$.

Initialization: $\tilde{\Theta}_{-1} = \mathbf{0}$; $\tilde{\Theta}_0^{(i)} = \nabla C^{(i)}(\mathbf{0})\eta_i^{-2}/n, \forall i \in [n]$; Iteration $k \leftarrow 1$.

while $\varepsilon(k) > \epsilon$ **do**

for path $p = 1$ to n_p **do**

for $t = 1$ to T **do**

$$u_t^{(i)} = M_k^{(i)} y_t^{(i)} + N_k^{(i)} \bar{y}_t;$$

$$x_{t+1}^{(i)} \leftarrow Ax_t^{(i)} + Bu_t^{(i)} + \bar{A}\bar{x}_t + w_t^{(i)}, \text{ for all } i \in [n] \text{ in parallel};$$

$$\hat{Q}_{i,p}^\pi(x_t^{(i)}, u_t^{(i)}) \leftarrow \sum_{s=t}^T \gamma^{s-t} c_s^{(i)};$$

end

end

 Compute $\hat{\nabla} C(\tilde{\Theta}_k^{(i)}) \leftarrow 1/n_p \sum_{p=1}^{n_p} \sum_{t=1}^T \hat{Q}_{i,p}^\pi(x_t^{(i)}, u_t^{(i)}) \cdot \nabla_{\tilde{\Theta}} \log \pi_{\tilde{\Theta}_k^{(i)}}(u_t^{(i)} | x_t^{(i)})$, for all $i \in [n]$

Communication and update: For all $i \in [n]$,

$$\begin{aligned} \tilde{\Theta}_{k+1}^{(i)} &\leftarrow \tilde{\Theta}_k^{(i)} - \frac{1}{(\eta_i^\theta)^2} \left(\frac{1}{n} \left(\nabla \hat{C}(\tilde{\Theta}_k^{(i)}) - \nabla \hat{C}(\tilde{\Theta}_{k-1}^{(i)}) \right) \right. \\ &\quad \left. - 2 \sum_{j:j \sim i} (\sigma_{ij}^\theta)^2 \tilde{\Theta}_k^{(j)} + (\gamma_i^\theta)^2 \left(\tilde{\Theta}_{k-1}^{(i)} - \tilde{\Theta}_k^{(i)} \right) + \sum_{j:j \sim i} (\sigma_{ij}^\theta)^2 \left(\tilde{\Theta}_{k-1}^{(j)} + \tilde{\Theta}_{k-1}^{(i)} \right) \right) \end{aligned} \quad (3.11)$$

$$\hat{\Theta} \leftarrow \tilde{\Theta}_k, k \leftarrow k + 1$$

end

MM [46, 59]. Following distributed constrained optimization approaches [24, 17], we set the error

$$\varepsilon(t) = \min_{s \in [t]} \left| 1/n \sum_{i=1}^n C(\tilde{\Theta}_s^{(i)}) - \tilde{C}(\tilde{\Theta}^*) \right| + \|\Omega_\theta \tilde{\mathcal{L}} \tilde{\Theta}_s\|^2 \quad (3.10)$$

for $t \in \bar{\mathcal{N}}$ to monitor consensus.

4. Theoretical Results and Analysis

In this section, we provide a sublinear global convergence rate achieved by MF-DPGM. To simplify further notations, we define $\Sigma_0^{(i)} := \mathbb{E} y_0^{(i)} y_0^{(i)\top}$, $\Sigma_0 := \mathbb{E} \bar{y}_0 \bar{y}_0^\top$, $\xi_i := \sigma_{\min}(\Sigma_0^{(i)})$, $\bar{\xi} := \sigma_{\min}(\Sigma_0)$, where $\sigma_{\min}(X)$ denotes the smallest singular value of X . Additionally, we have the following condition on edge and agent associated parameters, *i.e.* Ω_θ and Γ_θ defined in (3.3), to guarantee the development of our decentralized update.

Condition 4.1 . The parameters of MF-DPGM are chosen to satisfy for any $k \geq 1$ and any $i \in [n]$,

$$\frac{1}{2} (\Omega_\theta + \Gamma_\theta^2) \succeq \frac{(2c+1)\Phi_\theta}{n} + \frac{4\kappa}{n^2} \Phi_\theta \Gamma_\theta^{-2} \Phi_\theta, \quad c = \max\{1, 6\kappa\}, \quad \Gamma_\theta^2 \succeq \frac{\Phi_\theta \Gamma_\theta^{-2} \Phi_\theta}{n^2}, \quad (4.1)$$

where $\kappa = 1/\underline{\sigma}_{\min}(\Omega_\theta \tilde{\mathcal{L}} H_\theta^{-1} \tilde{\mathcal{L}}^\top \Omega_\theta)$, $\theta \in \{M, N\}$. Also, $\Phi_M = \text{diag}(\beta_1^M, \dots, \beta_n^M) \otimes I_{md} \in \mathbb{R}^{nmd \times nmd}$ and $\Phi_N = \text{diag}(\beta_1^N, \dots, \beta_n^N) \otimes I_{md}$ are almost-smoothness constants specified in Lemma D.2.

In particular, here κ encodes the network structure and c helps construct the potential function indicating the one-step progress of MF-DPGM. By properly choosing Ω_θ and Γ_θ , Condition 4.1 is met. Hence, the potential function (C.1) decreases and boundary condition (D.9) of two adjacent

iterates is satisfied by dynamically adjusting the adaptive area in (D.9), which is illustrated in D.3. Below we present the global convergence result of MF-DPGM.

Theorem 4.2. For $i, j \in [n]$, we choose the σ_{ij}^θ and γ_i^θ at timestep t to be

$$\sigma_{ij}^\theta, \gamma_i^\theta = \text{poly}\left(\frac{\xi\sigma_{\min}(Q_\theta)}{C(\tilde{\Theta}_0^{(i)})}, \frac{1}{\|A_\theta\|}, \frac{1}{\|A\|}, \frac{1}{\|B\|}, \frac{1}{\|R\|}, \sigma_{\min}(R), \|\tilde{\Theta}_t^{(j)}\|, \|\tilde{\Theta}_{t-1}^{(j)}\|, \|\tilde{\Theta}_{t-1}^{(i)}\|, \|\tilde{\Theta}_{t-1}^{(i)} - \tilde{\Theta}_t^{(i)}\|\right), \quad (4.2)$$

where $Q_M = Q$, $Q_N = Q + \bar{Q}$, $A_M = A + BM$, and $A_N = A + \bar{A} + BN$. Then under Condition 4.1, within the timestep t of Algorithm 1, the iterates of MF-DPGM satisfies

$$\min_{s \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(\tilde{\Theta}_s^{(i)}) - \tilde{C}(\tilde{\Theta}^*) \right| + \|\Omega_\theta \tilde{\mathcal{L}} \tilde{\Theta}_s\|^2 \leq \underbrace{\frac{8\alpha_g C C'}{t}}_{\text{cost error bound}} + \underbrace{\frac{20C'}{t}}_{\text{consensus error bound}}, \quad (4.3)$$

where $\alpha_g = \max\{\alpha_g^M, \alpha_g^N\}$, with $\alpha_g^M = \|\Sigma_{M^*}\|/[\sigma_{\min}(\Sigma_0^{(i)})^2 \sigma_{\min}(R)]$ and $\alpha_g^N = \|\Sigma_{N^*}\|/[\sigma_{\min}(\Sigma_0)^2 \sigma_{\min}(R)]$, is a problem related constant for gradient domination, which also appears in Lemma D.3, $\tilde{\Theta}^*$ indicates the optimal policy parameters, and

$$C' = \tilde{C}(\tilde{\Theta}_0) - \tilde{C}(\tilde{\Theta}^*) + \frac{2\nabla C(\mathbf{0})^\top \Phi_\theta^{-1} \nabla C(\mathbf{0})}{n}, \quad C = 4 \sum_{(i,j):i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \quad (4.4)$$

are absolute constants. Here $\nabla C(\mathbf{0}) \in \mathbb{R}^{nmd}$.

Proof. See Appendix D.5 for a detailed proof. \square

As many other decentralized algorithms, our convergence results also rely on the smoothness property (Lipschitzness of the gradient) of the objective. However, we notice that the optimization landscape for each agent is not strictly smooth due to unstability of $A + BM$ (resp. $A + \bar{A} + BN$). At the boundary between stable and unstable policies, the cost function rapidly becomes infinity, which violates the traditional smoothness conditions. To address this issue, we regulate the parameters Ω_θ and Γ_θ for a small stepsize to set the next iterate $\tilde{\Theta}_{k+1}^{(i)}$ sufficiently close to the current one, so that $\Sigma_{\tilde{\Theta}_{k+1}^{(i)}} \approx \Sigma_{\tilde{\Theta}_k^{(i)}} + \mathcal{O}\left(\|\tilde{\Theta}_{k+1}^{(i)} - \tilde{\Theta}_k^{(i)}\|\right)$ for $k \geq 0$. Specifically, we adopt large enough σ_{ij} 's and γ_i 's by observing that feasible Ω_θ and Γ_θ multiplied by identical large scaling factors in (4.1) still meet the condition. Using tools from [13], we show in Lemma D.2 and C.1 the continuity property of the cost functions, state trajectories, and gradients of cost corresponding to $\tilde{\Theta}_k$ within an adaptive area of policies for each agent.

The main upshot of Theorem 4.2 is to characterize in detail the first global sublinear convergence rate for the overall error including the cost and consensus components. More precisely, MF-DPGM not only drives the average cost of agents to the optimal cost over time, but also aim at the same optimal policy for each agent. An alternative way to display the rates is distributedly showing for each agent with slightly different constants where C' is decomposed accordingly. Moreover, the optimization of $\tilde{\Theta}$ can be split into the development for local and mean-field policy respectively, which is frequently observed in further theoretical illustration. The difference of updates of these two policies sits between the only one mean-field state associated with $N^{(i)}$'s and corresponding multiple states for $M^{(i)}$'s. The constant C also reveals the dependency on the graph of error bounds, leading to guidance for parameter choice ($\Gamma_\theta, \Omega_\theta$) to optimize the rates. For example, setting

$\sigma_{ij}^2 = \frac{\tau^2 \sqrt{\beta_i \beta_j}}{\sqrt{d_i d_j}}$, $\Gamma^2 = \tau^2 \Phi$, where $\tau^2 = \frac{80 \max\{\sigma_{\max}(Z), 1\}}{n \min\{\sigma_{\min}(\mathcal{L}_G), 1\}}$ with \mathcal{L}_G being the generalized Laplacian, $Z_{ii} = \sum_{k:k \sim i} \frac{\sqrt{\beta_i \beta_k}}{\sqrt{d_i d_k}}$, results in a tight rate for complete graphs. See [53] for more details.

5. Proof Sketch

We sketch the proof of results in Section 4. Define $\Sigma_0^{(i)} := \mathbb{E} y_0^{(i)} y_0^{(i)\top}$, $\Sigma_0 := \mathbb{E} \bar{y}_0 \bar{y}_0^\top$, $\xi_i := \sigma_{\min}(\Sigma_0^{(i)})$, $\bar{\xi} = \sigma_{\min}(\mathbb{E}_{y_0^{(i)} \sim \mathcal{D}} \bar{y}_0 \bar{y}_0^\top)$, where $\sigma_{\min}(X)$ denotes the smallest singular value of X .

Geometry of cost functions. As mentioned in Section 4, Theorem 4.2 requires moderate smoothness of the landscape of cost functions. Based on the almost Lipschitzness of positive definite matrix P_θ parameterizing the optimal cost from a state going forward, and almost Lipschitzness of the Σ_θ which plays a key role in cost function. Lemma D.2 for each cost function with an almost Lipschitz gradient and Lemma D.3 for each cost function almost dominated by the gradient are crucial in bounding the one-step progress in Lemma D.6 and D.7.

Progress of decentralized iterations. To estimate the one-step progress of Algorithm 1, we construct the auxiliary potential function U in both primal ($\tilde{\Theta}_k$) and dual (Λ_k) variables as follows,

$$U_{k+1} = U(\tilde{\Theta}_{k+1}, \tilde{\Theta}_k, \Lambda_{k+1}) = \mathcal{A}_{k+1} + \frac{2\kappa}{n^2} \|\Gamma_\theta^{-1} \Phi_\theta(\tilde{\Theta}_{k+1} - \tilde{\Theta}_k)\|^2 \quad (5.1)$$

$$+ \frac{c}{2} \left(\|\Omega_\theta \tilde{\mathcal{L}} \tilde{\Theta}_{k+1}\|^2 + \|\tilde{\Theta}_{k+1} - \tilde{\Theta}_k\|_{H_\theta + \Phi_\theta/n}^2 \right),$$

where c is a constant chosen according to Condition 4.1 and \mathcal{A}_{k+1} is the augmented Lagrangian defined in (3.6). The decrease of the potential function U at each step has a nonnegative lower bound, as illustrated in the following lemma.

Lemma 5.1. Suppose that parameters of MF-DPGM are chosen according to (4.1). Then it holds that

$$U_k - U_{k+1} \geq \frac{1}{4} \|\tilde{\Theta}_{k+1} - \tilde{\Theta}_k\|_{\Omega_\theta + \Gamma_\theta^2}^2 + \kappa \|\mathbb{V}_{k+1}\|_{H_\theta}^2 \quad (5.2)$$

for any $k \geq 0$, where $\mathbb{V}_{k+1} := (\tilde{\Theta}_{k+1} - \tilde{\Theta}_k) - (\tilde{\Theta}_k - \tilde{\Theta}_{k-1})$. Moreover, we have the following bounds,

$$U_k \leq U_0 \leq \tilde{C}(\tilde{\Theta}_0) + \frac{2\nabla \tilde{C}(\mathbf{0})^\top \Phi^{-1} \nabla \tilde{C}(\mathbf{0})}{n}, \quad U_{k+1} \geq \tilde{C}(\tilde{\Theta}^*) > -\infty \quad (5.3)$$

for any $k \geq 1$, where $\mathbf{0} \in \mathbb{R}^{nmd}$ is the all-zero vector, and $\nabla \tilde{C}(\mathbf{0}) = 1/n \cdot (\nabla C^{(1)}(\mathbf{0}), \dots, \nabla C^{(n)}(\mathbf{0}))$.

Proof. See the appendix of [53] for a detailed proof. \square

Note that decreasing potential function in Lemma C.1 also tracks stability of decentralized LQR in the optimization process. We start by optimality condition (3.8) and derive upper bounds for objective gradient norms by the distances of primal variables. Then Lemma C.1 is applied reducing the bounds to differences of adjacent potential functions. Similarly, the consensus error is controlled using Lemma D.6 and D.7, where almost smoothness of costs (D.10) is involved, and processed by Lemma C.1 to keep the same difference terms as those of gradient norms. Finally, combining two bounds of similar structure with Lemma D.3 we establish Theorem 4.2. See D.5 for a detailed proof.

6. Experiments

We introduce numerical experiments using Algorithm 1 on both synthetic data and real data.

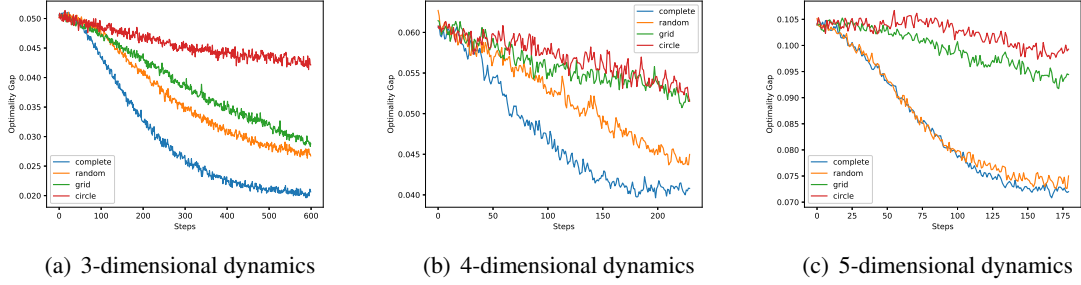


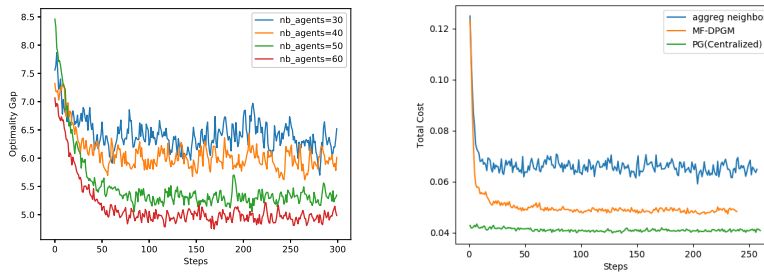
Figure 1: Simulation results (convergence curves) on complete (blue), random (orange), grid (green) and circle (red) networks, as well as different dynamics $d \in \{3, 4, 5\}$.

Experiments on synthetic data. We consider three multi-agent systems with different dimensions under the setting of LQR. Specifically speaking, we consider the settings where the number of agents $n = 25$, the dimension of the dynamics parameters d varies in $\{3, 4, 5\}$, and the dimension of action space $m = d$. For each multi-agent system, we consider four different structures of the communication network \mathcal{G} : complete graph, grid graph, circle graph, and random graph. See Appendix B for details on the experiment setup. We present our numerical results using Algorithm 1 in Figure 1, where the x-axis is the number of timesteps, while the y-axis measures the optimality gap defined in (3.10).

From Figure 1, we see that using our decentralized Algorithm 1, as the number of steps increases, the optimality gap decreases. This not only implies that the local policy of each agent converges to an optimal policy, but also shows that the policies of all agents will be eventually close enough. Also, for a fixed system, we see that different structures of communication networks lead to different convergence speed of Algorithm 1. Specifically speaking,

For fixed dynamics, we see that different topologies of communication networks lead to different convergence speed of Algorithm 1. Generally speaking, a complete graph (blue curve) as the communication network gives the fastest convergence, while a circle graph (red curve) leads to the slowest convergence. Note that the complete graph has the most links among all four graphs, while the circle graph has the least links among all graphs. This implies that our decentralized algorithm performs better as the number of links in the communication network increases. Also, we see from Figure 1 that the dimension of dynamics also have an impact on the number of iterations to achieve convergence. Specifically speaking, the system with lower dimension enjoys a faster convergence, while the system with higher dimension converges more slowly.

In the meanwhile, we plot the training curves in Figure 2(a) for different numbers of agents to show the effectiveness of our method at different scales, where we also justify that the mean-field approximation works better as the size of the population grows by the increasing performance with the larger population, as the effect of mean-field phenomenon is more significant in larger systems. In addition, although it is unfair to compare against the centralized setting, where the controller has access to the costs of all agents and updates their policies simultaneously and identically, we plot the comparison in Figure 2(b) for identification with $n = 25$. We reproduce a baseline of decentralized policy gradient with gossip matrices [45], which aggregates the updated policies from the neighborhood (aggreg. neighbor). Note that our algorithm compares favorably against this baseline and is competitive with the centralized one, which gives a better performance achievable.



(a) Convergence curves on different populations n over a circle graph. (b) Comparison with baselines on a circle graph.

Figure 2: Simulation results with different size of populations and comparison with baselines.

Experiments on real data. We test our algorithm on building energy regulation of Heating Ventilation and Air Conditioning (HVAC) systems for a multi-zone building as in [33], where the task involves controlling the temperatures of 4 rooms (and thus 4 agents) with varying outdoor temperatures. The control for each room is represented by a scalar that indicates the air flow rate of the cooling system in that room. The underlying dynamics of the room temperatures are linear and stochastic, which also indicates a symmetric system. In particular, the system is permutation invariant; the matrices modeling the dynamics are symmetric. The cost function with respect to Room i takes the form of $c_i(t) = \frac{1}{2}(x_i(t) - x_{target})^2 + \alpha u_i(t)$ where $x_i(t)$ and $u_i(t)$ are the temperature and the control of Room i at time t and α is a scalar. For simplicity, we set all environment hyperparameters to be the same as in [33]. The target temperature x_{target} is set to be 22 for both training and testing.

For more technical details and experimental results, see the supplementary material at http://lewis-algo.com/files/mf_lqr_v1.pdf.

7. Conclusion

In this paper, focusing on a simple yet fundamental setting LQR, we first formulate the decoupled policy gradient under mean-field effect by reparameterization. Then, we propose a decentralized policy gradient method where each agent updates the local policy by combining policies from neighborhood with its own policy. It is the first decentralized policy learning algorithm to improve searching efficiency and alleviate communication overhead, potentially applicable to complex mean-field models. In addition, we quantify the non-convex problem geometry by several almost continuity results, which is combined with one-step progresses for our algorithm to establish a sublinear global convergence rate for LQR. Additional experiments justify our theoretical results and show a promising performance.

References

- [1] Jalal Arabneydi and Aditya Mahajan. Linear quadratic mean field teams: Optimal and approximately optimal decentralized solutions. *arXiv preprint arXiv:1609.00056*, 2016.
- [2] Jonathan Baxter and Peter L Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001.
- [3] Pascal Bianchi and Jérémie Jakubowicz. Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization. *IEEE transactions on automatic control*, 58(2):391–405, 2012.

- [4] Sergio Bittanti, Alan J Laub, and Jan C Willems. *The Riccati Equation*. Springer Science & Business Media, 2012.
- [5] Michael Bowling and Manuela Veloso. Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136(2):215–250, 2002.
- [6] Steven J Bradtke, B Erik Ydstie, and Andrew G Barto. Adaptive linear quadratic control using policy iteration. In *Proceedings of 1994 American Control Conference-ACC'94*, volume 3, pages 3475–3479. IEEE, 1994.
- [7] Lucian Bu, Robert Babu, Bart De Schutter, et al. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 38(2):156–172, 2008.
- [8] Duncan S Callaway and Ian A Hiskens. Achieving controllability of electric loads. *Proceedings of the IEEE*, 99(1):184–199, 2010.
- [9] René Carmona, Mathieu Laurière, and Zongjun Tan. Linear-quadratic mean-field reinforcement learning: convergence of policy gradient methods. *arXiv preprint arXiv:1910.04295*, 2019.
- [10] Fan RK Chung and Fan Chung Graham. *Spectral graph theory*. Number 92. American Mathematical Soc., 1997.
- [11] Peter Corke, Ron Peterson, and Daniela Rus. Networked robots: Flying robot navigation using a sensor net. In *Robotics research. The eleventh international symposium*, pages 234–243. Springer, 2005.
- [12] Jorge Cortes, Sonia Martinez, Timur Karatas, and Francesco Bullo. Coverage control for mobile sensing networks. *IEEE Transactions on robotics and Automation*, 20(2):243–255, 2004.
- [13] Maryam Fazel, Rong Ge, Sham M Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. *arXiv preprint arXiv:1801.05039*, 2018.
- [14] Jakob Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 122–130. International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- [15] Rainer Gemulla, Erik Nijkamp, Peter J Haas, and Yannis Sismanis. Large-scale matrix factorization with distributed stochastic gradient descent. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 69–77. ACM, 2011.
- [16] Davood Hajinezhad, Mingyi Hong, and Alfredo Garcia. Zeroth order nonconvex multi-agent optimization over networks. *arXiv preprint arXiv:1710.09997*, 2017.
- [17] Mingyi Hong, Zhi-Quan Luo, and Meisam Razaviyayn. Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems. *SIAM Journal on Optimization*, 26(1):337–364, 2016.

- [18] Mingyi Hong, Zhi-Quan Luo, and Mesiam Razaviyayn. Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems. *CONVERGENCE*, 26(1): 337–364, 2016.
- [19] Mingyi Hong, Davood Hajinezhad, and Ming-Min Zhao. Prox-pda: The proximal primal-dual algorithm for fast distributed nonconvex optimization and learning over networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1529–1538. JMLR. org, 2017.
- [20] Junling Hu and Michael P Wellman. Nash q-learning for general-sum stochastic games. *Journal of machine learning research*, 4(Nov):1039–1069, 2003.
- [21] Minyi Huang, Peter E Caines, and Roland P Malhamé. Individual and mass behaviour in large population stochastic wireless power control problems: centralized and nash equilibrium solutions. In *42nd IEEE International Conference on Decision and Control (IEEE Cat. No. 03CH37475)*, volume 1, pages 98–103. IEEE, 2003.
- [22] Minyi Huang, Roland P Malhamé, Peter E Caines, et al. Large population stochastic dynamic games: closed-loop mckean-vlasov systems and the nash certainty equivalence principle. *Communications in Information & Systems*, 6(3):221–252, 2006.
- [23] Edward Hughes, Joel Z Leibo, Matthew Phillips, Karl Tuyls, Edgar Dueñez-Guzman, Antonio García Castañeda, Iain Dunning, Tina Zhu, Kevin McKee, Raphael Koster, et al. Inequity aversion improves cooperation in intertemporal social dilemmas. In *Advances in Neural Information Processing Systems*, pages 3326–3336, 2018.
- [24] Dušan Jakovetić, José MF Moura, and Joao Xavier. Linear convergence rate of a class of distributed augmented lagrangian algorithms. *IEEE Transactions on Automatic Control*, 60(4): 922–936, 2014.
- [25] Zhanhong Jiang, Aditya Balu, Chinmay Hegde, and Soumik Sarkar. Collaborative deep learning in fixed topology networks. In *Advances in Neural Information Processing Systems*, pages 5904–5914, 2017.
- [26] Raghavendra V Kulkarni and Ganesh Kumar Venayagamoorthy. Particle swarm optimization in wireless-sensor networks: A brief survey. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 41(2):262–267, 2010.
- [27] Jean-Michel Lasry and Pierre-Louis Lions. Games à middle field. ii - finite horizon and control ô the optimal. *Math é Math Accounts*, 343(10):679–684, 2006.
- [28] Jean-Michel Lasry and Pierre-Louis Lions. Mean field games. *Japanese journal of mathematics*, 2(1):229–260, 2007.
- [29] Martin Lauer and Martin Riedmiller. An algorithm for distributed reinforcement learning in cooperative multi-agent systems. In *In Proceedings of the Seventeenth International Conference on Machine Learning*. Citeseer, 2000.

- [30] Donghwan Lee, Hyungjin Yoon, and Naira Hovakimyan. Primal-dual algorithm for distributed reinforcement learning: distributed gtd. In *2018 IEEE Conference on Decision and Control (CDC)*, pages 1967–1972. IEEE, 2018.
- [31] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-agent reinforcement learning in sequential social dilemmas. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, pages 464–473. International Foundation for Autonomous Agents and Multiagent Systems, 2017.
- [32] Laurent Lessard and Sanjay Lall. Optimal control of two-player systems with output feedback. *IEEE Transactions on Automatic Control*, 60(8):2129–2144, 2015.
- [33] Yingying Li, Yujie Tang, Runyu Zhang, and Na Li. Distributed reinforcement learning for decentralized linear quadratic control: A derivative-free policy optimization approach. *arXiv preprint arXiv:1912.09135*, 2019.
- [34] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 5330–5340, 2017.
- [35] Gabriel M Lipsa and Nuno C Martins. Remote state estimation with communication costs for first-order lti systems. *IEEE Transactions on Automatic Control*, 56(9):2013–2025, 2011.
- [36] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine learning proceedings 1994*, pages 157–163. Elsevier, 1994.
- [37] Michael L Littman. Value-function reinforcement learning in markov games. *Cognitive Systems Research*, 2(1):55–66, 2001.
- [38] Steven KC Lo. A collaborative multi-agent message transmission mechanism in intelligent transportation system—a smart freeway example. *Information Sciences*, 184(1):246–265, 2012.
- [39] Daria Madjidian and Leonid Mirkin. Distributed control with low-rank coordination. *IEEE Transactions on Control of Network Systems*, 1(1):53–63, 2014.
- [40] Dhruv Malik, Ashwin Pananjady, Kush Bhatia, Koulik Khamaru, Peter L Bartlett, and Martin J Wainwright. Derivative-free methods for policy optimization: Guarantees for linear quadratic systems. *arXiv preprint arXiv:1812.08305*, 2018.
- [41] Laetitia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. Independent reinforcement learners in cooperative markov games: a survey regarding coordination problems. *The Knowledge Engineering Review*, 27(1):1–31, 2012.
- [42] Liviu Panait and Sean Luke. Cooperative multi-agent learning: The state of the art. *Autonomous agents and multi-agent systems*, 11(3):387–434, 2005.
- [43] James A Preiss, Sébastien MR Arnold, Chen-Yu Wei, and Marius Kloft. Analyzing the variance of policy gradient estimators for the linear-quadratic regulator. *arXiv preprint arXiv:1910.01249*, 2019.

- [44] Michael Rabbat and Robert Nowak. Distributed optimization in sensor networks. In *Proceedings of the 3rd international symposium on Information processing in sensor networks*, pages 20–27. ACM, 2004.
- [45] Dominic Richards and Patrick Rebeschini. Optimal statistical rates for decentralised non-parametric regression with linear speed-up. In *Advances in Neural Information Processing Systems*, pages 1214–1225, 2019.
- [46] R Tyrrell Rockafellar. Augmented lagrangians and applications of the proximal point algorithm in convex programming. *Mathematics of operations research*, 1(2):97–116, 1976.
- [47] William H Sandholm. *Population games and evolutionary dynamics*. MIT press, 2010.
- [48] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- [49] Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [50] Yoav Shoham, Rob Powers, and Trond Grenager. Multi-agent reinforcement learning: a critical survey. *Web manuscript*, 2003.
- [51] Yoav Shoham, Rob Powers, and Trond Grenager. If multi-agent learning is the answer, what is the question? *Artificial Intelligence*, 171(7):365–377, 2007.
- [52] H Eugene Stanley. *Phase transitions and critical phenomena*. Clarendon Press, Oxford, 1971.
- [53] Haoran Sun and Mingyi Hong. Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms. In *2018 52nd Asilomar Conference on Signals, Systems, and Computers*, pages 38–42. IEEE, 2018.
- [54] 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.
- [55] Stephen Tu and Benjamin Recht. Least-squares temporal difference learning for the linear quadratic regulator. *arXiv preprint arXiv:1712.08642*, 2017.
- [56] Hirofumi Uzawa. Iterative methods for concave programming. *Studies in linear and nonlinear programming*, 6:154–165, 1958.
- [57] Hoi-To Wai, Tsung-Hui Chang, and Anna Scaglione. A consensus-based decentralized algorithm for non-convex optimization with application to dictionary learning. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3546–3550. IEEE, 2015.
- [58] Hoi-To Wai, Zhuoran Yang, Princeton Zhaoran Wang, and Mingyi Hong. Multi-agent reinforcement learning via double averaging primal-dual optimization. In *Advances in Neural Information Processing Systems*, pages 9649–9660, 2018.

- [59] Stephen J Wright. Implementing proximal point methods for linear programming. *Journal of optimization Theory and Applications*, 65(3):531–554, 1990.
- [60] Bin Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of Probability*, pages 94–116, 1994.
- [61] Kaiqing Zhang, Zhuoran Yang, and Tamer Basar. Networked multi-agent reinforcement learning in continuous spaces. In *2018 IEEE Conference on Decision and Control (CDC)*, pages 2771–2776. IEEE, 2018.
- [62] Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Başar. Fully decentralized multi-agent reinforcement learning with networked agents. *arXiv preprint arXiv:1802.08757*, 2018.
- [63] M Zhu and S Martinez. An approximate dual subgradient algorithm for distributed cooperative non-convex constrained optimization. *IEEE Transactions on Automatic Control*, submitted, 2010.

Appendix A. Policy Gradient with Reparametrized States

In this section, we provide the details for reparametrizing states in mean-field dynamics in Section 3. We characterize the optimal cost from a state going forward (see Eq. (A.6)) with (algebraic) Riccati equations [4] governed by policy parameters (*i.e.*, K, L). Note that in the formulation of the previous section, it is hard to directly obtain the optimal control $u_t^{(i)}$'s depending on both $x_t^{(i)}$'s and \bar{x}_t 's correlated in a single dynamics. Hence we adopt the following reparameterization method to derive decoupled Riccati equations.

$$u_t^{(i)} = Kx_t^{(i)} + L\bar{x}_t = M(x_t^{(i)} - \bar{x}_t) + N\bar{x}_t \triangleq My_t^{(i)} + N\bar{y}_t. \quad (\text{A.1})$$

where we call $M = K$ the optimal *local policy* and $N = K + L$ the optimal *mean-field policy*. Below we show how to derive policy gradient to update these policy parameters separately. We rewrite the dynamic equations below:

$$y_{t+1}^{(i)} = (A + BM)y_t^{(i)} + w_t^{(i)} - \bar{w}_t, \quad \bar{y}_{t+1} = [A + \bar{A} + BN]\bar{y}_t + \bar{w}_t, \quad (\text{A.2})$$

where $\bar{w}_t = 1/n \sum_{i=1}^n w_t^{(i)}$. Then we have

$$\begin{aligned} C(\Theta) &= \mathbb{E}_{\mathbf{y}_0 \sim \mathcal{D}, \mathbf{w}} \sum_{t=0}^{\infty} \gamma^t \left(\sum_{i=1}^n y_t^{(i)\top} (Q + M^\top RM) y_t^{(i)} + n \bar{y}_t^\top (Q + \bar{Q} + N^\top RN) \bar{y}_t \right) \\ &= \sum_{i=1}^n \text{Tr}[(Q + M^\top RM) \Sigma_M^{(i)}] + n \text{Tr}[(\bar{Q} + N^\top RN) \Sigma_N], \end{aligned} \quad (\text{A.3})$$

where $\Sigma_M^{(i)} = \mathbb{E}_{\mathbf{y}_0 \sim \mathcal{D}, \mathbf{w}} [\sum_{t=0}^{\infty} \gamma^t y_t^{(i)} y_t^{(i)\top}]$, $\Sigma_N = \mathbb{E}_{\bar{\mathbf{y}}_0 \sim \bar{\mathcal{D}}, \mathbf{w}} [\sum_{t=0}^{\infty} \gamma^t \bar{y}_t \bar{y}_t^\top]$, $\bar{Q} = Q + \bar{Q}$, \mathcal{D} denotes the *i.i.d.* initial state distribution of $y_0^{(i)}$, and $\bar{\mathcal{D}}$ indicates the distribution of mean-field \bar{y}_0 . We will omit the explicit distributions for expectations without ambiguity. With an abuse of notations, let $C(M^{(i)}) = \text{Tr}[(Q + M^\top RM) \Sigma_M^{(i)}]$, $C(N) = n \text{Tr}[(\bar{Q} + N^\top RN) \Sigma_N]$, $\underline{C}(M) = \min_{i \in [n]} \text{Tr}[(Q + M^\top RM) \Sigma_M^{(i)}]$. From Bellman equation for cost functions, it follows that

$$\begin{aligned} C_\Theta(\mathbf{y}_0) &= \sum_{i=1}^n y_0^{(i)\top} (Q + M^\top RM) y_0^{(i)} + n \bar{y}_0^\top \bar{P}_N \bar{y}_0 \\ &\quad + \gamma \mathbb{E}_{w_0} C_\Theta(\text{diag}(A + BM) \mathbf{y}_0 + \mathbf{w}_0 - \bar{\mathbf{w}}_0), \end{aligned} \quad (\text{A.4})$$

where $C_\Theta(y)$ is the value function with initial state y , and $\text{diag}(X)$ denotes the $nd \times nd$ block-diagonal matrix with matrix $X \in \mathbb{R}^d$ being the block elements. We assume that the value function takes a quadratic form $C_\Theta(\mathbf{y}_0) = \sum_{i=1}^n y_0^{(i)\top} P_M y_0^{(i)} + n \bar{y}_0^\top \bar{P}_N \bar{y}_0 + \alpha_\Theta$, where $P_M, \bar{P}_N \in \mathbb{R}^{d \times d}$ and we ignore the cross terms as $\mathbb{E}[y_0^\top \bar{P}_N \bar{y}_0] = 0$. Let $A_y = A + BM$, $\bar{A}_y = A + \bar{A} + BN$, then by Bellman equation (A.4), we have

$$\begin{aligned} \sum_{i=1}^n y_0^{(i)\top} P_M y_0^{(i)} + n \bar{y}_0^\top \bar{P}_N \bar{y}_0 + \alpha_\Theta &= \sum_{i=1}^n y_0^{(i)\top} [\gamma A_y^\top P_M A_y + Q + M^\top RM] y_0^{(i)} \\ &\quad + n \bar{y}_0^\top [Q + \bar{Q} + \bar{N} RN + \gamma \bar{A}_y^\top \bar{P}_N \bar{A}_y] \bar{y}_0 + \gamma \frac{n+1}{n} \text{Tr} P_M + \gamma \text{Tr} \bar{P}_N + \gamma \alpha_\Theta, \end{aligned} \quad (\text{A.5})$$

which implies

$$\begin{aligned} P_M &= Q + M^\top RM + \gamma A_y^\top P_M A_y, \\ \alpha_\Theta &= \frac{\gamma}{1-\gamma} \left(\frac{n+1}{n} \text{Tr} P_M + \bar{P}_N \right), \\ \bar{P}_N &= Q + \bar{Q} + N^\top RN + \gamma \bar{A}_y^\top \bar{P}_N \bar{A}_y. \end{aligned} \quad (\text{A.6})$$

Hence M and N are decoupled for Riccati equations after transformation. The gradient with respect to M in (A.4) gives

$$\nabla_M C_\Theta(\mathbf{y}_0) = (2RM + 2\gamma B^\top P_M A_y) \sum_{i=1}^n y_0^{(i)} y_0^{(i)\top} + \gamma \mathbb{E} \nabla_M C_\Theta(\mathbf{y}_1). \quad (\text{A.7})$$

By applying recursion $y_{t+1} = (A + BM)y_t + w_t^{(i)} - \bar{w}_t$ ($t \geq 0$) to (A.7) iteratively, we can finally obtain

$$\begin{aligned} \nabla_M C_\Theta(\mathbf{y}_0) &= 2 \left(RM + \gamma B^\top P_M A_y \right) \sum_{i=1}^n \mathbb{E}_{\mathbf{y}_0, \mathbf{w}} \sum_{t=0}^{\infty} \gamma^t y_t^{(i)} y_t^{(i)\top} \\ &= 2 \left(RM + \gamma B^\top P_M A_y \right) \sum_{i=1}^n \Sigma_M^{(i)} \\ &\triangleq 2 \Xi_M \sum_{i=1}^n \Sigma_M^{(i)}. \end{aligned} \quad (\text{A.8})$$

Similarly, we have

$$\nabla_N C_\Theta(\mathbf{y}_0) = 2 \left(RN + \gamma B^\top \bar{P}_N \bar{A}_y \right) \Sigma_N \triangleq 2 \Xi_N \Sigma_N. \quad (\text{A.9})$$

The Σ_N in $\nabla_N C_\Theta(\mathbf{y}_0)$ is a shared expectation term for all the agents resulted from mean-field states.

Appendix B. Experiment Setup and Additional Details

In this section, We provide additional configuration details and analysis of the experimental results in Section 6.

Experiment setup. In the experiment we intend to demonstrate the convergence performance of our algorithm under different graph structures and different dynamics. We consider policy learning with global consensus as follows:

$$\min_{\tilde{\Theta}} \tilde{C}(\tilde{\Theta}) := \frac{1}{n} \sum_{i=1}^n C(\tilde{\Theta}^{(i)}), \quad \text{s.t.} \quad \tilde{\Theta}^{(i)} = \tilde{\Theta}^{(j)} \quad \text{for} \quad (i, j) \in \mathcal{E}, \quad (\text{B.1})$$

where for all i ,

$$x_{t+1}^{(i)} = Ax_t^{(i)} + Bu_t^{(i)} + \bar{A}\bar{x}_t + w_t^{(i)}, \quad (\text{B.2})$$

$$c_t^{(i)} = x_t^{(i)\top} Q x_t^{(i)} + u_t^{(i)\top} R u_t^{(i)} + \bar{x}_t^\top \bar{Q} \bar{x}_t. \quad (\text{B.3})$$

The state transition matrix A, B, \bar{A} is generated by first sampling a random uniform matrix from 0 to 1 and then tossing a biased coin with probability 0.7 for setting each element zero to keep sparsity for computational efficiency. Also, such transition matrices may lead to smoother loss surfaces. For the reward function we adopt diagonal matrices with each element on the diagonal from 0 to 1 and sparse perturbation for off-diagonal entries. We test on four different graphs as the communication network \mathcal{G} : 1) complete graph 2) grid graph 3) circle graph and 4) random graph. The random graph is essentially an Erdos-Rényi graph generated with connectivity of 0.25. Alternatively, for each pair (i, j) with $i \in [n], j \in [n]$, we toss a coin to decide whether there will be an edge between agent i and agent j . As we toss coins two times for both (i, j) and (j, i) , an equivalent random graph is obtained with 0.75 probability for each edge to vanish. We set $\Omega^2 = \Gamma^2 = I$ as identity matrices to better understand the popular average of neighborhood scheme. The results of convergence curves are presented in Figure 1(a)- 1(c), where y-axis denotes the error measure $\varepsilon(t) = \min_{s \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(\tilde{\Theta}_s^{(i)}) - \tilde{C}(\tilde{\Theta}^*) \right| + \|\Omega \tilde{\mathcal{L}}\{\tilde{\Theta}_k\}_{(1)}\|^2$ for $t \in \bar{\mathcal{N}}$.

For the impact of the topology of the communication graphs on the convergence rate, when choosing Ω and Γ according to Section 4, we can easily verify that the condition (4.1) holds, so that constant C of the bound in Theorem 4.2 can be reparameterized as follows,

$$C \leq \frac{320 \max\{\sigma_{\max}(Z), 1\}}{\min\{\underline{\sigma}_{\min}(\mathcal{L}_G), 1\}} \sum_{i \sim j} \left(\frac{\sqrt{\beta_i \beta_j}}{\sqrt{d_i d_j n}} + \frac{\bar{\beta}}{4} \right), \quad (\text{B.4})$$

where $\bar{\beta} = 1/n \sum_{i=1}^n \beta_i$; see more details in [53]. Hence the global rate is connected to the algebraic summary, *i.e.*, spectral gap, that captures the connectivity of the communication network, which quantifies how the connectivity of four different graphs impacts convergence performance of MF-DPGM.

For the impact of the dimension of dynamics on the convergence, high dimensions have a higher chance to introduce high variance in $\Sigma_0^{(i)}$ depending on the initial random states, and different cost dynamics also account a change in α_g . In addition, we observe that performance gap between complete graph and other two deterministic graphs increases in higher dimension configuration, which implies an interplay between communication structure (\mathcal{C}) and system dynamics (α_g) in the convergence rate.

We consider two scenarios, fixed and varying outdoor temperatures, and use the same training and testing settings similar to what is used in [33]. For the fixed outdoor temperature scenario, the outdoor temperature is set to 30; for the varying one, the outdoor temperature is given by some real temperature history data. To train the decentralized controller in the varying outdoor temperature scenario, we follow the setting of [33] and simultaneously optimize on 16 environments with different fixed outdoor temperatures (temperatures given in [33]). The training under the fixed outdoor temperature scenario is similar except that we no longer need to sample from the environments as the outdoor temperature is given. At each policy update step, we randomly sample one environment to train the controller on. We use standard SGD as the optimizer. We set the learning rate to 10^{-7} and batch size to 10. We do early stopping based on the average training cost. At test time, we run the trained controller with the test outdoor temperature data given in [33]. The results are shown in Figure 3 and Figure 4, in which the black dotted line is the outdoor temperature. For the fixed outdoor temperature scenario, we plot the test trajectory generated by models trained with different number of iterations to show how the performance improves with respect to the number of iterations.

Appendix C. Proof Sketch

We sketch the proof of results in Section 4. Define $\Sigma_0^{(i)} := \mathbb{E} y_0^{(i)} y_0^{(i)\top}$, $\Sigma_0 := \mathbb{E} \bar{y}_0 \bar{y}_0^\top$, $\xi_i := \sigma_{\min}(\Sigma_0^{(i)})$, $\bar{\xi} = \sigma_{\min}(\mathbb{E}_{y_0^{(i)} \sim \mathcal{D}} \bar{y}_0 \bar{y}_0^\top)$, where $\sigma_{\min}(X)$ denotes the smallest singular value of X . We denote by $\underline{\sigma}_{\min}(X)$ the second smallest eigenvalue of X .

Geometry of cost functions. As mentioned in Section 4, Theorem 4.2 requires moderate smoothness of the landscape of cost functions. Based on the almost Lipschitzness of positive definite matrix P_θ parameterizing the optimal cost from a state going forward, and almost Lipschitzness of the Σ_θ which plays an key role in cost function. Lemma D.2 for each cost function with an almost Lipschitz gradient and Lemma D.3 for each cost function almost dominated by the gradient are crucial in bounding the one-step differences in Lemma D.6 and D.7.

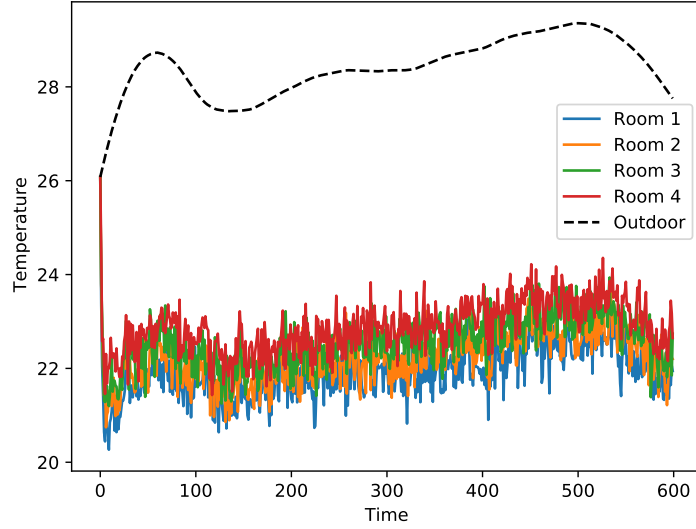


Figure 3: Performance of the trained decentralized controller on HVAC dynamics with varying temperatures.

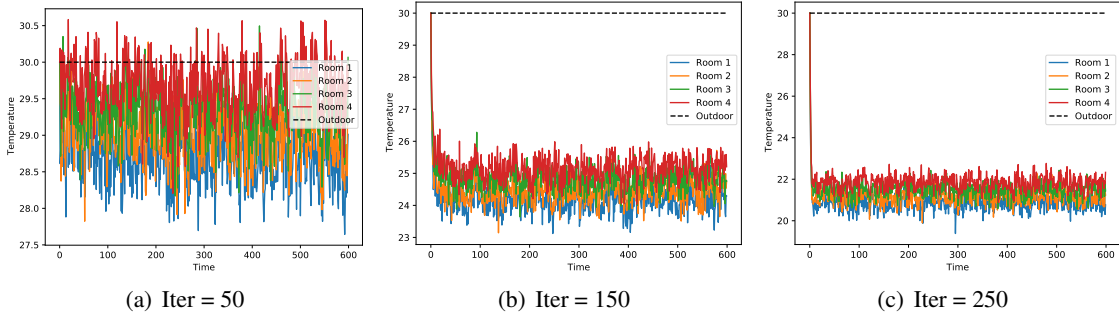


Figure 4: Performance of the trained decentralized controller on HVAC dynamics with fixed temperature (Temp = 30).

Progress of decentralized iterations. To estimate one-step progress of Algorithm 1, we construct the auxiliary potential function U to show certain monotonicity along the solution path in both primal ($\tilde{\Theta}$) and dual (Λ) variables:

$$U_{k+1} = U_c \left(\tilde{\Theta}_{k+1}, \tilde{\Theta}_k, \Lambda_{k+1} \right) := \mathcal{A}_{k+1} + \frac{2\kappa}{n^2} \left\| \Gamma^{-1} \Phi \left(\tilde{\Theta}_{k+1} - \tilde{\Theta}_k \right) \right\|^2 \quad (\text{C.1})$$

$$+ \frac{c}{2} \left(\left\| \Omega \tilde{\mathcal{L}} \{ \tilde{\Theta}_{k+1} \}_{(1)} \right\|^2 + \left\| \tilde{\Theta}_{k+1} - \tilde{\Theta}_k \right\|_{H + \tilde{\mathcal{L}}/n}^2 \right),$$

where c is a constant chosen according to Condition 4.1 and \mathcal{A}_{k+1} is the augmented Lagrangian defined in (3.6). The decrement of the potential function at each step is identified with a metric between two adjacent primal iterates in the following lemma.

Lemma C.1. When the parameters of MF-DPGM are chosen to satisfy (4.1) then it holds that

$$U_k - U_{k+1} \geq \frac{1}{4} \left\| \tilde{\Theta}_{k+1} - \tilde{\Theta}_k \right\|_{\tilde{\Omega} + \Gamma^2}^2 + \kappa \|\mathbb{V}_{k+1}\|_H^2 \quad (\text{C.2})$$

for any $k \geq 0$, where $\mathbb{V}_{k+1} := \left(\tilde{\Theta}_{k+1} - \tilde{\Theta}_k \right) - \left(\tilde{\Theta}_k - \tilde{\Theta}_{k-1} \right)$. Moreover, we have the following bounds:

$$U_k \leq U_0 \leq \tilde{C}(\tilde{\Theta}_0) + \frac{2\nabla\tilde{C}(\mathbf{0})^\top \Phi^{-1} \nabla\tilde{C}(\mathbf{0})}{n}, \quad U_{k+1} \geq \underline{C} > -\infty \quad (\text{C.3})$$

for any $k \geq 1$, where $\mathbf{0} \in \mathbb{R}^{nmd}$ is the all zero vector, and $\nabla\tilde{C}(\mathbf{0}) = 1/n(\nabla C^{(1)}(\mathbf{0}), \dots, \nabla C^{(n)}(\mathbf{0}))$.

Proof. Detailed proof can be found in the appendix of [53]. \square

Note that decreasing potential function in Lemma C.1 also tracks stability of decentralized LQR in the optimization process. We start by optimality condition (3.8) and derive upper bounds for objective gradient norms by the distances of primal variables. Then Lemma C.1 is applied reducing the bounds to differences of adjacent potential functions. Similarly, the consensus error is controlled using Lemma D.6 and D.7, where almost smoothness of costs (D.10) is involved, and processed by Lemma C.1 to keep the same difference terms as those of gradient norms. Finally, combining two bounds of similar structure with Lemma D.3 we establish Theorem 4.2. See D.5 for a detailed proof.

Appendix D. Detailed Proof of Main Results

In this section, we develop detailed proofs for the main result in Section 4 and give complementary details for theoretical claims.

In the sequel, due to the similarity between evolutions of local policies M , and mean-field policy N , we mainly focus on $M^{(i)}$'s to state and prove the results, where similar results are straightforward up to constants in \bar{A} , \bar{Q} , etc. In some cases, to stress on the discrepancy we formulate both illustrations, or to establish a unified higher-level convergence results, we use compact tensor / vectorization representations, such as $\tilde{\Theta}$, \mathbf{M} for statements. Also, we omit superscript for agent i without confusion in a specific proof.

We first define the following operators on symmetric matrix X ,

$$\begin{aligned} \mathcal{H}_M(X) &= \sum_{t=0}^{\infty} \gamma^t (A + BM)^t X [(A + BM)^\top]^t, \\ \mathcal{H}_N(X) &= \sum_{t=0}^{\infty} \gamma^t (A + \bar{A} + BN)^t X [(A + \bar{A} + BN)^\top]^t. \end{aligned} \quad (\text{D.1})$$

We also recall that for $X \in \mathbb{R}^{d \times m}$, we define

$$\mathcal{F}_M(X) = \min \{ \xi \sigma_{\min}(Q) / [4(\|A + BX\| + 1) \cdot \|B\| \cdot C(X)], \|X\| \}, \quad (\text{D.2})$$

$$\mathcal{F}_N(X) = \min \{ \bar{\xi} \sigma_{\min}(Q + \bar{Q}) / [4(\|A + BX\| + 1) \cdot \|B\| \cdot C(X)], \|X\| \}, \quad (\text{D.3})$$

which are frequently used notations to simplify our conditions and proof. For notation convenience, define $\xi := \inf_{i \in [n]} [\sigma_{\min}(\mathbb{E}_{y_0^{(i)} \sim \mathcal{D}} y_0^{(i)} y_0^{(i)\top})]$, $\bar{\xi} := \sigma_{\min}(\mathbb{E}_{y_0^{(i)} \sim \mathcal{D}} \bar{y}_0 \bar{y}_0^\top)$, where $\sigma_{\min}(X)$ refers to the smallest singular value of matrix X . In addition, we define

$$\Sigma_0^{(i)} \triangleq \mathbb{E} y_0^{(i)} y_0^{(i)\top}, \quad \|\mathcal{H}_M\| \triangleq \sup_X \frac{\|\mathcal{H}_M(X)\|}{\|X\|}. \quad (\text{D.4})$$

It follows that the operator norms are bounded by the composite cost and the extremal singular values of cost matrices.

Lemma D.1 (Upper bounds of operators \mathcal{H}_M and \mathcal{H}_N). It holds that

$$\|\mathcal{H}_M\| \leq \frac{\underline{C}(M)}{\xi \sigma_{\min}(Q)}, \quad \|\mathcal{H}_N\| \leq \frac{C(N)}{\bar{\xi} \sigma_{\min}(Q + \bar{Q})}. \quad (\text{D.5})$$

Using the operator norm bounds and definitions, we show the continuity property of the cost functions, state trajectories, and gradients of cost corresponding to policy \mathbb{M} and \mathbb{N} with an adaptive area for each agent. Lemma D.4-D.2 quantify the problem geometry from different perspectives in distributed settings.

Proof. By the definition in (D.4) and $\Sigma_0 = \mathbb{E} \bar{y}_0 \bar{y}_0^\top$, we can obtain for any $i \in [n]$,

$$\mathcal{H}_M(\Sigma_0^{(i)}) = \Sigma_M^{(i)}, \quad \mathcal{H}_N(\Sigma_0) = \Sigma_N \quad (\text{D.6})$$

By the definition of the operator norm, for $x \in \mathbb{R}^d$ of unit vector norm and matrix X of unit spectral norm, for any $i \in [n]$ we have

$$\begin{aligned} x^\top (\mathcal{H}_M(X)) x &= \sum_{t=0}^{\infty} \text{Tr}([(A + BM)^\top]^t x x^\top (A + BM)^t X) \\ &= \sum_{t=0}^{\infty} \text{Tr}(\Sigma_0^{(i)1/2} [(A + BM)^\top]^t x x^\top (A + BM)^t \Sigma_0^{(i)1/2} \Sigma_0^{(i)-1/2} X \Sigma_0^{(i)-1/2}) \\ &\leq \sum_{t=0}^{\infty} \text{Tr}(\Sigma_0^{(i)1/2} [(A + BM)^\top]^t x x^\top (A + BM)^t \Sigma_0^{(i)1/2}) \|\Sigma_0^{(i)-1/2} X \Sigma_0^{(i)1/2}\| \\ &= x^\top \mathcal{H}_M(\Sigma_0^{(i)}) \|\Sigma_0^{(i)-1/2} X \Sigma_0^{(i)1/2}\| \\ &\stackrel{(a)}{\leq} \frac{\|\mathcal{H}_M(\Sigma_0^{(i)})\|}{\sigma_{\min}(\mathbb{E} x_0^{(i)} x_0^{(i)\top})} \\ &= \frac{\|\Sigma_M^{(i)}\|}{\xi}, \end{aligned} \quad (\text{D.7})$$

where (a) uses the property $\|\Sigma_0^{(i)}\| \geq \sigma_{\min}(\Sigma_0^{(i)})$. On the other hand, we can derive an upper bound on $\|\Sigma_M^{(i)}\|$ as follows

$$\begin{aligned} \|\Sigma_M^{(i)}\| &\leq \text{Tr}(\Sigma_M^{(i)}) \leq \frac{\text{Tr}(\Sigma_M^{(i)}) \sigma_{\min}(Q)}{\sigma_{\min}(Q)} \\ &\leq \frac{\text{Tr}(\Sigma_M^{(i)}(Q + M^\top R M))}{\sigma_{\min}(Q)} \\ &= \frac{C^{(i)}(M)}{\sigma_{\min}(Q)}. \end{aligned} \quad (\text{D.8})$$

Combining (D.17) and (D.18) and applying uniform lower bound of $C^{(i)}(M)$'s we have $\|\mathcal{H}_M\| \leq \frac{\underline{C}(M)}{\xi \sigma_{\min}(Q)}$. Similar computation gives the upper bound for the norm of \mathcal{H}_N . \square

D.1 Main Lemmas for the Geometry of Cost Functions

According to Section C, appropriate smoothness of the landscape of cost functions is studied to characterize convergence rates. The following lemma for each cost function with an almost Lipschitz

gradient, which is a result of the almost Lipschitzness of positive definite matrix P_θ and almost Lipschitzness of the Σ_θ , shows significance in bounding the one-step differences in Lemma D.6 and D.7.

Lemma D.2. (Almost smoothness of private cost functions) Assume that for each $i \in [n]$ and any $\widehat{M}^{(i)}, M^{(i)} \in \mathbb{R}^{m \times d}$ it holds that

$$\begin{aligned}\|\widehat{M}^{(i)} - M^{(i)}\| &\leq \mathcal{F}_M(M^{(i)}), \\ \|\widehat{N}^{(i)} - N^{(i)}\| &\leq \mathcal{F}_N(N^{(i)}),\end{aligned}\tag{D.9}$$

then for the i -th agent, we have

$$\begin{aligned}\|\nabla C(\widehat{M}^{(i)}) - \nabla C(M^{(i)})\| &\leq \beta_i^M \|\widehat{M}^{(i)} - M^{(i)}\|, \\ \|\nabla C(\widehat{N}^{(i)}) - \nabla C(N^{(i)})\| &\leq \beta_i^N \|\widehat{N}^{(i)} - N^{(i)}\|,\end{aligned}\tag{D.10}$$

where $\beta_i^M = \text{poly}\left(\mathcal{B}, \mathbb{E}\|y_0^{(i)}\|^2, \frac{C(M_0^{(i)})}{\xi_i \sigma_{\min}(Q)}\right)$, $\beta_i^N = \text{poly}\left(\mathcal{B}, \mathbb{E}\|\bar{y}_0\|^2, \frac{C(N_0^{(i)})}{\xi \sigma_{\min}(Q+Q)}\right)$ denotes the almost smoothness constants, and $\mathcal{B} = \{\|A\|, \|B\|, \|R\|, \sigma_{\min}^{-1}(R)\}$.

Proof. See Lemma 6 in [13] for a detailed proof. \square

Another crucial property to guarantee the global convergence of MF-DPGM is the gradient domination condition, where the difference of the current cost and optimal cost is bounded by the current gradient norm. We conclude this landscape in \mathbb{M} and \mathbb{N} for each agent in the lemma below.

Lemma D.3. (Gradient domination of cost functions) Suppose $(M^*; N^*)$ is the optimal policy for each agent, and $\Sigma_0^{(i)}$ is full rank. Then $C(M^{(i)}), C(N^{(i)})$ is gradient dominated for each i , that is,

$$\begin{aligned}C(M^{(i)}) - C(M^*) &\leq \alpha_g^M \|\nabla C(M^{(i)})\|^2, \\ C(N^{(i)}) - C(N^*) &\leq \alpha_g^N \|\nabla C(N^{(i)})\|^2,\end{aligned}\tag{D.11}$$

where α_g^M and α_g^N are geometry-dependent coefficients specified in Theorem 4.2.

Proof. See Appendix D.4 for a detailed proof. \square

As $\Sigma_M^{(i)} \succeq \Sigma_0^i$, the full-rank condition essentially prevents the denominator of α_g^M from going to zero, so that a stationary point ($\nabla C(M^{(i)}) = 0$) on the R.H.S. of (D.11) implies an optimal policy $M^{(i)}$. Although $\Sigma_{N^{(i)}} \succeq \Sigma_0$, the difference from single-agent setting is the absence of the assumption for Σ_0 . In fact, we have $\Sigma_0 = 1/n^2 \mathbb{E}\left(\sum y_0^{(i)}\right)\left(\sum y_0^{(i)}\right)^\top = 1/n^2 \mathbb{E}\sum_{i,j} y_0^{(i)} y_0^{(j)\top} = \mathbb{E}y_0^{(i)} y_0^{(i)\top}$, where *i.i.d.* initial state distributions and linearity of expectation are used. Hence, the only condition in Lemma D.3 suffices to guarantee gradient domination for both local and mean-field policies. Such detail reveals additional advantages of condition relaxation from mean-field symmetry besides dimensionality reduction.

D.2 Lemmas for Almost-Smoothness of Cost Functions

In this subsection, we present two almost continuity lemmas on which Lemma D.2 is based. One of them is the following almost Lipschitz continuity of P_θ parameterizing the cost functions.

Lemma D.4 (Almost Lipschitzness of P_θ (value function)). For any two $M^{(i)}$ and $M^{(i)'}$ close enough to each other, that is,

$$\|M^{(i)} - M^{(i)'}\| \leq \min \left\{ \frac{\xi_i \sigma_{\min}(Q)}{4(\|A + BM^{(i)}\| + 1)\|B\|C(M^{(i)})}, \|M^{(i)}\| \right\}, \quad (\text{D.12})$$

then

$$\|P_{M^{(i)'}} - P_{M^{(i)}}\| \leq \frac{6\|M^{(i)}\|\|R\|}{\xi_i^2 \sigma_{\min}^2(Q)} (\|M^{(i)}\|\|B\|\|A + BM^{(i)}\| + \|M^{(i)}\|\|B\| + 1) \cdot \|M^{(i)'} - M^{(i)}\|. \quad (\text{D.13})$$

Similarly, if

$$\|N^{(i)'} - N^{(i)}\| \leq \min \left\{ \frac{\bar{\xi} \sigma_{\min}(Q + \bar{Q})}{4(\|A + \bar{A} + BN^{(i)}\| + 1)\|B\|C(N^{(i)})}, \|N^{(i)}\| \right\}, \quad (\text{D.14})$$

then we have

$$\|P'_{N^{(i)}} - P_{N^{(i)}}\| \leq \frac{6\|N^{(i)}\|\|R\|}{\bar{\xi}^2 \sigma_{\min}^2(Q + \bar{Q})} (\|N^{(i)}\|\|B\|\|A + \bar{A} + BN^{(i)}\| + \|N^{(i)}\|\|B\| + 1) \|N^{(i)'} - N^{(i)}\|.$$

Proof. We first define the following operators on symmetric matrix X ,

$$\begin{aligned} \mathcal{H}_M(X) &= \sum_{t=0}^{\infty} \gamma^t (A + BM)^t X [(A + BM)^\top]^t, \\ \mathcal{H}_N(X) &= \sum_{t=0}^{\infty} \gamma^t (A + \bar{A} + BN)^t X [(A + \bar{A} + BN)^\top]^t, \\ \mathcal{J}_M(X) &= \gamma^t (A + BM) X (A + BM)^\top, \\ \mathcal{J}_N(X) &= \gamma^t (A + \bar{A} + BN) X (A + \bar{A} + BN)^\top. \end{aligned} \quad (\text{D.15})$$

Then we can rewrite the difference between P_M and $P_{M'}$ as

$$\begin{aligned} &\|P_{M'} - P_M\| \\ &= \|\mathcal{H}_{M'}(Q + M'^\top RM') - \mathcal{H}_M(Q + M^\top RM)\| \\ &\leq \|\mathcal{H}_{M'}(Q + M'^\top RM') - \mathcal{H}_M(Q + M'^\top RM') - (\mathcal{H}_M(Q + M^\top RM) - \mathcal{H}_M(Q + M'^\top RM'))\| \\ &\leq 2\|\mathcal{H}_M\|^2 \|\mathcal{J}_M - \mathcal{J}_{M'}\| \|(M')^\top RM'\| + \|\mathcal{H}_M\| \|M^\top RM - (M')^\top RM'\| \\ &\stackrel{(a)}{\leq} \|\mathcal{H}_M\| \left(\|(M')^\top RM' - M^\top RM\| + 2\|\mathcal{H}_M\| \|\mathcal{J}_M - \mathcal{J}_{M'}\| \|M^\top RM\| \right) \\ &+ \|\mathcal{H}_M\| \|M^\top RM - (M')^\top RM'\| \\ &= 2\|\mathcal{H}_M\|^2 \|\mathcal{J}_M - \mathcal{J}_{M'}\| \|M^\top RM\| + 2\|\mathcal{H}_M\| \|(M')^\top RM' - M^\top RM\|, \end{aligned} \quad (\text{D.16})$$

where (a) uses the triangle inequality for ℓ_2 -norm, and the assumption $\|\mathcal{H}_M\| \|\mathcal{J}_M - \mathcal{J}_{M'}\| \leq 1/2$ for the coefficient of $\|(M')^\top RM' - M^\top RM\|$. To bound the first term in (D.16), letting $\delta = M - M'$, we take the following decomposition for $\|\mathcal{J}_M - \mathcal{J}_{M'}\|$ for each matrix X ,

$$\begin{aligned} \|(\mathcal{J}_M - \mathcal{J}_{M'})(X)\| &= \|(A + BM)X(B\delta)^\top + (B\delta)X(A + BM)^\top - (B\delta)X(B\delta)^\top\| \\ &\leq 2\|(A + BM)\| \|X\| \|B\| \|\delta\| + \|B\|^2 \|\delta\|^2 \|X\|. \end{aligned} \quad (\text{D.17})$$

According to the definition of the spectral norm and the assumed condition on $\|M - M'\|$ (D.15), we are able to bound the first term as below.

$$\begin{aligned}
& 2\|\mathcal{H}_M\|^2\|\mathcal{J}_M - \mathcal{J}_{M'}\|\|M^\top RM\| \\
& \leq 2\|\mathcal{H}_M\|^2(2\|(A + BM)\|\|B\|\|M - M'\| + \|B\|^2\|M - M'\|^2)\|M^\top RM\| \\
& \leq 4\|\mathcal{H}_M\|^2\|B\|\|M - M'\| \left(\|(A + BM)\| + \frac{\sigma_{\min}(Q)\xi}{8C(\Theta)(\|A + BM\| + 1)} \right) \|M^\top RM\| \\
& \leq 4\|\mathcal{H}_M\|^2\|B\|(\|(A + BM)\| + 1)\|M^\top RM\|\|M - M'\|. \tag{D.18}
\end{aligned}$$

Note that $\|M' - M\| \leq \|M\|$, the second term in (D.16) can be bounded as By plugging (D.17) and (D.18) into (D.16), we can finally obtain the almost Lipschitzness result for P_M . Similarly, we can also derive an argument for P_N as (D.16) below: Then applying a slightly different upper bound for $\|\mathcal{H}_N\|$ lead to the result. \square

The next lemma quantifies a Lipschitz continuity condition for $\Sigma_\theta^{(i)}$. Due to the policy gradient structure, it plays an important role in bounding a part of the gradient difference of cost functions.

Lemma D.5 (Almost-Lipschitzness of Σ_θ). For each $i \in [n]$, if the following holds

$$\|M^{(i)} - M'^{(i)}\| \leq \left\{ \frac{\sigma_{\min}(Q)\xi_i}{4C(M^{(i)})\|B\|(\|A + BM^{(i)}\| + 1)}, \|M^{(i)}\| \right\}, \tag{D.19}$$

it follows that

$$\left\| \Sigma_{M'}^{(i)} - \Sigma_M^{(i)} \right\| \leq 4 \left(\frac{C(M^{(i)})}{\sigma_{\min}(Q)} \right)^2 \frac{\|B\|(\|A + BM^{(i)}\| + 1)}{\xi_i} \|M^{(i)} - M'^{(i)}\|. \tag{D.20}$$

Also, when

$$\|N^{(i)} - N'^{(i)}\| \leq \left\{ \frac{\sigma_{\min}(Q + \bar{Q})\bar{\xi}}{4C(N^{(i)})\|B\|(\|A + \bar{A} + BN^{(i)}\| + 1)}, \|N^{(i)}\| \right\}, \tag{D.21}$$

we have

$$\left\| \Sigma_{N'}^{(i)} - \Sigma_N^{(i)} \right\| \leq 4 \left(\frac{C(N^{(i)})}{\sigma_{\min}(Q + \bar{Q})} \right)^2 \frac{\|B\|(\|A - BN^{(i)}\| + 1)}{\bar{\xi}} \|N^{(i)} - N'^{(i)}\|. \tag{D.22}$$

Proof. See [13] for a detailed proof. \square

D.3 Adaptive Choice of Parameters Ω and Γ

In this section, we provide guidance to choose edge associated parameter Ω and agent associated parameter Γ in MF-DPGM algorithm, in order to meet the conditions of the adaptive area of (D.12) and (D.14). According to the communication and update step in Algorithm 1, for $M_t^{(i)}$ we have

$$\begin{aligned}
\|M_{t+1}^{(i)} - M_t^{(i)}\| &= \frac{1}{2\sum_{j:j\sim i}\sigma_{ij}^2 + \gamma_i^2} \left\| \frac{1}{n} \left(\nabla C(M_t^{(i)}) - \nabla C(M_{t-1}^{(i)}) \right) - 2 \sum_{j:j\sim i} \sigma_{ij}^2 M_t^{(j)} \right. \\
&\quad \left. + \gamma_i^2 \left(M_{t-1}^{(i)} - M_t^{(i)} \right) + \sum_{j:j\sim i} \sigma_{ij}^2 \left(M_{t-1}^{(j)} + M_{t-1}^{(i)} \right) \right\|. \tag{D.23}
\end{aligned}$$

By Theorem 6.1 in [53] and Lemma 28 [13] we obtain the desired stepsize in (4.2).

D.4 Proof of Lemma D.3

Now we proceed to prove the gradient domination lemma based on last two almost Lipschitzness results for cost functions, which is essential to control the one-step progress of MF-DPGM in both primal and dual variables.

Proof. By (A.8) we can split the left-hand-side of (D.10) into two terms

$$\begin{aligned} \|\nabla C(\widehat{M}^{(i)}) - \nabla C(M^{(i)})\|_F &= \|2\Xi_{\widehat{M}^{(i)}}\Sigma_{M^{(i)}} - 2\Xi_{M^{(i)}}\Sigma_{M^{(i)}}\|_F \\ &\leq 2\|\Xi_{M^{(i)}}(\Sigma_{\widehat{M}^{(i)}} - \Sigma_{M^{(i)}})\| + 2\|(\Xi_{\widehat{M}^{(i)}} - \Xi_{M^{(i)}})\Sigma_{\widehat{M}^{(i)}}\|. \end{aligned} \quad (\text{D.24})$$

Let \widehat{y}_t and \widehat{u}_t be the sequence induced by $\widehat{M}^{(i)}$. Note that $C(M^{*(i)}) \leq C(\widehat{M}^{(i)})$. Let $V_M(y) = \mathbb{E}_{\mathbf{w}} y^\top P_M y$, $Q_M(y, u) = y^\top Q y + u^\top R u + V_M((A + BM)y + \widehat{w})$, $A_M(y, u) = Q_M(y, u) - V_M(y)$. By using cost difference lemma [13] and Lemma D.4 and D.5 we have

$$C(M) - C(M^*) \geq C(M) - C(\widehat{M}) \quad (\text{D.25})$$

$$= -\mathbb{E} \sum_t A_M(\widehat{y}_t, \widehat{u}_t) \quad (\text{D.26})$$

$$= \mathbb{E} \sum_t \text{Tr} \left(\widehat{y}_t \widehat{y}_t^\top \Xi_M^\top (R + B^\top P_M B)^{-1} \Xi_M \right) \quad (\text{D.27})$$

$$\geq \text{Tr} \left(\Sigma_{\widehat{M}} \Xi_M^\top (R + B^\top P_M B)^{-1} \Xi_M \right) \quad (\text{D.28})$$

$$\geq \frac{\xi}{\|R + B^\top P_M B\|} \text{Tr} \left(\Xi_M^\top \Xi_M \right). \quad (\text{D.29})$$

Hence we have the norm bound

$$\|\Xi_M\|_F^2 \leq \frac{\|R + B^\top P_M B\|}{\xi} (C(M) - C(M^*)). \quad (\text{D.30})$$

Then, for the first term in (D.24), we adopt Lemma D.5 to derive the upper bound. For the latter term, we note that $\|\Sigma_{\widehat{M}}\| \leq \|\Sigma_M\| + \frac{C(M)}{\sigma_{\min}(\widehat{Q})}$ due to small norm of $\widehat{M} - M$. Combining with Lemma D.4 we can obtain the final bound in $\|\widehat{M} - M\|$. \square

Next we turn to formulate the one-step progress of dual variable controlling the consensus error in the following lemma.

Lemma D.6 (One-step progress of dual variable). For any $k \in \mathbb{N}$, it holds that

$$\|\Lambda_{k+1} - \Lambda_k\| \leq 2\kappa \left(\frac{\left\| \Gamma^{-1} \Phi \left(\widetilde{\Theta}_k - \widetilde{\Theta}_{k-1} \right) \right\|^2}{n^2} + \|\mathbb{V}_{k+1}\|_H^2 \right), \quad (\text{D.31})$$

where

$$\begin{aligned} \kappa &:= \frac{1}{\lambda_{\min}(\Omega F H^{-1} F^\top \Omega)}, \\ \mathbb{V}_{k+1} &:= \left(\widetilde{\Theta}_{k+1} - \widetilde{\Theta}_k \right) - \left(\widetilde{\Theta}_k - \widetilde{\Theta}_{k-1} \right). \end{aligned} \quad (\text{D.32})$$

Proof. See [53] for a detailed proof. \square

By this lemma, we are able to transform difference norms of dual variables into primal variables. \mathbb{V} also keep a second order difference corresponding to the requirement of past gradients and policies

by MF-DPGM. Then, we introduce the progress of augmented Lagrangian, which captures the dynamics of both primal and dual variables.

Lemma D.7 (One-step progress of augmented Lagrangian function). For all $k \geq 0$, the iterates in MF-DPGM gives

$$\begin{aligned} \mathcal{A}_{k+1} - \mathcal{A}_k &\leq -\frac{1}{2} \left\| \tilde{\Theta}_{k+1} - \tilde{\Theta}_k \right\|_{\tilde{\Omega} + 2\Gamma^2 - \Phi/n}^2 \\ &+ \kappa \left(\frac{2}{n^2} \left\| \Gamma^{-1} \Phi \left(\tilde{\Theta}_k - \tilde{\Theta}_{k-1} \right) \right\|^2 + 2 \left\| \mathbb{V}_{k+1} \right\|_H^2 \right). \end{aligned} \quad (\text{D.33})$$

Proof. See [53] for a detailed proof. \square

Again the progress bound of augmented Lagrangian is parameterized by primal first-order differences and second order differences. Converting to such uniform differences is helpful in the proof of the main theorem.

With all the lemmas above in place, now we are ready to prove the global convergence result of our novel decentralized MARL algorithm.

D.5 Proof of Theorem 4.2

Proof. Let $\mathbf{M} \in \mathbb{R}^{nmd}$ be the vectorization of tensor parameter \mathbb{M} , $\mathbf{1}$ be the block all one vector with block vectors in \mathbb{R}^{md} . From the update of the algorithm, we have the following optimality condition that for any $k \geq -1$, it holds that

$$\langle \mathbf{1}, \nabla \tilde{C}(\mathbf{M}_k) \rangle + \langle \mathbf{1}, H(\mathbf{M}_{k+1} - \mathbf{M}_k) \rangle = 0. \quad (\text{D.34})$$

By taking the square of both sides, we can obtain

$$\left\| \frac{1}{n} \sum_{i=1}^n \nabla C(\mathbf{M}_k^{(i)}) \right\|^2 = |\mathbf{1}^\top H(\mathbf{M}_{k+1} - \mathbf{M}_k)|^2. \quad (\text{D.35})$$

Using Cauchy-Schwarz inequality under metric matrix H , we have

$$\begin{aligned} |\mathbf{1}^\top H(\mathbf{M}_{k+1} - \mathbf{M}_k)|^2 &\leq \|\mathbf{M}_{k+1} - \mathbf{M}_k\|_H^2 \|\mathbf{1}\|_H^2 \\ &= (\mathbf{M}_{k+1} - \mathbf{M}_k)^\top H(\mathbf{M}_{k+1} - \mathbf{M}_k) \cdot \mathbf{1}^\top H \mathbf{1} \\ &\stackrel{(b)}{\leq} \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right) \|\mathbf{M}_{k+1} - \mathbf{M}_k\|_H^2, \end{aligned} \quad (\text{D.36})$$

when (b) result from the definition of H . Then we combine Lemma C.1, (D.35), and (D.36) to get

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\nabla C(\mathbf{M}_k^{(i)})\|_F^2 &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla C(\mathbf{M}_k^{(i)}) \right\|^2 \\ &\leq \|\mathbf{M}_{k+1} - \mathbf{M}_k\|_H^2 \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right) \\ &\leq 8 \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right) (U_k - U_{k+1}), \end{aligned} \quad (\text{D.37})$$

where we also use the fact that $H \preceq 2(\tilde{\Omega} + \Gamma^2)$.

For the error caused by constraint violence (inexact consensus), we know from Lemma D.6 and parameter setting of Γ in Lemma C.1 that

$$\begin{aligned}\|\Omega\tilde{\mathcal{L}}\mathbf{M}_{k+1}\|^2 &\leq \kappa \left(2\mathbf{V}_{k+1}^\top H\mathbf{V}_{k+1} + \frac{2}{n^2}\|\Gamma^{-1}\Phi(\mathbf{M}_{k+1} - \mathbf{M}_k)\|^2 \right) \\ &\leq 4\kappa \left(\frac{1}{n^2}\|\Gamma^{-1}\Phi(\mathbf{M}_{k+1} - \mathbf{M}_k)\|^2 + 2\|\mathbf{V}_{k+1}\|_H^2 \right).\end{aligned}\quad (\text{D.38})$$

Then, we look one step back and using Jensen's inequality to bound similar term for step k ,

$$\begin{aligned}\|\Omega\tilde{\mathcal{L}}\mathbf{M}_k\|^2 &\leq 2 \left(\|\Omega\tilde{\mathcal{L}}\mathbf{M}_{k+1}\|^2 + \|\Omega\tilde{\mathcal{L}}(\mathbf{M}_{k+1} - \mathbf{M}_k)\|^2 \right) \\ &\leq 8\kappa \left(\frac{1}{n^2}\|\Gamma^{-1}\Phi(\mathbf{M}_{k+1} - \mathbf{M}_k)\|^2 + 2\|\mathbf{V}_{k+1}\|_H^2 \right) + 2\|\Omega\tilde{\mathcal{L}}(\mathbf{M}_{k+1} - \mathbf{M}_k)\|^2 \\ &= 8\kappa \left(\frac{1}{n^2}\|\mathbf{M}_{k+1} - \mathbf{M}_k\|_{\Phi\Gamma^{-2}\Phi}^2 + 2\|\mathbf{V}_{k+1}\|_H^2 \right) + 2\|\mathbf{M}_{k+1} - \mathbf{M}_k\|_{\tilde{\mathcal{L}}\Omega^2\tilde{\mathcal{L}}}^2.\end{aligned}\quad (\text{D.39})$$

From the constraints in (4.1), we have

$$\Phi\Gamma^{-2}\Phi \preceq \frac{n^2}{8\kappa}(\tilde{\Omega} + \Gamma^2).\quad (\text{D.40})$$

Meanwhile, the definition of $\tilde{\Omega}$ gives

$$\tilde{\mathcal{L}}\Omega^2\tilde{\mathcal{L}} \preceq 2\tilde{\Omega}.\quad (\text{D.41})$$

Again using the step improvement in potential function U combined with (D.40) and (D.41), we have

$$\begin{aligned}\|\Omega\tilde{\mathcal{L}}\mathbf{M}_k\|^2 &\leq \|\mathbf{M}_{k+1} - \mathbf{M}_k\|_{\tilde{\Omega}+\Gamma^2}^2 + 16\kappa\|\mathbf{V}_{k+1}\|_H^2 + 4\|\mathbf{M}_{k+1} - \mathbf{M}_k\|_{\tilde{\Omega}}^2 \\ &\leq 5\|\mathbf{M}_{k+1} - \mathbf{M}_k\|_{\tilde{\Omega}+\Gamma^2}^2 + 16\kappa\|\mathbf{V}_{k+1}\|_H^2 \\ &\leq 20(U_k - U_{k+1}).\end{aligned}\quad (\text{D.42})$$

On the other hand, according to Lemma D.3 and the measure (left hand side of (4.3)) for convergence rate, we can derive the following inequality on the cost error,

$$\begin{aligned}t \cdot \min_{k \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(M_k^{(i)}) - C(M^*) \right| &\leq \sum_{k=1}^t \left| \frac{1}{n} \sum_{i=1}^n C(M_k^{(i)}) - C(M^*) \right| \\ &\leq \sum_{k=1}^t \frac{1}{n} \left(\sum_{i=1}^n |C(M_k^{(i)}) - C(M^*)| \right) \\ &\stackrel{(\text{D.11})}{\leq} \sum_{k=1}^t \frac{\|\Sigma_{M^*}\|}{n\sigma_{\min}(\Sigma)^2\sigma_{\min}(R)} \sum_{i=1}^n \|\nabla C(M_k^{(i)})\|_F^2 \\ &\stackrel{(\text{D.37})}{\leq} \sum_{k=1}^t \frac{8\|\Sigma_{M^*}\|(U_k - U_{k+1})}{\sigma_{\min}(\Sigma)^2\sigma_{\min}(R)} \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right).\end{aligned}\quad (\text{D.43})$$

We note that the summation indexed by k only operates on difference terms of adjacent potential functions $(U_k - U_{k+1})$'s. Combining the above results with the upper bound and lower bound of the

potential function in Lemma C.1, we have

$$\begin{aligned} \min_{k \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(M_k^{(i)}) - C(M^*) \right| &\leq \frac{8 \|\Sigma_{M^*}\| (U_1 - U_{t+1})}{t \sigma_{\min}(\Sigma)^2 \sigma_{\min}(R)} \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right) \\ &\leq \frac{8 \|\Sigma_{M^*}\| (U_0 - \inf_{\mathbb{M}} \tilde{C}(\mathbb{M}))}{t \sigma_{\min}(\Sigma)^2 \sigma_{\min}(R)} \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right). \end{aligned} \quad (\text{D.44})$$

Similarly, we can directly derive the consensus error bound in constants and t from (D.42) as follows,

$$\begin{aligned} \min_{k \in [t]} \|\Omega \tilde{\mathcal{L}} \mathbf{M}_k\|^2 &\leq \frac{1}{t} \cdot \sum_{k=1}^t \|\Omega \tilde{\mathcal{L}} \mathbf{M}_k\|^2 \\ &\leq \frac{20}{t} \sum_{k=1}^t (U_k - U_{k+1}) \\ &\leq \frac{20(U_0 - U_{t+1})}{t} \\ &\stackrel{(\text{C.3})}{\leq} \frac{20}{t} \left(\tilde{C}(\mathbb{M}_0) - \inf_{\mathbb{M}} \tilde{C}(\mathbb{M}) + \frac{2 \nabla C(\mathbf{0})^\top \Phi^{-1} \nabla C(\mathbf{0})}{n} \right). \end{aligned} \quad (\text{D.45})$$

Therefore, we have attained the cost error and consensus error bound respectively in some problem coonstants. As the first inequalities of both derivations come from the same argument, we can finally obtain the overall convergence rate as

$$\begin{aligned} \min_{s \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(M_s^{(i)}) - C(M^*) \right| + \|\Omega \tilde{\mathcal{L}} \mathbf{M}_k\|^2 &\leq \underbrace{\frac{8 \|\Sigma_{M^*}\| U_0 - \inf_{\mathbb{M}} \tilde{C}(\mathbb{M})}{\sigma_{\min}(\Sigma)^2 \sigma_{\min}(R)} \left(4 \sum_{(i,j), i \sim j} \sigma_{ij}^2 + \sum_{i=1}^n \gamma_i^2 \right)}_{\text{cost error bound}} \\ &\quad + \underbrace{\frac{20}{t} \left(\tilde{C}(\mathbb{M}_0) - \inf_{\mathbb{M}} \tilde{C}(\mathbb{M}) + \frac{2 \nabla C(\mathbf{0})^\top \Phi^{-1} \nabla C(\mathbf{0})}{n} \right)}_{\text{consensus error bound}} \\ &\leq \frac{2\mathcal{C}'}{t} (5 + 4\alpha_g \mathcal{C}). \end{aligned} \quad (\text{D.46})$$

Since $\alpha_g = \max\{\alpha_g^M, \alpha_g^N\}$ and we choose the commonly applied constant when using inequalities for \mathbf{M} and \mathbf{N} , we complete the proof for overall variables. \square

Analysis of model-free policy gradient estimator. Our main analysis is based on the exact policy gradient $\nabla C(\mathbb{M})$, while in practice we adopt an empirical version $\widehat{\nabla} C(\mathbb{M})$ for updates. As mentioned in Section 3.2, our proof can be easily tweaked to include the variance incurred by $\widehat{\nabla} C(\mathbb{M})$. Specifically, the objective on the left-hand side of (D.43) is changed into $C(\widehat{M}_k^{(i)})$ where

$\widehat{M}_k^{(i)}$ is the iterate obtained by unbiased estimated policy gradient, which gives

$$\begin{aligned}
& t \cdot \min_{k \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(\widehat{M}_k^{(i)}) - C(M^*) \right| \\
& \leq t \cdot \min_{k \in [t]} \left| \frac{1}{n} \sum_{i=1}^n (C(\widehat{M}_k^{(i)}) - C(M_k^{(i)})) \right| + t \cdot \min_{k \in [t]} \left| \frac{1}{n} \sum_{i=1}^n C(M_k^{(i)}) - C(M^*) \right|, \quad (\text{D.47})
\end{aligned}$$

where the first term can be bounded according to the almost Lipschitzness of the cost function by the $\|\widehat{M}_k^{(i)} - M_k^{(i)}\| = \|\widehat{M}_k^{(i)} - \mathbb{E}\widehat{M}_k^{(i)}\| = \mathbf{\Pi} \cdot \|\widehat{\nabla}C(M_k^{(i)}) - \mathbb{E}\widehat{\nabla}C(M_k^{(i)})\|$, where $\mathbf{\Pi}$ denotes the step-size of the policy gradient in Algorithm 1. Such standard variance can be further bounded by the variance of the REINFORCE estimator in [43]. On the other hand, the second term can still be bounded following the proof above. Combining these bounds we obtain the error bound with estimated policy gradients.

Appendix E. Analysis of Computation and Communication Complexities

In this section, we briefly conclude the storage and computation resources and communication overhead during running MF-DPGM.

Firstly, for each agent the method requires previously computed gradients from its own and past simulated states from the neighborhood. Therefore, each agent needs to store two gradient tensors and two policy tensors to avoid the computational overhead brought by re-evaluating trajectory-based policy gradients and states, which is more space efficient than another policy evaluation method [58]. To conclude, the whole system is supposed to store $(2nmd + 2nmd)$ real numbers at any iteration. From a view of the update step for each agent, each step involves summation of $m \times d$ matrices by number of neighbors, leading to an $\mathcal{O}(d_i md)$ computation complexity for agent i . On the other hand, according to the information exchanging round in the communication and update step, MF-DPGM requires $\mathcal{O}(2e)$ communications at each round and each exchange delivers md real numbers. Although the overhead is greatly alleviated compared to centralized scheme, it still cost much bandwidth when the policy is extremely complicated to parameterize. We are trying to infer state information from part of the neighborhood to further reduce communication as the future work.