

A PRINCIPLED PERMUTATION INVARIANT APPROACH TO MEAN-FIELD MULTI-AGENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-agent reinforcement learning (MARL) becomes more challenging in the presence of more agents, as the capacity of the joint state and action spaces grows exponentially in the number of agents. To address such a challenge of scale, we identify a class of cooperative MARL problems with permutation invariance, and formulate it as mean-field Markov decision processes (MDP). To exploit the permutation invariance therein, we propose the mean-field proximal policy optimization (MF-PPO) algorithm, at the core of which is a permutation-invariant actor-critic neural architecture. We prove that MF-PPO attains the globally optimal policy at a sublinear rate of convergence. Moreover, its sample complexity is independent of the number of agents. We validate the theoretical advantages of MF-PPO with numerical experiments in the multi-agent particle environment (MPE). In particular, we show that the inductive bias introduced by the permutation-invariant neural architecture enables MF-PPO to outperform existing competitors with a smaller number of model parameters, which is the key to its generalization performance.

1 INTRODUCTION

Multi-Agent Reinforcement Learning (Littman, 1994; Zhang et al., 2019) generalizes Reinforcement Learning (Sutton and Barto, 2018) to address the sequential decision-making problem of multiple agents maximizing their individual long term rewards while interacting with each other in a common environment. With breakthroughs in deep learning, MARL algorithms equipped with deep neural networks have seen significant empirical successes in various domains, including simulated autonomous driving (Shalev-Shwartz et al., 2016), multi-agent robotic control (Matarić, 1997; Kober et al., 2013), and E-sports (Vinyals et al., 2019).

Despite tremendous successes, MARL is notoriously hard to scale to the many-agent setting, as the size of the state-action space grows exponentially with respect to the number of agents. This phenomenon is recently described as the curse of many agents (Menda et al., 2018). To tackle this challenge, we focus on *cooperative MARL*, where agents work together to maximize their team reward (Panait and Luke, 2005). We identify and exploit a key property of cooperative MARL with homogeneous agents, namely the *invariance with respect to the permutation of agents*. Such permutation invariance can be found in many real-world scenarios with homogeneous agents, such as distributed control of multiple autonomous vehicles and team sports (Cao et al., 2013; Kalyanakrishnan et al., 2006), but also in scenarios with heterogeneous agent groups, where invariance holds within each group (Liu et al., 2019b). More importantly, we find that permutation invariance has significant practical implications, as the optimal value functions remain invariant when permuting the joint state-action pairs. Such an observation strongly advocates a permutation invariant design for learning, which helps reduce the effective search space of the policy/value functions from exponential dependence on the number of agents to polynomial dependence.

Several empirical methods have been proposed to incorporate permutation invariance into solving MARL problems. Liu et al. (2019b) implement a permutation invariant critic based on Graph Convolutional Network (GCN) (Kipf and Welling, 2017). Sunehag et al. (2017) propose value decomposition, which together with parameter sharing, leads to a joint critic network that is permutation invariant over agents. While these methods are based on heuristics, we are the first to provide theoretical principles for introducing permutation invariance as an inductive bias for learning value functions and policies in homogeneous systems. In addition, we adopt the DeepSet (Zaheer et al., 2017)

architecture, which is well suited for handling homogeneity of agents, with much simpler operations to induce permutation invariance and greater parameter efficiency.

To scale MARL algorithms in the presence of a large number, even infinitely many, agents, mean-field approximation has been explored to directly model the population behavior of the agents. Mean-field game considers large populations of rational agents that play a noncooperative game. Yang et al. (2017) consider a mean-field game with deterministic linear state transitions, and show that it can be reformulated as a mean-field MDP, where the mean-field state lies in finite-dimensional probability simplex. Yang et al. (2018) take a mean-field approximation over actions, such that the interaction for any given agent and the population is approximated by the interaction between the agent’s action and the averaged actions of its neighboring agents. However, the motivation for averaging over local actions remains unclear, and it generally requires a sparse graph over agents. In practice, properly identifying such structure also demands extensive prior knowledge. Mean-field control instead considers a central controller who aims to compute strategy to optimize the average payoff across the population. Carmona et al. (2019) motivate a mean-field MDP from the perspective of mean-field control. The mean-field state therein lies in a probability simplex and is thus continuous in nature. To enable the ensuing Q-learning algorithm, discretization of the joint state-action space is necessary. In addition, the dynamic programming principles of such mean-field control problem has been studied in (Gu et al., 2019). Gu et al. (2020) also propose a Q-learning type algorithm, where the state-action space is first discretized into an epsilon-net. The kernel regression operator is used to construct an estimate of the unknown Q-function from samples. Gu et al. (2021) propose a localized training, decentralized execution framework by locally grouping homogenous agents using their states. Wang et al. (2020) motivate a mean-field MDP from permutation invariance, but assume a central controller coordinating the actions of all the agents, and hence is restricted to handling the curse of many agents from the exponential blowup of the joint state space. Our formulation of mean-field approximation allows agents to make their own local actions without resorting to a centralized controller.

We propose a mean-field Markov decision process motivated from the permutation invariance structure of cooperative MARL, which can be viewed as a natural limit of finite-agent MDP by taking the number of agents to infinity. Such a mean-field MDP generalizes traditional MDP, with each state representing a distribution over the state space of a single agent. The mean-field MDP provides us a tractable formulation to model MDP with many agents, including an infinite number of agents. We further propose the Mean-Field Proximal Policy Optimization (MF-PPO) algorithm, at the core of which is a pair of permutation invariant actor and critic neural networks. These networks are implemented based on DeepSet (Zaheer et al., 2017), which uses convolutional type operations to induce permutation invariance over the set of inputs. We show that with sufficiently many agents, MF-PPO converges to the optimal policy of the mean-field MDP with a sublinear sample complexity independent of the number of agents. To support our theory, we conduct numerical experiments on the benchmark multi-agent particle environment (MPE) and show that our proposed method requires a smaller number of model parameters and attains better performance than multiple baselines.

Notations. We denote $\mathcal{P}(X)$ as the set of distribution on set X . δ_x denotes the Dirac measure supported at x . For $\mathbf{s} = (s_1, \dots, s_N)$, we use $\mathbf{s} \stackrel{\text{i.i.d.}}{\sim} p$ to denote that each s_i is independently sampled from distribution p . For $f : X \rightarrow \mathbb{R}$ and a distribution $\pi \in \mathcal{P}(X)$, we write $\langle f, \pi \rangle = \mathbb{E}_{a \sim \pi} f(a)$. We write $[m]$ in short for $\{1, \dots, m\}$, and Δ_d for the standard probability simplex in \mathbb{R}^d .

2 PROBLEM SETUP

We focus on studying multi-agent systems with cooperative, homogeneous agents, where the agents within the system are of similar nature and hence cannot be distinguished from each other. Specifically, we consider a discrete time control problem with N agents, formulated as a Markov decision process $(\mathcal{S}^N, \mathcal{A}^N, \mathbb{P}, r)$. We define the joint state space \mathcal{S}^N to be the Cartesian product of the finite state space \mathcal{S} for each agent, and similarly define the joint action space \mathcal{A}^N . The homogeneous nature of the system is reflected in the transition kernel \mathbb{P} and the shared reward r , which satisfies:

$$r(\mathbf{s}_t, \mathbf{a}_t) = r(\kappa(\mathbf{s}_t), \kappa(\mathbf{a}_t)), \quad \mathbb{P}(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) = \mathbb{P}(\kappa(\mathbf{s}_{t+1}) | \kappa(\mathbf{s}_t), \kappa(\mathbf{a}_t)) \quad (2.1)$$

for all $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S}^N \times \mathcal{A}^N$ and the permutation mapping $\kappa(\cdot) \in \mathbb{S}_N$, where \mathbb{S}_N is the set of all one-to-one mapping from $[N]$ to itself. In other words, it is the configuration, rather than individual identities, that affects the team reward, and the transition to the next configuration solely depends on the current configuration. See Figure 1 for detailed illustration. Such permutation invariance finds applications in many real-world scenarios, including distributed control of autonomous vehicles, and social economic systems (Zheng et al., 2020; Cao et al., 2013; Kalyanakrishnan et al., 2006).

Our goal is to find the optimal policy ν , where $\nu(\mathbf{s}) \in \Delta_{|\mathcal{A}^N|}$ for all $\mathbf{s} \in \mathcal{S}^N$, and maximize the expected discounted reward $V^\nu(\mathbf{s}) = (1 - \gamma)\mathbb{E}\{\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) | \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_t \sim \nu(\mathbf{s}_t), \forall t \geq 0\}$. Our first result shows that learning with permutation invariance advocates invariant network design.

Proposition 2.1. *For cooperative MARL satisfying (2.1), there exists an optimal policy ν^* that is permutation invariant, i.e., $\nu^*(\mathbf{s}, \mathbf{a}) = \nu^*(\kappa(\mathbf{s}), \kappa(\mathbf{a}))$ for any permutation mapping $\kappa(\cdot)$. In addition, for any permutation invariant policy ν , the value function $V(\cdot)$ and the state-action value function $Q(\cdot)$ is also permutation invariant, i.e., $V^\nu(\mathbf{s}) = V^\nu(\kappa(\mathbf{s}))$, $Q^\nu(\mathbf{s}, \mathbf{a}) = Q^\nu(\kappa(\mathbf{s}), \kappa(\mathbf{a}))$, where $Q^\nu(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}'}\{r(\mathbf{s}, \mathbf{a}) + \gamma V^\nu(\mathbf{s}')\}$.*

Proposition 2.1 has an important implication for architecture design, as it states that it suffices to search within the permutation invariant policy and value function classes. To the best of our knowledge, this is the first theoretical justification of permutation invariant network design for learning with homogeneous agents.

We focus on the factorized policy class with a parameter sharing scheme, where each agent makes its own decision without consolidating with others. Specifically, the joint policy ν can be factorized as $\nu(\mathbf{a}|\mathbf{s}) = \prod_{i=1}^N \mu(a_i|o_i)$, where $\mu(\cdot)$ denotes the shared local mapping and o_i denotes the local observation. Such a policy class is widely adopted in the celebrated centralized training – decentralized execution paradigm (Lowe et al., 2017), due to its light overhead in the deployment phase and favorable performances. However, directly learning such factorized policy remains challenging, as each agent needs to estimate its state-action value function, denoted as $Q^\nu(\mathbf{s}, \mathbf{a})$. The search space during learning is $(|\mathcal{S}| \times |\mathcal{A}|)^N$, scaling exponentially with respect to the number of agents. The large search space poses as a significant roadblock for efficient learning, and is coined as the *curse of many agents*.

To address the curse of many agents, we exploit the homogeneity of the system and take the mean-field approximation. We begin by taking the perspective of agent i , which is arbitrarily chosen from the N agents. We denote its state as s and the states of the rest of the agents by \mathbf{s}_r . One can verify that when permuting the state of all the other agents, the value function remains unchanged; additionally, we can further characterize the value function as a function of the local state and the empirical state distribution over the rest of agents.

Proposition 2.2. *For any permutation mapping $\kappa(\cdot)$, the value function satisfies $V^\nu(s, \mathbf{s}_r) = V^\nu(s, \kappa(\mathbf{s}_r))$. Additionally, there exists g_ν such that: $V^\nu(s, \mathbf{s}_r) = g_\nu(s, \hat{\mathbf{p}}_{\mathbf{s}_r})$, where $\hat{\mathbf{p}}_{\mathbf{s}_r} = \frac{1}{N} \sum_{s \in \mathbf{s}_r} \delta_s$ is the empirical distribution over the states of rest of the agents \mathbf{s}_r .*

For a system with a large number of agents (e.g., financial markets, social networks), the empirical state distribution can be seen as the concrete realization of the underlying population distribution of the agents. Motivated from this observation and Proposition 2.2, we formulate the following mean-field MDP that can be seen as the limit of finite-agent MDP in the presence of infinitely many homogeneous agents.

Definition 2.1 (mean-field MDP). *The mean-field MDP consists of elements of the following: state $(s, d_S) \in \mathcal{S} \times \mathcal{P}(\mathcal{S})$; action $\bar{a} \in \bar{\mathcal{A}} \subseteq \mathcal{A}^S$; reward $r(s, d_S, \bar{a})$; transition kernel $\mathbb{P}(s', d'_S | s, d_S, \bar{a})$.*

The mean-field MDP has an intimate connection with our previously discussed finite-agent MDP. Since the agents are homogeneous, the system is the same from any agent’s perspective. We choose any agent (referred to as representative agent), the state information of such an agent includes the local state s , and the mean-field state d_S . With state information, the agent selects a meta action $\bar{a} \in \bar{\mathcal{A}} \subseteq \mathcal{A}^S$, and uses such a meta action to make local decision $a = \bar{a}(s) \in \mathcal{A}$. We remark that such a modeling of decision process allows the agent to make decision on both its local information (local state s) and the global information (mean-field state d_S). From homogeneity we assume all the rest of the agents uses the same meta action \bar{a} to make their local actions. Note that different agents can still make different local actions due to their different local states, i.e., $\bar{a}(z) \neq \bar{a}(z')$ in general for

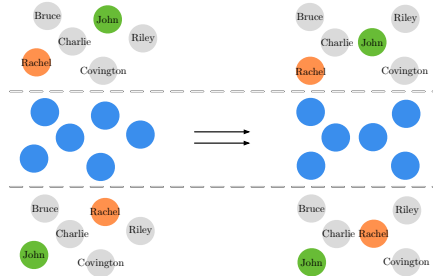


Figure 1: Illustration of permutation invariance. Exchanging identities of the agents (first and third row) does not change the transition of the system (second row).

$z \neq z' \in \mathcal{S}$. The joint state at the next timestep (s', d'_S) naturally depends on the current global state (s, d_S) and the meta action \bar{a} (since all the other agents use \bar{a} to generate their local actions), and is specified by the transition kernel $P(s', d'_S | s, d_S, \bar{a})$. In addition, the representative agents receives a reward $r(s, d_S, \bar{a})$, which depends on the local state and mean-field state, and the meta action \bar{a} .

Our goal is to learn efficiently a policy π , where $\pi(\cdot | s, d_S) \in \Delta_{|\bar{\mathcal{A}}|}$ for all $(s, d_S) \in \mathcal{S} \times \mathcal{P}(\mathcal{S})$, for maximized expected discounted reward. To facilitate discussions, we define the value function $V^\pi(s, d_S) = (1 - \gamma) \mathbb{E} \{ \sum_{t=0}^{\infty} \gamma^t r(s_t, d_{S,t}, \bar{a}_t) \}$, where $(s_0, d_{S,0}) = (s, d_S)$, $\bar{a}_t \sim \pi(s_t, d_{S,t})$, $\forall t \geq 0$; and Q-function $Q^\pi(s, d_S, \bar{a}) = (1 - \gamma) \mathbb{E} \{ \sum_{t=0}^{\infty} \gamma^t r(s_t, d_{S,t}, \bar{a}_t) \}$, where $(s_0, d_{S,0}) = (s, d_S)$, $\bar{a}_0 = \bar{a}$, $\bar{a}_t \sim \pi(s_t, d_{S,t})$. The optimal policy is denoted by $\pi^* \in \arg\max V^\pi(s, d_S)$.

Despite the intuitive analogy to finite-agent MDP, solving the mean-field MDP poses some unique challenges. In addition to having an unknown transition kernel and reward, the mean-field MDP takes a distribution as its state, which we do not have complete information of during training. In the following section, we propose our mean-field Neural Proximal Policy Optimization (MF-PPO) algorithm that, with a careful architecture design, can solve such mean-field MDP in a model-free fashion efficiently.

3 MEAN-FIELD PROXIMAL POLICY OPTIMIZATION

Our algorithm falls into the category of the actor-critic learning paradigm, consisting of alternating iterations of policy evaluation and improvement. The unique features of MF-PPO lie in the facts: (1) it exploits permutation invariance of the system, reducing search space of the actor/critic networks drastically and enables much more efficient learning; (2) it can handle a varying number of agents. For simplicity of exposition, we consider a fixed number of agents here.

Throughout the rest of the section, we assume that for any joint state $(s, d_S) \in \mathcal{S} \times \mathcal{P}(\mathcal{S})$, the agent has access to N i.i.d. samples $\{s_i\}_{i=1}^N$ from d_S . We denote concatenation of such samples as $\mathbf{s} \in \mathcal{S}^N$ and write $\mathbf{s} \stackrel{\text{i.i.d.}}{\sim} d_S$. MF-PPO maintains a pair of actor (denoted by F^A) and critic networks (denoted by F^Q), and uses the actor network to induce an energy-based policy $\pi(\bar{a} | s, d_S)$. Specifically, given state (s, d_S) , the actor network induces a distribution on set $\bar{\mathcal{A}}$ according to $\pi(\bar{a} | s, d_S) \propto \exp \{ \tau^{-1} F^A(s, d_S, \bar{a}) \}$, where τ denotes the temperature parameter. We use $\pi \propto \exp \{ F^A \}$ to denote the dependency of the policy on the energy function.

• **Permutation-invariant Actor and Critic.** We adopt a permutation invariant design of the actor and critic network. Specifically, given individual state $s \in \mathcal{S}$ and sampled states $\mathbf{s} \in \mathcal{S}^N$, the actor (resp. critic) network F^A (resp. F^Q) satisfies $F^A(s, \mathbf{s}, \bar{a}) = F^A(s, \kappa(\mathbf{s}), \bar{a})$ for any permutation mapping κ . With permutation invariance, the search space of the actor/critic network polynomially depends on the number of agents N .

Proposition 3.1. *The search space of a permutation invariance actor (critic) network is at the order of $(\sum_{k=1}^{\min\{|\mathcal{S}|, N\}} \binom{N-1}{k-1} \binom{|\mathcal{S}|}{k}) |\mathcal{S}| |\bar{\mathcal{A}}|$; Additionally, if $|\mathcal{S}| < N$, then the search space depends on N at the order of $N^{|\mathcal{S}|}$.*

Compared to architectures without permutation invariance, whose search space depends on N at the order of $(|\mathcal{S}| |\bar{\mathcal{A}}|)^N$, we can clearly see the search space of MF-PPO is exponentially smaller. Motivated by the characterization of the permutation invariant set function in Zaheer et al. (2017), the actor/critic network in MF-PPO takes the form of Deep Sets architecture, i.e., $F^A(s, \mathbf{s}, \bar{a}) = h(\sum_{s' \in \mathcal{S}} \phi(s, s', \bar{a}) / N)$. Both networks first aggregate local information by averaging over the output of a shared sub-network among agents, before feeding the aggregated information into a subsequent network $h(\cdot)$. See Figure 2 for detailed illustration. Effectively, by the average pooling layer and the preceding parameter sharing scheme, the network can keep its output unchanged when permuting the ordering of agents. Compared to a Graph Convolutional Neural Network (Kipf and Welling, 2017), which uses two sets of weights for the linear transformation layer, one for the agent itself and one for the aggregation state coming from the rest of the agents. The averaging operation is well suited for homogeneous agents and more parameter-efficient. It also naturally allows us to handle varying number of agents during training and evaluation, which is not readily achievable by GCN network.

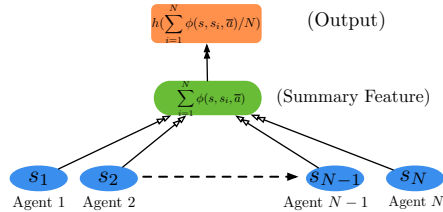


Figure 2: Illustration of a DeepSet parameterized critic network.

Naturally, the actor network when given the joint state-action pair (s, d_S, \bar{a}) is given by $F^A(s, d_S, \bar{a}) = h(\mathbb{E}_{s' \sim d_S} \phi(s, s', \bar{a}))$. We assume F^A is parameterized by a neural network with parameters $\alpha \in \mathbb{R}^D$, which is to be learned during training. We let the function class of all possible actor networks be denoted by \mathcal{F}^A . This same architecture design applies to the critic network, with learnable parameters denoted by $\theta \in \mathbb{R}^D$ and the function class denoted by \mathcal{F}^Q . MF-PPO then consists of successive iterations of policy evaluation and policy improvement described below.

• **Policy Evaluation.** At the k -th iteration of MF-PPO, we first update the critic network F^Q by minimizing the mean squared Bellman error while holding the actor network $F^{A,k}$ fixed. We denote the policy induced by the actor network as π_k , the stationary state distribution of policy π_k as ν_k , and the stationary state-action distribution as $\sigma_k(s, d_S, \bar{a}) := \nu_k(s, d_S) \pi_k(\bar{a} | s, d_S)$. Thus, we follow the update

$$\theta_k = \operatorname{argmin}_{\theta \in \mathbb{B}(\theta_0, R_\theta)} \mathbb{E}_{\sigma_k} \{ F_\theta^Q(s, d_S, \bar{a}) - [\mathcal{T}^{\pi_k} F_\theta^Q](s, d_S, \bar{a}) \}^2, \quad (3.1)$$

where $\mathbb{B}(\theta_0, R_\theta)$ denotes the Euclidean ball with radius R_θ centered at the initialized parameter θ_0 , and the Bellman evaluation operator \mathcal{T}^π is given by $\mathcal{T}^\pi F_\theta^Q(s, d_S, \bar{a}) = \mathbb{E} \{ (1 - \gamma)r(s, d_S, \bar{a}) + \gamma F_\theta^Q(s', d'_S, \bar{a}') \}$, where $(s', d'_S) \sim \mathbb{P}(\cdot | s, d_S, \bar{a})$, $\bar{a}' \sim \pi(\cdot | s', d'_S)$. We solve (3.1) by T-step temporal-difference (TD) update and output $\theta_k = \theta(T)$. At the t -th iteration of the TD update,

$$\begin{aligned} \theta(t + 1/2) &= \theta(t) - \eta \{ F_{\theta(t)}^Q(s, s, \bar{a}) - (1 - \gamma)r(s, d_S, \bar{a}) - \gamma F_{\theta(t)}^Q(s', s', \bar{a}') \} \nabla_\theta F_{\theta(t)}^Q(s, s, \bar{a}), \\ \theta(t + 1) &= \Pi_{\mathbb{B}(\theta_0, R_\theta)}(\theta(t + 1/2)), \end{aligned}$$

where we sample $(s, d_S, \bar{a}) \sim \sigma_k$, $(s', d'_S) \sim \mathbb{P}(\cdot | s, d_S, \bar{a})$, $\bar{a}' \sim \pi_k(\cdot | s', d'_S)$, $s \stackrel{\text{i.i.d.}}{\sim} d_S$, $s' \stackrel{\text{i.i.d.}}{\sim} d'_S$. We use $\Pi_X(\cdot)$ to denote the orthogonal projection onto set X , and η to denote the step size. Note that here for the simplicity of analyses, we sample the state-action pair (s, d_S, \bar{a}) independently from the stationary distribution. We remark that trajectory samples can also be used, which essentially requires bounding the bias of the gradient at each iteration due to the dependencies between trajectory samples, and we can readily apply the fast-mixing property of Markov chains to control such a bias (Bhandari et al., 2018). The details of policy evaluation are summarized in Algorithm 2, Appendix A.

• **Policy Improvement.** Following the policy evaluation step, MF-PPO updates its policy by updating the policy network F^A , which is the energy function associated with the policy. We update the policy network by

$$\pi_{k+1} = \operatorname{argmax}_{\pi \propto \exp\{F_{\theta_k}^A, k+1\}} \mathbb{E}_{\nu_k} \{ \langle F_{\theta_k}^Q(s, d_S, \cdot), \pi(\cdot | s, d_S) \rangle - v_k \text{KL}(\pi(\cdot | s, d_S) \| \pi_k(\cdot | s, d_S)) \}.$$

The update rule intuitively reads as increasing the probability for choosing action \bar{a} if it yields a higher value for critic network $F^Q(s, d_S, \bar{a})$, which can be viewed as a softened version of policy iteration (Bertsekas, 2011). Additionally, by controlling the proximity parameter v_k , we can control the softness of the update, with $v \rightarrow 0$ yielding the vanilla policy iteration. Moreover, without constraint of $F_{\theta_k}^A, k+1 \in \mathcal{F}^A$, such an update would have a nice closed form expression, and π_{k+1} itself is another energy-based policy.

Proposition 3.2. *Let $\pi_k \propto \exp(\tau_k^{-1} F_{\alpha_k}^A)$ denote the energy-based policy, then the update*

$$\bar{\pi}_{k+1} = \operatorname{argmax}_\pi \mathbb{E}_{\nu_k} \{ \langle F_{\theta_k}^Q(s, d_S, \cdot), \pi(\cdot | s, d_S) \rangle - v_k \text{KL}(\pi(\cdot | s, d_S) \| \pi_k(\cdot | s, d_S)) \}$$

yields $\bar{\pi}_{k+1} \propto \exp\{v_k^{-1} F_{\theta_k}^Q + \tau_k^{-1} F_{\alpha_k}^A\}$.

To take into account that the representable function of actor network resides in \mathcal{F}_A , we update the policy by projecting the energy function of $\bar{\pi}_{k+1}$ back to \mathcal{F}_A . Specifically, by denoting $\pi_{k+1} \propto \exp\{\tau_{k+1}^{-1} F_{\alpha_{k+1}}^A\}$, we recover the next actor network $F_{\alpha_{k+1}}^A$ (i.e., energy) by performing the following regression task

$$\alpha_{k+1}^* = \operatorname{argmin}_{\alpha \in \mathbb{B}(R_\alpha, \alpha_0)} \mathbb{E}_{\bar{\sigma}_k} \{ F_\alpha^A(s, d_S, \bar{a}) - \tau_{k+1} [v_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a})] \}^2, \quad (3.2)$$

where $\tilde{\sigma}_k = \nu_k \pi_0$. We approximately solve (3.2) via T-step stochastic gradient descent (SGD), and output $\alpha_{k+1} = \alpha(T) \approx \alpha_{k+1}^*$. At the t -th iteration of SGD,

$$\alpha(t+1/2) = \alpha(t) - \eta \nabla_{\alpha} F_{\alpha(t)}^A(s, \mathbf{s}, \bar{a}) \{ F_{\alpha(t)}^A(s, \mathbf{s}, a) - \tau_{k+1} (v_k^{-1} F_{\theta_k}^Q(s, \mathbf{s}, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, \mathbf{s}, \bar{a})) \},$$

$$\alpha(t+1) = \Pi_{\mathbb{B}(R_{\alpha}, \alpha_0)}(\alpha(t+1/2)),$$

where we sample $(s, d_S, \bar{a}) \sim \tilde{\sigma}_k$, and $\mathbf{s} \stackrel{\text{i.i.d.}}{\sim} d_S$, and η is the step size. The details are summarized in Algorithm 3 of Appendix A. Finally, we present the complete MF-PPO in Algorithm 1.

Algorithm 1 Mean-Field Neural Proximal Policy Optimization

Require: Mean-field MDP $(\mathcal{S} \times \mathcal{P}(\mathcal{S}), \bar{\mathcal{A}}, \mathbb{P}, r)$, discount factor γ ; number of outer iterations K , number of inner updates T ; policy update parameter v , step size η , projection radius R_{α}, R_{θ} .

Initialize: $\tau_0 \leftarrow 1$, $F^{A,0} \leftarrow 0$, $\pi_0 \propto \exp\{\tau_0^{-1} F^{A,0}\}$ (uniform policy).

for $k = 0, \dots, K-1$ **do**

Set temperature parameter $\tau_{k+1} \leftarrow v\sqrt{K}/(k+1)$, and proximity parameter $v_k \leftarrow v\sqrt{K}$

Solve (3.1) to update the critic network $F_{\theta_k}^Q$, using TD update (Algorithm 2)

Solve (3.2) to update the actor network for $F_{\alpha_{k+1}}^A$, using SGD update (Algorithm 3)

Update policy: $\pi_{k+1} \propto \exp\{\tau_{k+1}^{-1} F_{\alpha_{k+1}}^A\}$

end for

4 GLOBAL CONVERGENCE OF MF-PPO

We present the global convergence of MF-PPO algorithm for the two-layer permutation-invariant parameterization of the actor and critic networks. We remark that our analysis can be extended to multi-layer permutation-invariant networks, and we present the two-layer case here for simplicity of exposition. Specifically, the actor and critic networks take the form

$$F_{\alpha}^A(s, \mathbf{s}, \bar{a}) = \frac{1}{\sqrt{m_A N}} \sum_{j=1}^{m_A} \sum_{s' \in \mathbf{s}} u_j \sigma(\alpha_j^{\top}(s, s', \bar{a})), \quad F_{\theta}^Q(s, \mathbf{s}, \bar{a}) = \frac{1}{\sqrt{m_Q N}} \sum_{j=1}^{m_Q} \sum_{s' \in \mathbf{s}} v_j \sigma(\theta_j^{\top}(s, s', \bar{a})),$$

where m_A (resp. m_Q) is the width of the actor (resp. critic) network, and $\sigma(x) = \max\{x, 0\}$ denotes the ReLU activation. We randomly initialize u_j (resp. v_j) and first layer weights $\alpha_0 = [\alpha_{0,1}^{\top}, \dots, \alpha_{0,m_A}^{\top}]^{\top} \in \mathbb{R}^{d \cdot m_A}$ (resp. $\theta_0 \in \mathbb{R}^{d \cdot m_Q}$) by

$$u_j \sim \text{Unif}\{-1, +1\}, \alpha_{0,j} \sim \mathcal{N}(0, \mathbf{I}_d/d), \quad \forall j \in [m].$$

For ease of analysis, we take $m = m_A = m_Q$ and share the initialization of α_0 and θ_0 (resp. u_0 and v_0). Additionally, we keep u_j 's fixed during training, and α (resp. θ) within ball $\mathbb{B}(\alpha_0, R_A)$ (resp. $\mathbb{B}(\theta_0, R_Q)$) throughout training. We define the following function class which approximates the class of previously defined actor/critic network for large network width.

Definition 4.1. Given $R_{\beta} > 0$, define the function class over domain $\mathcal{S} \times \mathcal{S} \times \mathcal{A}$ by

$$\mathcal{F}_{\beta, m} = \{f_{\beta}(\cdot) \mid f_{\beta}(s, s', \bar{a}) = \frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \mathbb{1}\{\beta_{0,j}^{\top}(s, s', \bar{a}) > 0\} \beta_j^{\top}(s, s', \bar{a}), \|\beta - \beta_0\|_2 \leq R_{\beta}\},$$

where $v_j, \beta_{0,j}$ are random weights sampled according to

$$v_j \sim \text{Unif}\{-1, +1\}, \beta_{0,j} \sim \mathcal{N}(0, \mathbf{I}_d/d), \quad \forall j \in [m].$$

$\mathcal{F}_{\beta, m}$ also induces the function class over $\mathcal{S} \times \mathcal{P}(\mathcal{S}) \times \mathcal{A}$ given by

$$\mathcal{F}_{\beta, m}^{\mathcal{P}} = \{F(\cdot) \mid F(s, d_S, \bar{a}) = \mathbb{E}_{s' \sim d_S} f_{\beta}(s, s', \bar{a}), f_{\beta} \in \mathcal{F}_{\beta, m}\}.$$

It is well known that functions within $\mathcal{F}_{\beta, m}$ approximate functions within the reproducing kernel Hilbert space associated with kernel $K(x, y) = \mathbb{E}_{z \sim \mathcal{N}(0, \mathbf{I}_d/d)} \{\mathbb{1}(z^{\top} x > 0, z^{\top} y > 0)\}$ for a large network width m (Jacot et al., 2018; Chizat and Bach, 2018; Cai et al., 2019; Arora et al., 2019) and whose RKHS norm is bounded by R_{β} . For large R_{β} and m , $\mathcal{F}_{\beta, m}$ represents a rich class of functions. Additionally, functions within $\mathcal{F}_{\beta, m}^{\mathcal{P}}$ can be viewed as the mean-embedding of the joint state-action pair onto the RKHS space (Muandet et al., 2016; Song et al., 2009; Smola et al., 2007). Below, we make one important assumption, which assumes that $\mathcal{F}_{\beta, m}^{\mathcal{P}}$ is rich enough to represent the Q-function of all the policies within our policy class.

Assumption 1. For any policy π induced by $F_A \in \mathcal{F}^A$, we have $Q^\pi \in \mathcal{F}_{\theta, m_Q}^P$.

We remark that Assumption 1 can be relaxed into requiring that $\mathcal{F}_{\beta, m}^P$ has ϵ approximation error when parameterizing the set of Q-functions, with an additional ϵ term appearing in the convergence bound developed in Theorem 4.1 (Lan, 2021).

We define mild conditions stating boundedness of reward, and regularity of stationary distributions.

Assumption 2. Reward function $r(\cdot) \leq \bar{r}$ for some $\bar{r} > 0$. Additionally, there exists $c > 0$ such that $\mathbb{E} \{ \mathbb{1} (|z^\top(s, s', a)| \leq t) \} \leq c \cdot \frac{t}{\|z\|_2}$ for any $z \in \mathbb{R}^d$ and $t > 0$.

We measure the progress of MF-PPO in Algorithm 1 using the expected total reward

$$\mathcal{L}(\pi) = \mathbb{E}_{\nu^*} [V^\pi(s, d_S)] = \mathbb{E}_{\nu^*} \{ \langle Q^\pi(\cdot | s, d_S), \pi(\cdot | s, d_S) \rangle \}, \quad (4.1)$$

where ν^* is the stationary state distribution of the optimal policy π^* . We also denote σ^* as the stationary state-action distribution induced by π^* . Note that we have: $\mathcal{L}(\pi^*) = \mathbb{E}_{\nu^*} [V^{\pi^*}(s, d_S)] \geq \mathbb{E}_{\nu^*} [V^\pi(s, d_S)] = \mathcal{L}(\pi)$, for any policy π . Our main results are presented in the following theorem, showing that $\mathcal{L}(\pi_k)$ converges to $\mathcal{L}(\pi^*)$ at a sub-linear rate.

Theorem 4.1 (Global Convergence of MF-PPO). *Under Assumptions 1 and 2, the policies $\{\pi_k\}_{k=1}^K$ generated by Algorithm 1 satisfy*

$$\min_{0 \leq k \leq K} \{ \mathcal{L}(\pi^*) - \mathcal{L}(\pi_k) \} \leq \frac{v(\log |\bar{\mathcal{A}}| + \sum_{k=1}^{K-1} (\varepsilon_k + \varepsilon'_k))}{(1-\gamma)\sqrt{K}} + \frac{M}{(1-\gamma)v\sqrt{K}},$$

where $M = \mathbb{E}_{\nu^*} \left\{ \max_{\bar{a} \in \bar{\mathcal{A}}} (F_{\theta_0}^Q(s, s', \bar{a}))^2 \right\} + 2R_A^2$, ε_k and ε'_k are defined in Lemma 4.3. In particular, suppose at each iteration of MF-PPO, we observe $N = \Omega(K^3 R_A^4 + K R_Q^4)$ agents, and the actor/critic network satisfies $m_A = \Omega(K^6 R_A^{10} + K^4 R_A^{10} |\bar{\mathcal{A}}|^2)$, $m_Q = \Omega(K^2 R_Q^{10})$, and $T = \Omega(K^3 R_A^4 + K R_Q^4)$, then we have

$$\min_{0 \leq k \leq K} \{ \mathcal{L}(\pi^*) - \mathcal{L}(\pi_k) \} \leq \frac{v^2(\log |\bar{\mathcal{A}}| + \mathcal{O}(1)) + M}{(1-\gamma)v\sqrt{K}}.$$

Theorem 4.1 states that, given sufficiently many agents and a large enough actor/critic network, MF-PPO attains global optimality at a sublinear rate. Our result shows that when solving the mean-field MDP, having more agents serves as a blessing instead of a curse. In addition, as will be demonstrated in our proof sketch, there exists an inherent phase transition, where the final optimality gap is dominated by statistical error for a small number of agents (first phase); and by optimization error for a large number of agents (second phase).

The complete proof of Theorem 4.1 takes careful analysis on the error from policy evaluation (3.1) and the improvement step (3.2). The analysis on the outer iterations of MF-PPO can be overviewed as approximate mirror descent, which needs to take into account how the evaluation and improvement error interacts. Intuitively, the tuple $(\varepsilon_k, \varepsilon'_k)$ describes the effect of policy update when using approximate policy evaluation and policy improvement, and will be further clarified in the ensuing discussion. We present here the skeleton of our proof, and defer the technical detail to the appendix.

Proof Sketch. We first establish the convergence of the policy evaluation and improvement step.

Lemma 4.1 (Policy Evaluation). *Under the same assumptions in Theorem 4.1, let Q^{π_k} denote the Q-function of policy π_k , let $\epsilon_k = \mathbb{E}_{\text{init}, \sigma_k} \left(F_{\bar{\theta}(T)}^Q(\cdot) - Q^{\pi_k}(\cdot) \right)^2$ denote policy evaluation error, where $\bar{\theta}(T)$ is the output of Algorithm 2, we have $\epsilon_k = \mathcal{O} \left(\frac{R_Q^2}{T^{1/2}} + \frac{R_Q^{5/2}}{m_Q^{1/4}} + \frac{R_Q^2}{N^{1/2}} + \frac{R_Q^3}{m_Q^{1/2}} \right)$.*

Lemma 4.2 (Policy Improvement). *Under the same assumptions in Theorem 4.1, let $\epsilon'_{k+1} = \mathbb{E}_{\text{init}, \bar{\sigma}_k} \left\{ F_{\bar{\alpha}(T)}^A(\cdot) - \tau_{k+1} \left(v_k^{-1} F_{\theta_k}^Q(\cdot) + \tau_k^{-1} F_{\alpha_k}^A(\cdot) \right) \right\}^2$ denote policy optimization error, where $\bar{\alpha}(T)$ is the output of Algorithm 3, we have $\epsilon'_{k+1} = \mathcal{O} \left(\frac{R_A^2}{T^{1/2}} + \frac{R_A^{5/2}}{m_A^{1/4}} + \frac{R_A^2}{N^{1/2}} + \frac{R_A^3}{m_A^{1/2}} \right)$.*

Lemma 4.1 and 4.2 show that despite non-convexity, both policy evaluation and policy improvement steps converge approximately to the global optimal solution. In particular, for both policy evaluation steps and improvement steps, given networks with large width, for a small number of iterations T , the optimization error $\mathcal{O}(T^{-1/2})$ dominates the optimality gap; for a large number of iterations T , the statistical error $\mathcal{O}(N^{-1/2})$ dominates the optimality gap.

With Lemma 4.1 and 4.2, we illustrate the main argument for the proof of Theorem 4.1. Let us assume the ideal case when $\epsilon_k = \epsilon'_{k+1} = 0$. Note that for $\epsilon_k = 0$, we obtain the exact Q -function of policy π_k . For $\epsilon'_{k+1} = 0$, we obtain the ideal energy-based updated policy define in Proposition 3.2. That is,

$$\pi_{k+1} = \operatorname{argmax}_{\pi} \mathbb{E}_{\nu_k} \left\{ \langle Q^{\pi_k}(s, d_S, \cdot), \pi(\cdot|s, d_S) \rangle v_k \operatorname{KL}(\pi(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) \right\}. \quad (4.2)$$

Without function approximation, problem (4.2) can be solved by treating each joint state (s, d_S) independently, hence one can apply the well known three-point lemma in mirror descent (Chen and Teboulle, 1993) and obtain that, for all $(s, d_S) \in \mathcal{S} \times \mathcal{P}(\mathcal{S})$:

$$\begin{aligned} & \langle Q^{\pi_k}(s, d_S, \cdot), \pi^*(\cdot|s, d_S) - \pi_k(\cdot|s, d_S) \rangle \\ & \leq v_k \{ \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) - \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_{k+1}(\cdot|s, d_S)) - \operatorname{KL}(\pi_{k+1}(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) \} \\ & \quad + \langle Q^{\pi_k}(s, d_S, \cdot), \pi_{k+1}(\cdot|s, d_S) - \pi_k(\cdot|s, d_S) \rangle. \end{aligned}$$

From Lemma 6.1 in Kakade and Langford (2002), the expectation of the left hand side yields exactly $(1 - \gamma) \{ \mathcal{L}(\pi^*) - \mathcal{L}(\pi_k) \}$. Hence we have

$$\begin{aligned} & (1 - \gamma) \{ \mathcal{L}(\pi^*) - \mathcal{L}(\pi_k) \} \\ & \leq v_k \mathbb{E}_{\nu^*} \{ \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) - \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_{k+1}(\cdot|s, d_S)) \\ & \quad - \operatorname{KL}(\pi_{k+1}(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) \} + \mathbb{E}_{\nu^*} \langle Q^{\pi_k}(s, d_S, \cdot), \pi_{k+1}(\cdot|s, d_S) - \pi_k(\cdot|s, d_S) \rangle. \end{aligned}$$

Pinsker's Inequality $\operatorname{KL}(\pi_{k+1}(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) \geq \frac{1}{2} \|\pi_{k+1} - \pi_k\|_1^2$, combined with observation $\|Q^{\pi_k}(s, d_S, \cdot)\|_{\infty} \leq \bar{r}/(1 - \gamma)$, and basic inequality $-ax^2 + bx \leq b^2/(4a)$ for $a > 0$ gives us

$$\begin{aligned} & (1 - \gamma) \{ \mathcal{L}(\pi^*) - \mathcal{L}(\pi_k) \} \\ & \leq v_k \mathbb{E}_{\nu^*} \{ \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) - \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_{k+1}(\cdot|s, d_S)) \} + \frac{\bar{r}^2}{2v_k(1 - \gamma)^2}. \end{aligned}$$

By setting $v_k = \mathcal{O}(\sqrt{K})$, and telescoping the above inequality from $k = 0$ to $K - 1$, we obtain: $\min_{0 \leq k \leq K-1} \{ \mathcal{L}(\pi^*) - \mathcal{L}(\pi_k) \} = \mathcal{O}(1/\sqrt{K})$. Note that the key element in the global convergence of MF-PPO is the recursion defined in the previous inequality, which holds whenever we have an exact Q -function of the current policy and no function approximation is used when updating the next policy. Now MF-PPO conducts approximate policy evaluation $\epsilon_k > 0$, and after obtaining the approximate Q -function, conducts approximate policy improvement step $\epsilon'_{k+1} > 0$ with function approximation. In addition, the error of approximating the Q -function introduced in the evaluation step can be further compounded in the improvement step. Nevertheless, the previous inequality still holds approximately, with additional terms representing the policy evaluation/improvement errors.

Lemma 4.3 (Liu et al. (2019a)). *Let ϵ_k (evaluation error) and ϵ'_{k+1} (improvement error) be defined as in Lemma 4.1 and Lemma 4.2, respectively. We have:*

$$\begin{aligned} (1 - \gamma) (\mathcal{L}(\pi^*) - \mathcal{L}(\pi_k)) & \leq v_k \mathbb{E}_{\nu^*} \{ \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_k(\cdot|s, d_S)) - \operatorname{KL}(\pi^*(\cdot|s, d_S) \| \pi_{k+1}(\cdot|s, d_S)) \} \\ & \quad + v_k (\epsilon_k + \epsilon'_{k+1}) + v_k^{-1} M. \end{aligned} \quad (4.3)$$

where

$$\epsilon_k = \tau_{k+1}^{-1} \epsilon'_{k+1} \phi_{k+1}^* + v_k^{-1} \epsilon_k \psi_k^*, \quad \epsilon'_k = |\mathcal{A}| \tau_{k+1}^{-2} (\epsilon'_{k+1})^2, \quad M = \mathbb{E}_{\nu^*} \left\{ \max_{\bar{a} \in \bar{\mathcal{A}}} \left[F_{\theta_0}^Q(s, d_S, \bar{a}) \right]^2 \right\} + 2R_A^2.$$

In addition, ϕ_k^* and ψ_k^* are defined by:

$$\phi_k^* = \mathbb{E}_{\tilde{\sigma}_k} [|\mathrm{d}\pi^*/\mathrm{d}\pi_0 - \mathrm{d}\pi_k/\mathrm{d}\pi_0|^2]^{1/2}, \quad \psi_k^* = \mathbb{E}_{\sigma_k} [|\mathrm{d}\sigma^*/\mathrm{d}\sigma_k - \mathrm{d}(\nu^* \times \pi_k)/\mathrm{d}\sigma_k|^2]^{1/2}.$$

Finally, by telescoping inequality (4.3) from $k = 0$ to $K - 1$, we complete the proof of Theorem 4.1.

5 EXPERIMENTS

We perform experiments on the benchmark multi-agent particle environment (MPE) used in prior work (Lowe et al., 2017; Mordatch and Abbeel, 2018; Liu et al., 2019b). In the *cooperative navigation* task, N agents each with position $x_i \in \mathbb{R}^2$ must move to cover N fixed landmarks at positions $y_i \in \mathbb{R}^2$. They receive a team reward $R = -\sum_{i=1}^N \min_{j \in [N]} \|y_i - x_j\|_2$; In the *cooperative push* task, N agents with position $x_i \in \mathbb{R}^2$ work together to push a ball $x \in \mathbb{R}^2$ to a fixed landmark $y \in \mathbb{R}^2$. They receive a team reward $R = -\min_{j \in [N]} \|x_j - x\|_2 - \|x - y\|_2$. Both tasks involve homogeneous

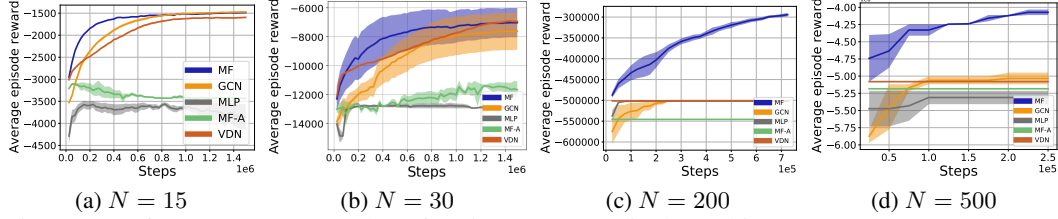


Figure 3: Performance versus number of environment steps in the multi-agent *cooperative navigation* task, over five independent runs per method. Points are taken every 1000 training episodes, with the first point taken after the first 1000, and is the average reward of 1000 evaluation episodes. MF significantly outperforms other critic representations for various number of agents.

agents, and all the agents share the same team reward. Note that MPE environment also models interaction between agents, including collision, and the collided agents receive negative rewards.

We instantiate our method **MF**, by parameterizing the centralized critic function using a DeepSet (Zaheer et al., 2017) network, with two hidden layers. We use a standard two-layer multi-layer perception (MLP) for the centralized actor network in all algorithms. The actor network outputs the mean and diagonal covariance of a Gaussian distribution over the joint action space. We refer interested readers to Appendix B for detailed configurations of hyperparameters.

We compare with two other critic representations: one that uses MLP for the centralized critic, labeled **MLP**, and another that uses a graph convolutional network for the critic (Liu et al., 2019b), labeled **GCN**. Note that the GCN representation is permutation invariant if one imposes a fully-connected graph for the agents in the MPE, but this invariance property does not hold for all graphs in general. We also compare with an extension of (Yang et al., 2018) to the case of continuous action spaces, labeled **MF-A**, in which each independent DDPG agent i has a decentralized critic $Q(s_i, a_i, \bar{a}_i)$ that takes in the mean of all other agents' actions $\bar{a}_i := \frac{1}{N-1} \sum_{j \neq i} a_j$. Finally, we include comparison with **VDN** (Sunehag et al., 2017), where the centralized critic network is the direct summation of local critic networks and thus being permutation invariant. Empirically, as we find that off-policy RL learns faster than on-policy RL in the MPE with higher agent number, regardless of the critic representation, we make all comparisons on top of MADDPG (Lowe et al., 2017). For a fair comparison of all critic representations, we ensure that all neural network architectures contain approximately the same number of trainable weights.

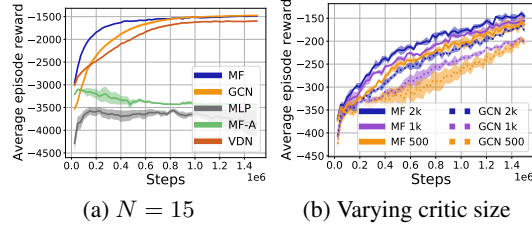


Figure 4: (a) Performance versus number of environment steps in the multi-agent *cooperative push* task. (b) MF outperforms GCN even with a fewer number of critic network parameters ($N = 3$).

For the cooperative navigation task, Figure 3 shows that the permutation invariant critic representation based on DeepSet enables MF to learn faster or reach a higher performance than all other representations and methods in the MPE with 15, 200, and 500 agents. For the cooperative push task, Figure 4a demonstrates a similar performance improvement provided by MF. In addition, we also demonstrate the superior parameter efficiency of MF compared to GCN. Figure 4b shows that MF consistently and significantly outperforms GCN as the number of parameters in their critic network varies over a range, with all other settings fixed. In particular, MF requires much fewer critic parameters to achieve higher performance than GCN.

Computational Improvements. Theorem 4.1 states that to obtain a small optimality gap in MF-PPO, one needs to compute the update on a large number of agents. We remark that with the dual embedding techniques developed in Dai et al. (2017), one can avoid computation on all the agents by sampling a small number of agents to compute the update. This technique could be readily incorporated into MF-PPO to improve its computational efficiency.

Conclusion. We propose a principled approach to exploit agent homogeneity and permutation invariance through the mean-field approximation in MARL. Our results are also the first to show the global convergence of MARL algorithms with neural networks as function approximators. This is in sharp contrast to current practices, which are mostly heuristic methods without convergence guarantees.

REFERENCES

- ARORA, S., DU, S. S., HU, W., LI, Z. and WANG, R. (2019). Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*.
- BERTSEKAS, D. P. (2011). Approximate policy iteration: A survey and some new methods. *Journal of Control Theory and Applications* **9** 310–335.
- BHANDARI, J., RUSSO, D. and SINGAL, R. (2018). A finite time analysis of temporal difference learning with linear function approximation. In *Conference On Learning Theory*. PMLR.
- BLOEM-REDDY, B. and TEH, Y. W. (2019). Probabilistic symmetry and invariant neural networks. *arXiv preprint arXiv:1901.06082*.
- CAI, Q., YANG, Z., LEE, J. D. and WANG, Z. (2019). Neural temporal-difference learning converges to global optima. In *Advances in Neural Information Processing Systems*.
- CAO, Y., YU, W., REN, W. and CHEN, G. (2013). An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Transactions on Industrial informatics* **9** 427–438.
- CARMONA, R., LAURIÈRE, M. and TAN, Z. (2019). Model-free mean-field reinforcement learning: mean-field mdp and mean-field q-learning. *arXiv preprint arXiv:1910.12802*.
- CHEN, G. and TEBOULLE, M. (1993). Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization* **3** 538–543.
- CHIZAT, L. and BACH, F. (2018). A note on lazy training in supervised differentiable programming. *arXiv preprint arXiv:1812.07956* **8**.
- DAI, B., SHAW, A., LI, L., XIAO, L., HE, N., LIU, Z., CHEN, J. and SONG, L. (2017). Sbed: Convergent reinforcement learning with nonlinear function approximation. *arXiv preprint arXiv:1712.10285*.
- GU, H., GUO, X., WEI, X. and XU, R. (2019). Dynamic programming principles for learning mfcs. *arXiv preprint arXiv:1911.07314*.
- GU, H., GUO, X., WEI, X. and XU, R. (2020). Mean-field controls with q-learning for cooperative marl: Convergence and complexity analysis. *arXiv preprint arXiv:2002.04131*.
- GU, H., GUO, X., WEI, X. and XU, R. (2021). Mean-field multi-agent reinforcement learning: A decentralized network approach. *arXiv preprint arXiv:2108.02731*.
- JACOT, A., GABRIEL, F. and HONGLER, C. (2018). Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*.
- KAKADE, S. and LANGFORD, J. (2002). Approximately optimal approximate reinforcement learning. In *ICML*, vol. 2.
- KALYANAKRISHNAN, S., LIU, Y. and STONE, P. (2006). Half field offense in robocup soccer: A multiagent reinforcement learning case study. In *Robot Soccer World Cup*. Springer.
- KIPF, T. N. and WELING, M. (2017). Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*.
- KOBER, J., BAGNELL, J. A. and PETERS, J. (2013). Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research* **32** 1238–1274.
- LAN, G. (2021). Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *arXiv preprint arXiv:2102.00135*.
- LITTMAN, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *Machine Learning Proceedings 1994*. Elsevier, 157–163.

- LIU, B., CAI, Q., YANG, Z. and WANG, Z. (2019a). Neural proximal/trust region policy optimization attains globally optimal policy. *arXiv preprint arXiv:1906.10306*.
- LIU, I.-J., YEH, R. A. and SCHWING, A. G. (2019b). Pic: Permutation invariant critic for multi-agent deep reinforcement learning. *arXiv preprint arXiv:1911.00025*.
- LOWE, R., WU, Y., TAMAR, A., HARB, J., ABBEEL, O. P. and MORDATCH, I. (2017). Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*.
- MATARIĆ, M. J. (1997). Reinforcement learning in the multi-robot domain. In *Robot colonies*. Springer, 73–83.
- MENDA, K., CHEN, Y.-C., GRANA, J., BONO, J. W., TRACEY, B. D., KOCHENDERFER, M. J. and WOLPERT, D. (2018). Deep reinforcement learning for event-driven multi-agent decision processes. *IEEE Transactions on Intelligent Transportation Systems* **20** 1259–1268.
- MORDATCH, I. and ABBEEL, P. (2018). Emergence of grounded compositional language in multi-agent populations. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- MUANDET, K., FUKUMIZU, K., SRIPERUMBUDUR, B. and SCHÖLKOPF, B. (2016). Kernel mean embedding of distributions: A review and beyond. *arXiv preprint arXiv:1605.09522*.
- PANAIT, L. and LUKE, S. (2005). Cooperative multi-agent learning: The state of the art. *Autonomous agents and multi-agent systems* **11** 387–434.
- SCHULMAN, J., WOLSKI, F., DHARIWAL, P., RADFORD, A. and KLIMOV, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- SHALEV-SHWARTZ, S., SHAMMAH, S. and SHASHUA, A. (2016). Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*.
- SMOLA, A., GRETTON, A., SONG, L. and SCHÖLKOPF, B. (2007). A hilbert space embedding for distributions. In *International Conference on Algorithmic Learning Theory*. Springer.
- SONG, L., HUANG, J., SMOLA, A. and FUKUMIZU, K. (2009). Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual International Conference on Machine Learning*.
- SUNEHAG, P., LEVER, G., GRUSLYS, A., CZARNECKI, W. M., ZAMBALDI, V., JADERBERG, M., LANCTOT, M., SONNERAT, N., LEIBO, J. Z., TUYSLS, K. ET AL. (2017). Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*.
- SUTTON, R. S. and BARTO, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- VINYALS, O., BABUSCHKIN, I., CZARNECKI, W. M., MATHIEU, M., DUDZIK, A., CHUNG, J., CHOI, D. H., POWELL, R., EWALDS, T., GEORGIEV, P. ET AL. (2019). Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature* **575** 350–354.
- WANG, L., YANG, Z. and WANG, Z. (2020). Breaking the curse of many agents: Provable mean embedding q-iteration for mean-field reinforcement learning. In *International Conference on Machine Learning*. PMLR.
- YANG, J., YE, X., TRIVEDI, R., XU, H. and ZHA, H. (2017). Learning deep mean field games for modeling large population behavior. *arXiv preprint arXiv:1711.03156*.
- YANG, Y., LUO, R., LI, M., ZHOU, M., ZHANG, W. and WANG, J. (2018). Mean field multi-agent reinforcement learning. *arXiv preprint arXiv:1802.05438*.
- ZAHEER, M., KOTTUR, S., RAVANBAKSH, S., POZOS, B., SALAKHUTDINOV, R. R. and SMOLA, A. J. (2017). Deep sets. In *Advances in neural information processing systems*.

ZHANG, K., YANG, Z. and BAŞAR, T. (2019). Multi-agent reinforcement learning: A selective overview of theories and algorithms. *arXiv preprint arXiv:1911.10635* .

ZHENG, S., TROTT, A., SRINIVASA, S., NAIK, N., GRUESBECK, M., PARKES, D. C. and SOCHER, R. (2020). The ai economist: Improving equality and productivity with ai-driven tax policies. *arXiv preprint arXiv:2004.13332* .

Supplementary Material for A Permutation Invariant Approach to Deep Mean-Field Multi-Agent Reinforcement Learning

A ALGORITHMS FOR POLICY EVALUATION/IMPROVEMENT

Algorithm 2 Policy Evaluation via TD

Require: Mean-field MDP $(\bar{\mathcal{S}}, \bar{\mathcal{A}}, \mathbb{P}, r, \gamma)$, sample $\{(s_t, \mathbf{s}_t, \bar{a}_t, s'_t, \mathbf{s}'_t, \bar{a}'_t)\}$, number of iterations T
Initialize: $w_j(0) \sim \text{Unif}\{-1, +1\}$, $\theta_j(0) \sim \mathcal{N}(0, \mathbf{I}_d/d)$, $\forall j \in [m_Q]$
Set step size $\eta \leftarrow T^{-1/2}$
for $t = 0, \dots, T-1$ **do**
 Set $(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') \leftarrow (s_t, \mathbf{s}_t, \bar{a}_t, s'_t, \mathbf{s}'_t, \bar{a}'_t)$
 $\theta(t+1/2) = \theta(t) - \eta \left\{ F_{\theta(t)}^Q(s, \mathbf{s}, \bar{a}) - (1 - \gamma)r(s, \mathbf{d}_{\mathcal{S}}, \bar{a}) - \gamma F_{\theta(t)}^Q(s, \mathbf{s}', \bar{a}') \right\} \nabla_{\theta} F_{\theta(t)}^Q(s, \mathbf{s}, \bar{a})$
 $\theta(t+1) = \Pi_{\mathbb{B}(R_{\theta}, \theta_0)}(\theta(t+1/2))$
end for
Take ergodic average: $\bar{\theta}_T \leftarrow \frac{1}{T} \sum_{t=0}^{T-1} \theta(t)$
Output: $F_{\bar{\theta}_T}^Q$

Algorithm 3 Policy Improvement via SGD

Require: Mean-field MDP $(\bar{\mathcal{S}}, \bar{\mathcal{A}}, \mathbb{P}, r, \gamma)$, sample $\{(s_t, \mathbf{s}_t, \bar{a}_t)\}$, number of iterations T
Initialize: $v_j(0) \sim \text{Unif}\{-1, +1\}$, $\alpha_j(0) \sim \mathcal{N}(0, \mathbf{I}_d/d)$, $\forall j \in [m_A]$
Set step size $\eta \leftarrow T^{-1/2}$
for $t = 0, \dots, T-1$ **do**
 Set $(s, \mathbf{s}, \bar{a}) \leftarrow (s_t, \mathbf{s}_t, \bar{a}_t)$
 $\alpha(t+1/2) = \alpha(t) - \eta \left\{ F_{\alpha(t)}^A(s, \mathbf{s}, \bar{a}) - \tau_{k+1}^{-1} F_{\theta_k}^Q(s, \mathbf{s}, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, \mathbf{s}, \bar{a}) \right\} \nabla_{\alpha} F_{\alpha(t)}^A(s, \mathbf{s}, \bar{a})$
 $\alpha(t+1) = \Pi_{\mathbb{B}(R_{\alpha}, \alpha_0)}(\alpha(t+1/2))$
end for
Take ergodic average: $\bar{\alpha}_T \leftarrow \frac{1}{T} \sum_{t=0}^{T-1} \alpha(t)$
Output: $F_{\bar{\alpha}_T}^A$

B EXPERIMENTAL DETAILS

◇ Additional Experiment

We further demonstrate the superiority of our proposed MF method on the problem with huge number of agents. As can be seen from Figure 5, the MF method significantly outperforms other alternatives for cooperative navigation task with 1000 agents. The performance gap is even more significant than the experiments presented in Section 5.

◇ **Environment.** We used the open-source code for the multi-agent particle environments provided by Liu et al. (2019b), which itself is based on the original code by Lowe et al. (2017), without any modification. Please refer to Liu et al. (2019b, Appendix B) for complete details.

◇ **Computing infrastructure and runtime.** Experiments were run on Intel Xeon 6154 CPUs, using one core for each independent policy training session. Average training time was approximately 4 hours for $N = 15$ and $N = 30$ with $1.5e6$ steps, and 12 hours for $N = 200$ with $7.25e5$ steps.

B.1 IMPLEMENTATION

◇ **GCN**¹. For experiments based on MADDPG, we re-ran the open-source code provided by Liu et al. (2019b) without modification. For experiments based on PPO, we used the GCN network in Kipf and Welling (2017) as the critic in PPO. Performance of GCN reported in this work are the results of our experimental runs.

◇ **MF**. In the PPO-based implementation, the joint policy is parameterized by a multi-layer perceptron (MLP) with two hidden layers of size 128 and ReLU activation, which takes in the concatenation of all agents’ observations and outputs the mean and diagonal covariance of a joint Gaussian policy. In the MADDPG-based implementation, each decentralized actor (i.e., policy) network is an MLP with two hidden layers of size 187 and ReLU activation, which takes in each agent’s individual observation and outputs the agent’s real-valued action vector. The centralized critic has the form $Q(s) = f_3(f_2(\sum_{i=1}^N f_1(s_i)))$, where f_1 and f_2 are hidden layers of size $h = 190/205/187/187/187$ (for $N = 3/6/15/30/200$ agents, respectively, to have approximately equal number of trainable critic parameters as the GCN critic used in Liu et al. (2019b)) with ReLU activation, and f_3 is a linear layer with output size 1.

◇ **MLP**. The centralized critic is an MLP with two hidden layers of size 138,187,75,42,7 for the case of $N = 3, 6, 15, 30, 200$ agents, respectively, such that total number of trainable parameters is approximately equal to that in GCN (Liu et al., 2019b). The joint policy has the same representation as the one in **MF** for both PPO-based and MADDPG-based implementations.

◇ **Hyperparameters**. We used the Adam optimizer for all algorithms. For experiments based on PPO, the common hyperparameters across all algorithms are: discount factor $\gamma = 0.99$, KL divergence target $d_{\text{targ}} = 0.01$, KL divergence initial coefficient $\beta = 3$, and GAE $\lambda = 0.95$ (see Schulman et al. (2017)), policy entropy loss coefficient 0.01, max gradient norm 0.5, and actor learning rate 10^{-4} . Each PPO training step uses 5 epochs with minibatch size 2. We did a sweep of critic learning rate for each critic architecture, choosing 10^{-3} for MF, 10^{-4} for GCN and MLP, and 10^{-2} for MF-A.

For experiments based on MADDPG, the common hyperparameters are: actor learning rate 0.01, critic learning rate 0.01, batch size 1024, discount factor $\gamma = 0.95$, replay buffer size 10^5 , 8 actor and critic updates per 100 environment steps, and target network update rate $\tau = 0.01$.

C PROOFS IN SECTION 2

Proof of Proposition 2.1. We first prove the first part of claim. We begin by showing that the optimal Q-function Q^* is permutation invariant, i.e., $Q^*(s, \mathbf{a}) = Q^*(\kappa(s), \kappa(\mathbf{a}))$ for all permutation mapping κ . Note that Q^* is the unique solution to Bellman optimality condition over the space of Q-function

$$\begin{aligned} Q(s, \mathbf{a}) &= \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, \mathbf{a})} \left\{ r(s, \mathbf{a}) + \max_{\mathbf{a}'} Q(s', \mathbf{a}') \right\} \\ &= r(s, \mathbf{a}) + \sum_{s'} \mathbb{P}(s' | s, \mathbf{a}) \max_{\mathbf{a}'} Q(s', \mathbf{a}'). \end{aligned} \quad (\text{C.1})$$

¹We label this architecture as “GCN” to distinguish it from alternative ways to instantiate a permutation invariant critic.

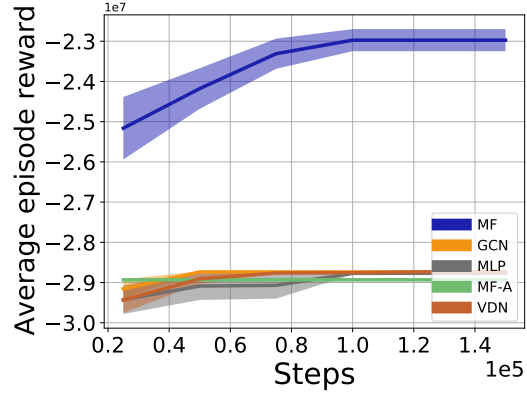


Figure 5: Performance versus number of environment steps in the multi-agent *cooperative navigation* task, with $N = 1000$ number of agents.

Or equivalently, from the permutation invariance of r and \mathbb{P} defined in (2.1)

$$\begin{aligned} Q(\kappa(\mathbf{s}), \kappa(\mathbf{a})) &= r(\kappa(\mathbf{s}), \kappa(\mathbf{a})) + \sum_{\mathbf{s}'} \mathbb{P}(\kappa(\mathbf{s}') | \kappa(\mathbf{s}), \kappa(\mathbf{a})) \max_{\mathbf{a}'} Q(\kappa(\mathbf{s}'), \kappa(\mathbf{a}')) \\ &= r(\mathbf{s}, \mathbf{a}) + \sum_{\mathbf{s}'} \mathbb{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) \max_{\mathbf{a}'} Q(\kappa(\mathbf{s}'), \kappa(\mathbf{a}')). \end{aligned}$$

That is, $Q(\kappa(\mathbf{s}), \kappa(\mathbf{a}))$ is also a solution to the Bellman optimality condition (C.1). From the uniqueness of the solution, we have $Q^*(\mathbf{s}, \mathbf{a}) = Q^*(\kappa(\mathbf{s}), \kappa(\mathbf{a}))$. Hence the optimal policy $\nu^*(\mathbf{s}) = \operatorname{argmax}_{\mathbf{a}} Q^*(\mathbf{s}, \mathbf{a})$ is permutation invariant.

For second part of the proposition, we show the detailed proof for permutation invariance of value function V . For any permutation invariance policy ν , let $R^\nu(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \nu} R(\mathbf{s}, \mathbf{a})$, and $\mathbb{P}^\nu(\mathbf{s}' | \mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \nu} \mathbb{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a})$ denote the key elements of the induced Markov reward process. One can clearly see that from permutation invariance defined in (2.1), we have that both R^ν and $\mathbb{P}^\nu(\mathbf{s}' | \mathbf{s})$ is permutation invariant. We define the k -step value function V_k^ν as the expected reward from time 0 to time k , then:

$$\begin{aligned} V_1^\nu(\mathbf{s}) &= R^\nu(\mathbf{s}), \\ V_k^\nu(\mathbf{s}) &= R^\nu(\mathbf{s}) + \gamma \mathbb{E}_{\mathbf{s}' \sim \mathbb{P}^\nu(\cdot | \mathbf{s})} V_{k-1}^\nu(\mathbf{s}'), \quad \forall k > 1. \end{aligned}$$

We can see from above recursion that V_k^ν is permutation invariant for all $k \geq 1$. From $V^\nu = \lim_{k \rightarrow \infty} V_k^\nu$, we can conclude that V^ν is also permutation invariant.

For permutation invariance of Q-function, recall that $Q^\nu(\mathbf{s}, \mathbf{a})$ is the unique solution of the following Bellman evaluation equation

$$\begin{aligned} Q(\mathbf{s}, \mathbf{a}) &= \mathbb{E}_{\mathbf{s}' \sim \mathbb{P}(\cdot | \mathbf{s}, \mathbf{a})} \{r(\mathbf{s}, \mathbf{a}) + \mathbb{E}_{\mathbf{a}' \sim \nu(\cdot | \mathbf{s}')} Q(\mathbf{s}', \mathbf{a}')\} \\ &= r(\mathbf{s}, \mathbf{a}) + \sum_{\mathbf{s}', \mathbf{a}'} \mathbb{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) Q(\mathbf{s}', \mathbf{a}') \nu(\mathbf{s}' | \mathbf{a}'). \end{aligned} \quad (\text{C.2})$$

From the permutation invariance of r , \mathbb{P} and ν , we have Q^ν is also a solution to (C.2), as

$$Q(\kappa(\mathbf{s}), \kappa(\mathbf{a})) = r(\mathbf{s}, \mathbf{a}) + \sum_{\mathbf{s}', \mathbf{a}'} \mathbb{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) Q(\kappa(\mathbf{s}'), \kappa(\mathbf{a}')) \nu(\mathbf{s}' | \mathbf{a}').$$

Since Q^ν is the unique solution to (C.2), we have $Q^\nu(\mathbf{s}, \mathbf{a}) = Q^\nu(\kappa(\mathbf{s}), \kappa(\mathbf{a}))$. □

Proof of Proposition 2.2. The proof for showing $V^\nu(s, \mathbf{s}_r) = V^\nu(s, \kappa(\mathbf{s}_r))$ follows the exact same ingredients as in the proof for Proposition 2.1. Note that from the permutation invariance, we can show $\sum_{t \geq 0} \gamma^t r(s^t, \mathbf{s}_r^t, \bar{a}^t) |_{\mathbf{s}_r^0 = \mathbf{s}} \stackrel{d}{=} \sum_{t \geq 0} \gamma^t r(s^t, \mathbf{s}_r^t, \bar{a}^t) |_{\mathbf{s}_r^0 = \kappa(\mathbf{s})}$ for any permutation mapping κ of order N . That is, the conditional distribution of future reward conditioned upon current joint state \mathbf{s}_r is invariant w.r.t permutation of \mathbf{s}_r . Hence from Theorem 11 of Bloem-Reddy and Teh (2019), we know that there exists a function g such that $\sum_{t \geq 0} \gamma^t r(s^t, \mathbf{s}_r^t) |_{\mathbf{s}_r^0 = \mathbf{s}} \stackrel{d}{=} g(s, \eta, \hat{p}_{\mathbf{s}_r})$, where $\eta \sim \text{Unif}([0, 1])$. Hence we conclude with $V^\nu(s, \mathbf{s}_r) = E_\eta g(s, \eta, \hat{p}_{\mathbf{s}_r})$. □

D PROOFS IN SECTION 3

Proof of Proposition 3.1. For any permutation invariant actor/critic function $F(s, \mathbf{s}, \bar{a}) = F(s, \kappa(\mathbf{s}), \bar{a})$, we consider the size of tabular representation for such functions and denote it as the size for the search space. We say \mathbf{s} and \mathbf{s}' are *permutation equivalent* if $\mathbf{s} = \kappa(\mathbf{s}')$ for some permutation mapping κ . We can observe that (i) the permutation equivalent is a valid equivalent relation defined over the space of all possible \mathbf{s} ; (ii) \mathbf{s} and \mathbf{s}' are permutation equivalent if and only if for every value v in \mathcal{S} , v occurs the same number of times in \mathbf{s} and \mathbf{s}' . One can verify that for \mathbf{s} containing N elements, with each element s' taking values in \mathcal{S} , then the number of equivalence classes over the space of all possible \mathbf{s} , induced by *permutation equivalent* relation is $\sum_{k=1}^{\min\{|\mathcal{S}|, N\}} \binom{N-1}{k-1} \binom{|\mathcal{S}|}{k}$. The claim of Proposition 3.1 follows immediately. □

Proof of Proposition 3.2. Note that the maximization problem can be solved separately for each $s, d_{\mathcal{S}}$. Hence we can equivalently solve

$$\bar{\pi}(\cdot | s, d_{\mathcal{S}}) \max_{\pi(\cdot | s, d_{\mathcal{S}}) \in \mathbb{R}^{\mathcal{A}}} \left[\left\langle F_{\theta_k}^Q(s, d_{\mathcal{S}}, \cdot), \pi(\cdot | s, d_{\mathcal{S}}) \right\rangle - v_k \text{KL}(\pi_\alpha(\cdot | s, d_{\mathcal{S}}) \| \pi(\cdot | s, d_{\mathcal{S}})) \right],$$

which is an optimization problem over finite-dimensional parameter $\pi(\cdot|s, d_S) \in \mathbb{R}^{\bar{\mathcal{A}}}$. Setting the derivative w.r.t. $\pi(\bar{a}|s, d_S)$ for each $\bar{a} \in \bar{\mathcal{A}}$ equal to zero, we have

$$F_{\theta_k}^Q(s, d_S, \bar{a}) - v_k(\log\left(\frac{\pi}{\pi_k(\bar{a}|s, d_S)}\right) - 1) = 0,$$

or equivalently, $\bar{\pi}(a|s, d_S) \propto \pi_k(a|s, d_S) \exp\{v_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a})\}$. Plug in the definition of $\pi_k \propto \exp(\tau_k^{-1} F_{\alpha_k}^A)$ completes the proof. \square

E PROOFS IN SECTION 4

We provide a unified analysis of the policy evaluation and policy improvement step. The policy evaluation and improvement steps can be described as minimizing the following mean-squared loss:

$$\beta(t+1) = \underset{\beta \in \mathbb{B}(R_\beta, \beta_0)}{\operatorname{argmin}} \mathbb{E}_{(s, d_S, a) \sim \rho} [(F_\beta(s, d_S, \bar{a}) - \zeta_{F_\beta}(s, d_S, \bar{a}))^2]. \quad (\text{E.1})$$

where the operator ζ maps function F_β to

$$\zeta_{F_\beta}(s, d_S, \bar{a}) = \mathbb{E}_{s', d'_S \sim \mathbb{P}(\cdot|s, d_S, \bar{a}), \bar{a}' \sim \pi(\cdot|s', d'_S)} [\tau \xi(s, d_S, \bar{a}) + \mu F_\beta(s', d'_S, \bar{a}')]. \quad (\text{E.2})$$

Formulation (E.1) includes policy improvement and policy evaluation as special cases:

Policy Improvement. This corresponds to $\rho = \tilde{\sigma}_k$, $\xi(s, d_S, \bar{a}) = \tau_{k+1}(v_k^{-1} F_{\beta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a}))$, $\tau = 1, \mu = 0$,

Policy Evaluation. This corresponds to $\rho = \sigma_k$, $\tau = 1 - \gamma$, $\mu = \gamma$, $\xi(s, d_S, \bar{a}) = r(s, d_S, \bar{a})$. Note that $\zeta_{F_\theta}(s, d_S, \bar{a})$ equals to the Bellman operator: $\zeta_{F_\theta}(s, d_S, \bar{a}) = \mathcal{T}^\pi F_\theta(s, d_S, \bar{a})$.

To solve problem (E.1), we consider the following generic update rule:

$$\beta(t+1/2) = (F_{\beta(t)}(s, d_S, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu F_{\beta(t)}(s', d'_S, \bar{a}')) \nabla_\beta F_{\beta(t)}(s, d_S, \bar{a}), \quad (\text{E.3})$$

$$\beta(t+1) = \Pi_{\mathbb{B}^0(R_\beta)}(\beta(t+1/2)). \quad (\text{E.4})$$

where $(s, d_S, \bar{a}) \sim \rho$, $(s', d'_S) \sim \mathbb{P}(\cdot|s, d_S, \bar{a})$, $\bar{a}' \sim \pi(\cdot|s', d'_S)$.

Instead of knowing s, d_S exactly, we only have access to s, d_S via N independent samples represented as the states of N agents. Hence, we perform the following update:

$$\beta(t+1/2) = (F_{\beta(t)}(s, \mathbf{s}, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu F_{\beta(t)}(s', \mathbf{s}', \bar{a}')) \nabla_\beta F_{\beta(t)}(s, \mathbf{s}, \bar{a}), \quad (\text{E.5})$$

$$\beta(t+1) = \Pi_{\mathbb{B}^0(R_\beta)}(\beta(t+1/2)). \quad (\text{E.6})$$

where $(s, d_S, \bar{a}) \sim \rho$, $(s', d'_S) \sim \mathbb{P}(\cdot|s, d_S, \bar{a})$, $\bar{a}' \sim \pi(\cdot|s', d'_S)$; and $\mathbf{s} \stackrel{\text{i.i.d.}}{\sim} d_S$, $\mathbf{s}' \stackrel{\text{i.i.d.}}{\sim} d'_S$.

Policy Improvement Updates. For $\rho = \tilde{\sigma}_k$, $\tau = 1, \mu = 0, R_\beta = R_\alpha$ and $\xi(s, d_S, \bar{a}) = \tau_{k+1}(\beta_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a}))$, we recover the policy improvement update in Algorithm 2.

Policy Evaluation Updates. For $\rho = \sigma_k$, $\tau = (1 - \gamma), \mu = \gamma, \xi(s, d_S, \bar{a}) = r(s, d_S, \bar{a}), R_\beta = R_\theta$, we recover the policy evaluation update in Algorithm 3.

We define a few more notations before we proceed

$$(\text{residual}): \delta_{\beta(t)}(s, d_S, \bar{a}, s', d'_S, \bar{a}') = F_{\beta(t)}(s, d_S, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu F_{\beta(t)}(s', d'_S, \bar{a}'),$$

$$(\text{stochastic semi-gradient}): g_{\beta(t)}(s, d_S, \bar{a}, s', d'_S, \bar{a}') = \delta_{\beta(t)}(s, d_S, \bar{a}, s', d'_S, \bar{a}') \nabla_\beta F_{\beta(t)}(s, d_S, \bar{a}),$$

$$(\text{population semi-gradient}): \bar{g}_{\beta(t)} = \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} g(s, d_S, \bar{a}, s', d'_S, \bar{a}'),$$

where $(s, d_S, \bar{a}) \sim \rho$, $(s', d'_S) \sim \mathbb{P}(\cdot|s, d_S, \bar{a})$, $\bar{a}' \sim \pi_k(\cdot|s', d'_S)$.

We also define the associated finite sample approximation:

$$\hat{\delta}_{\beta(t)}(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') = F_{\beta(t)}(s, \mathbf{s}, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu F_{\beta(t)}(s, \mathbf{s}', \bar{a}'),$$

$$\hat{g}_{\beta(t)}(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') = \hat{\delta}_{\beta(t)}(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') \nabla_\beta F_{\beta(t)}(s, \mathbf{s}, \bar{a}),$$

where $\mathbf{s} \stackrel{\text{i.i.d.}}{\sim} d_S$, $\mathbf{s}' \stackrel{\text{i.i.d.}}{\sim} d'_S$. Note that given \hat{g} is not an unbiased estimator of g , as $\mathbb{E}_{\mathbf{s}, \mathbf{s}'} \hat{g}(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') \neq g(s, d_S, \bar{a}, s', d'_S, \bar{a}')$.

E.1 LINEARIZATION AT INITIALIZATION

For a two layer ReLU network defined as $f_{\beta,u}(s, x, a) = \frac{1}{\sqrt{m}} \sum_{j=1}^m u_j \sigma(\beta_j^\top(s, x, a))$, where $\sigma(x) = \max\{0, x\}$, we define its linearization at the initialization

$$f_{\beta}^0(s, x, a) = \frac{1}{\sqrt{m}} \sum_{j=1}^m u_j \mathbb{1}\{\beta_j(0)^\top(s, x, a) > 0\} \beta_j^\top(s, x, a).$$

Note that $f_{\beta}^0(s, x, \bar{a}) = \langle \nabla_{\beta} f_{\beta(0)}(s, x, a), \beta \rangle$, which equivalently means f_{β}^0 is the local linearization of f_{β} at its initialization $\beta(0)$ (note we do not train the second layer u).

We define

$$F_{\beta}^0(s, \bar{a}) = \frac{1}{N} \sum_{i=1}^N f_{\beta}^0(s, s_i, \bar{a}), \quad F_{\beta}^0(s, d_S, \bar{a}) = \mathbb{E}_{x \sim d_S} f_{\beta}^0(s, x, \bar{a}),$$

and similarly define $\delta^0, g^0, \bar{g}^0, \hat{\delta}^0$, and \hat{g}^0 , with everything related to F_{β}, f_{β} replaced by F_{β}^0, f_{β}^0 .

Our main objectives are to show that the objective in (E.1) can be replaced by replacing everything related to F with its local linearization F^0 , and the linearized stochastic semi-gradient \hat{g}^0 remains close to the population semi-gradient \bar{g} when the number of hidden units m and number of agents are large enough. By doing so, we can conclude that learning with linearized stochastic semi-gradient \hat{g} is sufficient to solve (E.1) to high accuracy, when the number of hidden units m and number of agents are large enough. For the ease of exposition, we also make the following assumptions,

Assumption 3. We assume $\|(s, x, \bar{a})\|_2 \leq 1$ for all $(s, x, \bar{a}) \in \mathcal{S} \times \mathcal{S} \times \mathcal{A}$.

Assumption 4. The function $\xi(s, d_S, \bar{a})$ satisfies: for each $(s, d_S) \in \bar{\mathcal{S}}$ and any action a , we have $(\xi(s, d_S, \bar{a}))^2 \leq \mathbb{E}_{x \sim d_S} \tau_1 (f_{\beta(0)}(s, x, \bar{a}))^2 + \tau_2 R_{\beta}^2 + \tau_3$.

Remark: We will verify that with $(\tau_1, \tau_2, \tau_3) = (4, 4, 0)$, Assumption 4 holds for the policy evaluation step. With $(\tau_1, \tau_2, \tau_3) = (0, 0, \bar{r}^2)$, Assumption 4 holds for the policy improvement step.

We first bound the difference between $F_{\beta}^0(s, d_S, \bar{a})$ and $F_{\beta}(s, d_S, \bar{a})$.

Lemma E.1. With Assumption 2 and 3, we have

$$\mathbb{E}_{\text{init}, \rho} |F_{\beta}^0(s, d_S, \bar{a}) - F_{\beta}(s, d_S, \bar{a})|^2 \leq \mathcal{O}(R_{\beta}^3 m^{-1/2}).$$

Proof. Be definition of $F_{\beta}^0(s, d_S, \bar{a})$ and $F_{\beta}(s, d_S, \bar{a})$ and Jensen's inequality, we have

$$\mathbb{E}_{\text{init}, s, d_S, \bar{a}} |F_{\beta}^0(s, d_S, \bar{a}) - F_{\beta}(s, d_S, \bar{a})| \leq \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} |f_{\beta}^0(s, x, \bar{a}) - f_{\beta}(s, x, \bar{a})|.$$

We can further bound the right hand side with

$$\begin{aligned} & \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} |f_{\beta}^0(s, x, \bar{a}) - f_{\beta}(s, x, \bar{a})| \\ &= \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \frac{1}{\sqrt{m}} \sum_{j=1}^m |\mathbb{1}\{\beta_j(0)^\top(s, x, \bar{a}) > 0\} - \mathbb{1}\{\beta_j^\top(s, x, \bar{a}) > 0\}| \beta_j^\top(s, x, \bar{a})| \\ &\leq \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \frac{1}{\sqrt{m}} \sum_{j=1}^m |\mathbb{1}\{\beta_j(0)^\top(s, x, \bar{a}) > 0\} - \mathbb{1}\{\beta_j^\top(s, x, \bar{a}) > 0\}| \|\beta_j - \beta_j(0)\|_2, \end{aligned}$$

where the last inequality comes from that fact that $\mathbb{1}\{\beta_j(0)^\top(s, x, \bar{a}) > 0\} \neq \mathbb{1}\{\beta_j^\top(s, x, \bar{a}) > 0\}$ implies $|\beta_j^\top(s, x, \bar{a})| \leq |(\beta_j - \beta_j(0))^\top(s, x, \bar{a})| \leq \|\beta_j - \beta_j(0)\|_2 \|(s, x, \bar{a})\|_2 \leq \|\beta_j - \beta_j(0)\|_2$.

Applying Cauchy Schwartz inequality, we have

$$\begin{aligned}
& \mathbb{E}_{\text{init}, s, d_S, a} |F_\beta^0(s, d_S, \bar{a}) - F_\beta(s, d_S, \bar{a})|^2 \\
& \leq \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \frac{1}{m} \left(\sum_{j=1}^m |\mathbb{1}\{\beta_j(0)^\top(s, x, \bar{a}) > 0\} - \mathbb{1}\{\beta_j^\top(s, x, \bar{a}) > 0\}| \right) \left(\sum_{j=1}^m \|\beta_j - \beta_j(0)\|_2^2 \right) \\
& \leq \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \frac{R_\beta^2}{m} \sum_{j=1}^m \mathbb{1}\{|\beta_j(0)^\top(s, x, \bar{a})| < \|\beta_j(0) - \beta_j\|_2\} \\
& \leq \mathbb{E}_{\text{init}} \frac{cR_\beta^2}{m} \sum_{j=1}^m \frac{\|\beta_j(0) - \beta_j\|_2}{\|\beta_j(0)\|_2} \\
& \leq \mathbb{E}_{\text{init}} \frac{cR_\beta^2}{m} \left(\sum_{j=1}^m \|\beta_j(0) - \beta_j\|_2^2 \right)^{1/2} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} = \mathcal{O} \left(\frac{cR_\beta^3}{m^{1/2}} \right),
\end{aligned}$$

where in the third inequality we use Assumption 2. \square

Next, we bound the difference between \bar{g} and \bar{g}^0 .

Lemma E.2. *With Assumption 2, 3 and 4, we have*

$$\mathbb{E}_{\text{init}} \|\bar{g}_\beta - \bar{g}_\beta^0\|_2^2 \leq \mathcal{O} \left(\frac{R_\beta^3}{m^{1/2}} \right).$$

Proof. With decomposition

$$\begin{aligned}
g_\beta - g_\beta^0 &= (F(s, d_S, \bar{a}) - F(s, d_S, \bar{a}) - \mu(F(s', d'_S, \bar{a}') - F(s', d'_S, \bar{a}')) \nabla F(s, d_S, \bar{a}) \\
&\quad + (F(s, d_S, \bar{a}) - \tau\xi(s, d_S, \bar{a}) - \mu F(s', d'_S, \bar{a}')) (\nabla_\beta F(s, d_S, \bar{a}) - \nabla_\beta F(s, d_S, \bar{a})),
\end{aligned}$$

we apply basic inequality $\|a + b\|_2^2 \leq 2(\|a\|_2^2 + \|b\|_2^2)$ and obtain

$$\begin{aligned}
\|\bar{g}_\beta - \bar{g}_\beta^0\|_2^2 &\leq 2 \left(\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} |F(s, d_S, \bar{a}) - F(s, d_S, \bar{a}) - \mu(F(s', d'_S, \bar{a}') - F(s', d'_S, \bar{a}'))| \|\nabla F(s, d_S, \bar{a})\|_2 \right)^2 + \\
&\quad 2 \left(\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} |F(s, d_S, \bar{a}) - \tau\xi(s, d_S, \bar{a}) - \mu F(s', d'_S, \bar{a}')| \|\nabla_\beta F(s, d_S, \bar{a}) - \nabla_\beta F(s, d_S, \bar{a})\|_2 \right)^2.
\end{aligned} \tag{E.7}$$

We have $\nabla_\beta F(s, d_S, \bar{a}) = \frac{1}{\sqrt{m}} (\mathbb{1}\{\beta_1(0)^\top(s, x, \bar{a}) > 0\}(s, x, \bar{a}), \dots, \mathbb{1}\{\beta_m(0)^\top(s, x, \bar{a}) > 0\}(s, x, \bar{a}))$, and with the assumption that $\|(s, x, \bar{a})\|_2 \leq 1$, we have $\|\nabla_\beta F(s, d_S, \bar{a})\|_2 \leq 1$. In addition, we have

$$\mathbb{E}_{\text{init}, s, d_S, \bar{a}, s', d'_S, \bar{a}'} |F(s, d_S, \bar{a}) - F(s, d_S, \bar{a}) - \mu(F(s', d'_S, \bar{a}') - F(s', d'_S, \bar{a}'))|^2 = \mathcal{O}(R_\beta^3 m^{-1/2}),$$

as implied by Lemma E.1. Hence the first term in (E.7) is of order $\mathcal{O}(R_\beta^3 m^{-1/2})$.

To bound the second term in (E.7), by Cauchy-Schwartz inequality we have

$$\begin{aligned}
& \left(\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} |F(s, d_S, \bar{a}) - \tau\xi(s, d_S, \bar{a}) - \mu F(s', d'_S, \bar{a}')| \|\nabla_\beta F(s, d_S, \bar{a}) - \nabla_\beta F(s, d_S, \bar{a})\|_2 \right)^2 \\
& \leq \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} |F(s, d_S, \bar{a}) - \tau\xi(s, d_S, \bar{a}) - \mu F(s', d'_S, \bar{a}')|^2 \mathbb{E}_{s, d_S, a} \|\nabla_\beta F(s, d_S, \bar{a}) - \nabla_\beta F(s, d_S, \bar{a})\|_2^2.
\end{aligned}$$

From Jensen's inequality we have

$$\begin{aligned}
|F(s, d_S, \bar{a}) - \tau\xi(s, d_S, \bar{a}) - \mu F(s', d'_S, \bar{a}')|^2 &\leq \mathbb{E} |f_\beta^0(s, x, \bar{a}) - \tau\xi(s, d_S, \bar{a}) - \mu f_\beta^0(s', x', \bar{a}')|^2 \\
&\leq 3\mathbb{E} [(f_\beta^0(s, x, \bar{a}))^2 + (\tau\xi(s, d_S, \bar{a}))^2 + (\mu f_\beta^0(s', x', \bar{a}'))^2].
\end{aligned}$$

Note that $\nabla f_\beta^0 = \nabla f_{\beta(0)}^0$ for all β , we have

$$f_\beta^0(s, x, \bar{a}) \leq f_{\beta(0)}^0(s, x, \bar{a}) + \|\nabla_\beta f_\beta^0\|_2 \|\beta - \beta(0)\|_2 \leq f_{\beta(0)}^0(s, x, \bar{a}) + R_\beta.$$

We have $3 \left(f_{\beta}^0(s, x, \bar{a})^2 + (\mu f_{\beta}^0(s', x', \bar{a}'))^2 \right) \leq 6 \left(f_{\beta(0)}(s, x, \bar{a})^2 + (\mu f_{\beta(0)}(s', x', \bar{a}'))^2 \right) + 12R_{\beta}^2$.

Note that from Assumption 4: $\xi(s, d_S, \bar{a})^2 \leq \tau_1 \mathbb{E}_{s,a} \left[(f_{\beta(0)}(s, x, \bar{a}))^2 + \tau_2 R_{\beta}^2 + \tau_3 \right]$. Hence

$$\begin{aligned} & \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} |F(s, d_S, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu F(s', d'_S, \bar{a}')|^2 \\ & \leq \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} \mathbb{E}_{x, x'} [6(f_{\beta(0)}(s, x, \bar{a})^2 + 6\mu^2 f_{\beta(0)}(s', x', \bar{a}')^2 + 12R_{\beta}^2 + 3\tau^2 \xi(s, d_S, \bar{a})^2)] \\ & = \mathbb{E}_{s, d_S, x, \bar{a}} [6f_{\beta(0)}(s, x, \bar{a})^2 + 6\mu^2 f_{\beta(0)}(s, x, \bar{a})^2 + 12R_{\beta}^2 + 3\tau^2 \xi(s, d_S, \bar{a})^2] \\ & \leq \mathbb{E}_{s, d_S, x, \bar{a}} [6f_{\beta(0)}(s, x, \bar{a})^2 + 6\mu^2 f_{\beta(0)}(s, x, \bar{a})^2 + 12R_{\beta}^2 + 3\tau^2 \tau_1 f_{\beta(0)}(s, x, \bar{a})^2 + 3\tau^2 \tau_2 R_{\beta}^2 + 3\tau^2 \tau_3] . \end{aligned} \quad (\text{E.8})$$

On the other hand, we have

$$\begin{aligned} & \mathbb{E}_{s, d_S, a} \|\nabla_{\beta} F(s, d_S, \bar{a}) - \nabla_{\beta} F(s, d_S, \bar{a})\|_2^2 \\ & \leq \frac{1}{m} \mathbb{E}_{s, d_S, a, x} \| (\mathbb{1}\{\beta_1(0)^{\top}(s, x, \bar{a}) > 0\} - \mathbb{1}\{\beta_1^{\top}(s, x, \bar{a}) > 0\}, \dots, \mathbb{1}\{\beta_m(0)^{\top}(s, x, \bar{a}) > 0\} - \mathbb{1}\{\beta_m^{\top}(s, x, \bar{a}) > 0\}) \|_2^2 \\ & \quad \cdot \|(s, x, \bar{a})\|_2^2 \\ & \leq \mathbb{E}_{s, d_S, \bar{a}, x} \frac{1}{m} \sum_{j=1}^m (\mathbb{1}\{\beta_j(0)^{\top}(s, x, \bar{a}) > 0\} - \mathbb{1}\{\beta_j^{\top}(s, x, \bar{a}) > 0\})^2 \\ & \leq \mathbb{E}_{s, d_S, \bar{a}, x} \frac{1}{m} \sum_{j=0}^m \mathbb{1}\{\beta_j(0)^{\top}(s, x, \bar{a}) \leq \|\beta_j(0) - \beta_j\|_2\} \\ & \leq \frac{c}{m} \sum_{j=1}^m \frac{\|\beta_j(0) - \beta_j\|_2}{\|\beta_j(0)\|_2} \\ & \leq \frac{c}{m} \left(\left\| \sum_{j=1}^m \beta_j(0) - \beta_j \right\|_2^2 \right)^{1/2} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} \\ & \leq \frac{cR_{\beta}}{m} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} , \end{aligned} \quad (\text{E.10})$$

where in the fourth inequality we use Assumption 1, and in the final inequality we use $\left(\left\| \sum_{j=1}^m \beta_j(0) - \beta_j \right\|_2^2 \right)^{1/2} \leq R_{\beta}$.

Combining (E.8) and (E.10), to bound the second term in (E.7), it remains to bound the following:

$$\begin{aligned} & \mathbb{E}_{s, d_S, \bar{a}, x} \left\{ f_{\beta(0)}(s, x, \bar{a})^2 \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} \left(\frac{cR_{\beta}}{m} \right) \right\} \\ & = \frac{cR_{\beta}}{m^2} \mathbb{E}_{s, d_S, \bar{a}, x} \left(\sum_{j=1}^m \sigma^2(\beta_j(0)^{\top}(s, x, \bar{a})) + \sum_{k \neq l}^m u_k u_l \sigma(\beta_k(0)^{\top}(s, x, \bar{a})) \sigma(\beta_l(0)^{\top}(s, x, \bar{a})) \right) \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} \\ & \leq \frac{cR_{\beta}}{m^2} \mathbb{E}_{s, d_S, \bar{a}, x} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^2 + \sum_{k \neq l}^m u_k u_l \sigma(\beta_k(0)^{\top}(s, x, \bar{a})) \sigma(\beta_l(0)^{\top}(s, x, \bar{a})) \right) \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} , \end{aligned}$$

Taking expectation with respect to initialization and noticing that $\mathbb{E}_{\text{init}} \{u_k u_l\} = 0$ for $k \neq l$, we have

$$\begin{aligned}
& \frac{cR_\beta}{m^2} \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^2 + \sum_{k \neq l}^m u_k u_l \sigma(\beta_k(0)^\top(s, x, \bar{a})) \sigma(\beta_l(0)^\top(s, x, \bar{a})) \right) \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} \\
&= \frac{cR_\beta}{m^2} \mathbb{E}_{\text{init}} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^2 \right) \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right)^{1/2} \\
&\leq \frac{cR_\beta}{m^2} \mathbb{E}_{\text{init}}^{1/2} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^2 \right)^2 \mathbb{E}_{\text{init}}^{1/2} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^{-2} \right) \\
&= \mathcal{O} \left(\frac{cR_\beta}{m^{1/2}} \right).
\end{aligned}$$

Then the the second term in (E.7) is at the order of $\mathcal{O} \left(\max \left\{ \frac{R_\beta^3}{m^{1/2}}, \frac{R_\beta}{m^{1/2}} \right\} \right) = \mathcal{O} \left(\frac{R_\beta^3}{m^{1/2}} \right)$.

□

We then bound the the variance of stochastic semi-gradient $g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}')$.

Lemma E.3. *With Assumption 4, there exists $\varsigma^2 = \mathcal{O}(R_\beta^2)$, such that:*

$$\mathbb{E}_{\text{init}, s, d_S, \bar{a}, s', d'_S, \bar{a}'} \|g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \bar{g}_\beta\|_2^2 \leq \varsigma^2.$$

Proof. We have

$$\begin{aligned}
\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} \|g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \bar{g}_\beta\|_2^2 &\leq \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} \|g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}')\|_2^2 \\
&\leq \mathbb{E}_{s, s, d_S, s', d'_S, \bar{a}'} \|\delta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \nabla_\beta F_{\beta(t)}(s, d_S, \bar{a})\|_2^2.
\end{aligned}$$

With $\nabla_\beta F_{\beta(t)}(s, d_S, \bar{a}) = \mathbb{E}_x \nabla_\beta f_{\beta(t)}(s, x, \bar{a}) = \mathbb{E}_x \frac{1}{\sqrt{m}} (\mathbb{1}\{\beta_j^\top(s, x, \bar{a}) > 0\}, \dots, \mathbb{1}\{\beta_j^\top(s, x, \bar{a}) > 0\}) (s, x, \bar{a})$, we have $\|\nabla_\beta F_{\beta(t)}(s, d_S, \bar{a})\|_2 \leq 1$. Then

$$\begin{aligned}
& \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} \|g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \bar{g}_\beta\|_2^2 \\
&\leq \mathbb{E}_{s, s, d_S, s', d'_S, \bar{a}'} [\delta^2(s, d_S, \bar{a}, s', d'_S, \bar{a}')] \\
&= \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} (F_\beta(s, d_S, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu F_\beta(s', d'_S, \bar{a}'))^2 \\
&\leq \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}', x, x'} (f_\beta(s, x, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu f_\beta(s', x', \bar{a}'))^2 \\
&\leq 3 \mathbb{E}_{s, d_S, \bar{a}, x} \left[(1 + \mu^2) (f_\beta(s, x, \bar{a}))^2 + (\tau \xi(s, d_S, \bar{a}))^2 \right] \\
&\leq 3 \mathbb{E}_{s, d_S, \bar{a}, x} \left[(1 + \mu^2) (f_{\beta(0)}(s, x, \bar{a}))^2 + \tau_1 (\tau f_{\beta(0)}(s, x, \bar{a}))^2 \right] + \mathcal{O}(R_\beta^2) + \mathcal{O}(\tau_3).
\end{aligned}$$

where in the last inequality we use $|f_{\beta(0)}(s, x, \bar{a}) - f_\beta(s, x, \bar{a})| \leq \|\beta - \beta(0)\|_2 \leq R_\beta$ and Assumption

4. Finally, note that

$$\begin{aligned}
\mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} [f_{\beta(0)}(s, x, \bar{a})^2] &= \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \frac{1}{m} \left(\sum_{j=1}^m \sigma^2(\beta_j(0)^\top(s, x, \bar{a})) + \sum_{k \neq l}^m u_k u_l \sigma(\beta_k(0)^\top(s, x, \bar{a})) \sigma(\beta_l(0)^\top(s, x, \bar{a})) \right) \\
&= \mathbb{E}_{\text{init}, s, d_S, \bar{a}, x} \frac{1}{m} \left(\sum_{j=1}^m \sigma^2(\beta_j(0)^\top(s, x, \bar{a})) \right) \\
&\leq \mathbb{E}_{\text{init}} \frac{1}{m} \left(\sum_{j=1}^m \|\beta_j(0)\|_2^2 \right) = 1.
\end{aligned}$$

The claim follows immediately. □

We then bound the variance of $\hat{g}(s, \mathbf{s}, a, s', \mathbf{s}', \bar{a}')$ for fixed d_S, d'_S .

Lemma E.4. *For any $s, d_S, \bar{a}, s', d'_S, \bar{a}'$, let $\mathbf{s} \stackrel{\text{i.i.d.}}{\sim} s, d_S$ and $\mathbf{s}' \stackrel{\text{i.i.d.}}{\sim} s', d'_S$, and $|\mathbf{s}| = |\mathbf{s}'| = N$, we have*

$$\mathbb{E}_{\text{init}, \mathbf{s}, \mathbf{s}'} \|\hat{g}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}')\|_2^2 \leq \mathcal{O}\left(\frac{R_\beta^2}{N}\right).$$

Proof. We have

$$\begin{aligned} & \hat{g}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \\ &= \hat{\delta}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') \nabla_\beta \sum_{s \in \mathbf{s}} f_\beta(s, x, \bar{a})/N - \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \nabla_\beta \mathbb{E}_{x \sim d_S} f_\beta(s, x, \bar{a}) \\ &= \left(\hat{\delta}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right) \nabla_\beta \sum_{s \in \mathbf{s}} f_\beta(s, x, \bar{a})/N \\ & \quad + \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \left(\sum_{s \in \mathbf{s}} \nabla_\beta f_\beta(s, x, \bar{a})/N - \nabla_\beta \mathbb{E}_{x \sim d_S} f_\beta(s, x, \bar{a}) \right). \end{aligned} \quad (\text{E.11})$$

For the first term, note that $\|\sum_{s \in \mathbf{s}} \nabla_\beta f_\beta(s, x, \bar{a})/N\|_2 \leq 1$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{s}, \mathbf{s}'} \left\| \left(\hat{\delta}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right) \sum_{s \in \mathbf{s}} \nabla_\beta f_\beta(s, x, \bar{a})/N \right\|_2^2 \\ & \leq \mathbb{E}_{\mathbf{s}, \mathbf{s}'} \left(\hat{\delta}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right)^2 \\ & \leq \mathbb{E}_{x, x'} (f_\beta(s, x, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu f_\beta(s', x', \bar{a}'))^2 / N \\ & \leq 3 \mathbb{E}_{x, x'} \left[(f_{\beta(0)}(s, x, \bar{a}))^2 + \mu^2 (f_{\beta(0)}(s', x', \bar{a}'))^2 + \tau_1 (\tau f_{\beta(0)}(s, x, \bar{a}))^2 \right] / N + \mathcal{O}(R_\beta^2/N) + \mathcal{O}(\tau_3). \end{aligned}$$

where the last inequality we use the similar argument in Lemma E.3. Now as we have shown in Lemma E.3, we have $\mathbb{E}_{\text{init}, s'} [f_{\beta(0)}(s, x, \bar{a})^2] \leq 1$. Hence the first term in (E.11) is at the order of $\mathcal{O}(R_\beta^2/N)$.

To bound the second term in (E.11), from Jensen's inequality we have

$$\begin{aligned} & \mathbb{E}_{\text{init}} \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}')^2 \\ & \leq \mathbb{E}_{\text{init}, x, x'} (f_\beta(s, x, \bar{a}) - \tau \xi(s, d_S, \bar{a}) - \mu f_\beta(s', x', \bar{a}'))^2 \\ & \leq 3 \mathbb{E}_{\text{init}, x, x'} \left[(f_{\beta(0)}(s, x, \bar{a}))^2 + \mu^2 (f_{\beta(0)}(s', x', \bar{a}'))^2 + \tau_1 (\tau f_{\beta(0)}(s, x, \bar{a}))^2 \right] + \mathcal{O}(R_\beta^2) + \mathcal{O}(\tau_3) \\ & = \mathcal{O}(R_\beta^2). \end{aligned}$$

On the other hand, we have

$$\mathbb{E}_{\mathbf{s}} \left\| \sum_{s \in \mathbf{s}} \nabla_\beta f_\beta(s, x, \bar{a})/N - \nabla_\beta \mathbb{E}_{x \sim d_S} f_\beta(s, x, \bar{a}) \right\|_2^2 \leq \mathbb{E}_{x \sim d_S} \|\nabla_\beta f_\beta(s, x, \bar{a})\|_2^2 / N \leq 1/N.$$

Therefore, taking square on both sides of (E.11) and using Cauchy Schwartz inequality, we have

$$\begin{aligned} & \mathbb{E}_{\text{init}, \mathbf{s}, \mathbf{s}'} \|\hat{g}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}')\|_2^2 \\ & \leq 2 \mathbb{E}_{\text{init}, \mathbf{s}, \mathbf{s}'} \left(\hat{\delta}_\beta(s, \mathbf{s}, \bar{a}, s', \mathbf{s}', \bar{a}') - \delta_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right)^2 \left\| \sum_{s \in \mathbf{s}} \nabla_\beta f_\beta(s, x, \bar{a})/N \right\|_2^2 \\ & \quad + 2 \mathbb{E}_{\text{init}, \mathbf{s}, \mathbf{s}'} \{ \delta_\beta^2(s, d_S, \bar{a}, s', d'_S, \bar{a}') \} \mathbb{E}_{\text{init}, x} \left\| \sum_{s \in \mathbf{s}} \nabla_\beta f_\beta(s, x, \bar{a})/N - \nabla_\beta \mathbb{E}_{x \sim d_S} f_\beta(s, x, \bar{a}) \right\|_2^2 \\ & = \mathcal{O}(R_\beta^2/N). \end{aligned}$$

□

With Lemma E.3 and Lemma E.4, we can now bound the different between \hat{g}_β and \bar{g}_β .

Lemma E.5. *With Assumption 4, there exists $\varsigma^2 = \mathcal{O}(R_\beta^2)$, such that*

$$\mathbb{E}_{\text{init}, s, d_S, \bar{a}, s', d'_S, \bar{a}', s, s'} \|\hat{g}_\beta - \bar{g}_\beta\|_2^2 \leq \varsigma^2 + \mathcal{O}(R_\beta^2/N).$$

Proof. We have

$$\begin{aligned} \mathbb{E} \|\hat{g}_\beta - \bar{g}_\beta\|_2^2 &\leq 2\mathbb{E} \|\hat{g}_\beta(s, s, a, s', s', a') - g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}')\|_2^2 \\ &\quad + 2\mathbb{E} \|g_\beta(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \bar{g}_\beta\|_2^2 \end{aligned}$$

Now apply Lemma E.3 and Lemma E.4, the claim follow immediately. \square

We are now ready to show the convergence of using updates (E.3) and (E.4) to solve problem (E.1). We denote $F_{\beta_*}^0$ to be the function within the class $\mathcal{F}_{R,m}$, which satisfies the stationary condition:

$$\Pi_{\mathbb{B}(R_\beta, \beta_0)}(\beta_* - \bar{g}_{\beta_*}^0) = \beta_*, \quad (\text{E.12})$$

where $\Pi_X(\cdot)$ denotes the Euclidean projection onto set X . Condition (E.12) is equivalent to:

$$\langle \bar{g}_{\beta_*}^0, \beta - \beta_* \rangle \geq 0.$$

Note that $\bar{g}_{\beta_*}^0 = \mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} \left\{ \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') \nabla F_{\beta_*}^0(s, d_S, \bar{a}) \right\}$, and $\langle \nabla F_{\beta_*}^0(s, d_S, \bar{a}), \beta - \beta_* \rangle = F_\beta^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a})$. Hence we have

$$\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}'} \left\{ \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') (F_{\beta, u}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a})) \right\} \geq 0. \quad (\text{E.13})$$

With the definition of operator ζ_F in (E.2), we know that

$$\mathbb{E}_{s', d'_S, \bar{a}'} \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') = F_{\beta_*}^0(s, d_S, \bar{a}) - \zeta_{F_{\beta_*}^0}(s, d_S, \bar{a}). \quad (\text{E.14})$$

From (E.13) and (E.14), we have

$$\langle F_{\beta_*}^0 - \zeta_{F_{\beta_*}^0}, F_\beta^0 - F_{\beta_*}^0 \rangle_\rho \geq 0,$$

which is equivalent to that $F_{\beta_*}^0 = \Pi_{F_{R_\beta, m}}(\zeta_{F_{\beta_*}^0})$. That is, $F_{\beta_*}^0$ is the projection (with metric defined w.r.t $\langle \cdot, \cdot \rangle_\rho$), after applying operator ζ to itself.

Theorem E.1. *Let $\{\beta(t)\}_{t=0}^{T-1}$ be generated by updates (E.5) and (E.6). Define $\bar{\beta}_T = \frac{1}{T} \sum_{t=0}^{T-1} \beta(t)$, we have*

$$\mathbb{E}_{\text{init}, s, d_S, \bar{a}} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\bar{\beta}_T}(s, d_S, \bar{a}) \right]^2 \leq \mathcal{O} \left(\frac{R_\beta^2}{T^{1/2}} + \frac{R_\beta^{5/2}}{m^{1/4}} + \frac{R_\beta^2}{N^{1/2}} + \frac{R_\beta^3}{m^{1/2}} \right).$$

Proof. Conditioned on the t -th iteration, we have

$$\begin{aligned} &\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}', s, s'} [\|\beta(t+1) - \beta_*\|_2^2 | \beta(t)] \\ &= \|\Pi_{\mathbb{B}_{R_\beta}^0}(\beta - \eta \hat{g}_{\beta(t)}) - \Pi_{\mathbb{B}_{R_\beta}^0}(\beta_* - \eta \bar{g}_{\beta_*}^0)\|_2^2 \\ &\leq \mathbb{E} \|\beta(t) - \eta \hat{g}_{\beta(t)} - \beta_* - \eta \bar{g}_{\beta_*}^0\|_2^2 \\ &= \mathbb{E} (\|\beta(t) - \beta_*\|_2^2 - \eta \langle \beta(t) - \beta_*, \hat{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \rangle + \eta^2 \|\hat{g}_{\beta(t)} - \bar{g}_{\beta_*}^0\|_2^2), \end{aligned} \quad (\text{E.15})$$

where the inequality comes from the non-expansive property of projection.

We first consider the second term in (E.15), we have

$$\begin{aligned} &\mathbb{E}_{s, d_S, \bar{a}, s', d'_S, \bar{a}', s, s'} \langle \beta(t) - \beta_*, \hat{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \rangle \\ &= \mathbb{E} \left[\langle \beta(t) - \beta_*, \hat{g}_{\beta(t)} - g_{\beta(t)} \rangle + \langle \beta(t) - \beta_*, g_{\beta(t)} - \bar{g}_{\beta(t)} \rangle + \langle \beta(t) - \beta_*, \bar{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \rangle \right] \\ &= \mathbb{E} \langle \beta(t) - \beta_*, \hat{g}_{\beta(t)} - g_{\beta(t)} \rangle + \langle \beta(t) - \beta_*, \bar{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \rangle \\ &\geq \mathbb{E} \langle \beta(t) - \beta_*, \bar{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \rangle - R_\beta (\mathbb{E} \|\hat{g}_{\beta(t)} - g_{\beta(t)}\|_2^2)^{1/2} \\ &\geq \mathbb{E} \langle \beta(t) - \beta_*, \bar{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \rangle + \langle \beta(t) - \beta_*, \bar{g}_{\beta(t)}^0 - \bar{g}_{\beta_*}^0 \rangle - R_\beta (\mathbb{E} \|\hat{g}_{\beta(t)} - g_{\beta(t)}\|_2^2)^{1/2}, \end{aligned} \quad (\text{E.16})$$

where in the first equality we use definition $\mathbb{E}g_{\beta(t)} = \bar{g}_{\beta(t)}$. We further have

$$\left\langle \beta(t) - \beta_*, \bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0 \right\rangle \geq -R_\beta \|\bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0\|_2. \quad (\text{E.17})$$

Using again $\bar{g}_\beta^0 = \mathbb{E} \left\{ \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') \nabla F_\beta^0(s, d_S, \bar{a}) \right\}$, $\nabla F_\beta^0(s, d_S, \bar{a}) = \nabla F_{\beta_*}^0(s, d_S, \bar{a})$ for all β , and $\left\langle \nabla F_\beta^0(s, d_S, \bar{a}), \beta - \beta_* \right\rangle = F_\beta^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a})$, we have

$$\begin{aligned} & \mathbb{E}_{\text{init}, s, d_S, \bar{a}, s', d'_S, \bar{a}'} \left\langle \beta(t) - \beta_*, \bar{g}_{\beta(t)}^0 - \bar{g}_{\beta_*}^0 \right\rangle \\ &= \mathbb{E} \left(\delta_{\beta(t)}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right) \left(F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) \right) \\ &= \mathbb{E} \left(F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) - \mu \left[F_{\beta(t)}^0(s', d'_S, \bar{a}') - F_{\beta_*}^0(s', d'_S, \bar{a}') \right] \right) \left(F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) \right) \\ &\geq (1 - \mu) \mathbb{E}_{\text{init}, s, d_S, \bar{a}} \left(F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\beta(t)}^0(s, d_S, \bar{a}) \right)^2, \end{aligned} \quad (\text{E.18})$$

where in the last line we used Cauchy-Schwartz inequality, together with the fact that (s, d_S, \bar{a}) and (s', d'_S, \bar{a}') share the same distribution. That is

$$\begin{aligned} & \mathbb{E}_{\text{init}, s, d_S, \bar{a}, s', d'_S, \bar{a}'} \left[F_{\beta(t)}^0(s', d'_S, \bar{a}') - F_{\beta_*}^0(s', d'_S, \bar{a}') \right] \left[F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) \right] \\ &\leq \left(\mathbb{E} \left[F_{\beta(t)}^0(s', d'_S, \bar{a}') - F_{\beta_*}^0(s', d'_S, \bar{a}') \right]^2 \mathbb{E} \left[F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) \right]^2 \right)^{1/2} \\ &= \mathbb{E}_{\text{init}, s, d_S, \bar{a}} \left[F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) \right]^2. \end{aligned}$$

Combining (E.16), (E.17), and (E.18), the second term in (E.15) can be bounded by

$$\begin{aligned} & -\eta \mathbb{E} \left\langle \beta(t) - \beta_*, \hat{g}_{\beta(t)} - \bar{g}_{\beta_*}^0 \right\rangle \\ &\leq -\eta \mathbb{E} \left[-R_\beta \|\bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0\|_2 + (1 - \mu) \mathbb{E} \left(F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\beta(t)}^0(s, d_S, \bar{a}) \right)^2 \right] \\ &\quad + \eta R_\beta \left(\mathbb{E} \|\hat{g}_{\beta(t)} - g_{\beta(t)}\|_2^2 \right)^{1/2}. \end{aligned} \quad (\text{E.19})$$

To bound the third term in (E.15), we have

$$\mathbb{E} \|\hat{g}_{\beta(t)} - \bar{g}_{\beta_*}^0\|_2^2 \leq 3 \mathbb{E} \left(\|\hat{g}_{\beta(t)} - \bar{g}_{\beta(t)}\|_2^2 + \|\bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0\|_2^2 + \|\bar{g}_{\beta(t)}^0 - \bar{g}_{\beta_*}^0\|_2^2 \right). \quad (\text{E.20})$$

We can use Lemma E.5 and Lemma E.2 to control the first two terms in (E.20). We proceed to bound the third one, note that $\nabla F_\beta^0(s, d_S, \bar{a}) = \nabla F_{\beta(0)}^0(s, d_S, \bar{a})$ for any β , and $\|\nabla F_{\beta(0)}^0(s, d_S, \bar{a})\|_2 \leq 1$.

$$\begin{aligned} \|\bar{g}_{\beta(t)}^0 - \bar{g}_{\beta_*}^0\|_2^2 &\leq \mathbb{E} \left(\delta_{\beta(t)}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right) \left\| \nabla F_{\beta(0)}^0(s, d_S, \bar{a}) \right\|_2^2 \\ &\leq \mathbb{E} \left[\left(\delta_{\beta(t)}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right) \left\| \nabla F_{\beta(0)}^0(s, d_S, \bar{a}) \right\|_2 \right]^2 \\ &\leq \mathbb{E} \left[\delta_{\beta(t)}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') - \delta_{\beta_*}^0(s, d_S, \bar{a}, s', d'_S, \bar{a}') \right]^2 \\ &\leq \mathbb{E} \left(F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) - \mu \left[F_{\beta(t)}^0(s', d'_S, \bar{a}') - F_{\beta_*}^0(s', d'_S, \bar{a}') \right] \right)^2 \\ &\leq 2(1 + u^2) \mathbb{E} \left[F_{\beta(t)}^0(s, d_S, \bar{a}) - F_{\beta_*}^0(s, d_S, \bar{a}) \right]^2, \end{aligned} \quad (\text{E.21})$$

where in the last inequality follows from (s, d_S, \bar{a}) and (s', d'_S, \bar{a}') sharing the same distribution.

Combining (E.19), (E.20) and (E.21), conditioned on $\beta(t)$, we have the following bound for (E.15):

$$\begin{aligned} & \mathbb{E} \left[\|\beta(t+1) - \beta_*\|_2^2 | \beta(t) \right] \\ &\leq \|\beta(t) - \beta_*\|_2^2 - (\eta(1 - \mu) - 6(1 + \mu^2 \eta^2) \mathbb{E} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\beta(t)}^0(s, d_S, \bar{a}) \right]^2) \\ &\quad + \eta R_\beta \left(\mathbb{E} \|\hat{g}_{\beta(t)} - g_{\beta(t)}\|_2^2 \right)^{1/2} + \eta R_\beta \|\bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0\|_2 \\ &\quad + 3\eta^2 \mathbb{E} \left(\|\hat{g}_{\beta(t)} - \bar{g}_{\beta(t)}\|_2^2 + \|\bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0\|_2^2 \right). \end{aligned}$$

From Lemma E.5, we have $\mathbb{E}\|\hat{g}_{\beta(t)} - g_{\beta(t)}\|_2^2 = \mathcal{O}(R_\beta^2/N)$. From Lemma E.2, we have $\mathbb{E}\|\bar{g}_{\beta(t)} - \bar{g}_{\beta(t)}^0\|_2 = \mathcal{O}(R_\beta^{3/2}/m^{1/4})$. From Lemma E.5, we have $\mathbb{E}\|\hat{g}_{\beta(t)} - \bar{g}_{\beta(t)}\|_2^2 = \mathcal{O}(R_\beta^2)$. Hence

$$\begin{aligned} & \mathbb{E}\|\beta(t+1) - \beta_*\|_2^2 \\ & \leq \|\beta(t) - \beta_*\|_2^2 - (\eta(1-\mu) - 6(1+\mu^2)\eta^2) \mathbb{E} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\beta(t)}^0(s, d_S, \bar{a}) \right]^2 \\ & \quad + \eta \mathcal{O} \left(\frac{R_\beta^2}{N^{1/2}} + \frac{R_\beta^{5/2}}{m^{1/4}} \right) + \eta^2 \mathcal{O}(R_\beta^2). \end{aligned}$$

Re-arrange and telescope, we have

$$\begin{aligned} & \sum_{t=0}^{T-1} \mathbb{E} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\beta(t)}^0(s, d_S, \bar{a}) \right]^2 \\ & \leq (\eta(1-\mu) - 6(1+\mu^2)\eta^2)^{-1} \mathbb{E} \left[\|\beta_0 - \beta_*\|_2^2 + \eta T \mathcal{O} \left(\frac{R_\beta^2}{N^{1/2}} + \frac{R_\beta^{5/2}}{m^{1/4}} \right) + T \eta^2 \mathcal{O}(R_\beta^2) \right]. \end{aligned} \tag{E.22}$$

Notice that the left hand side of (E.22) is convex with respect to $\beta(t)$, we can take average of $\{\beta(t)\}_{t=0}^{T-1}$. Define $\bar{\beta}_T = \frac{1}{T} \sum_{t=0}^{T-1} \beta(t)$, with Jensen's inequality we have

$$\begin{aligned} & \mathbb{E} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\bar{\beta}_T}^0(s, d_S, \bar{a}) \right]^2 \\ & \leq (\eta(1-\mu) - 6(1+\mu^2)\eta^2)^{-1} \mathbb{E}_{\text{init}} \left[\frac{\|\beta_0 - \beta_*\|_2^2}{T} + \eta \mathcal{O} \left(\frac{R_\beta^2}{N^{1/2}} + \frac{R_\beta^{5/2}}{m^{1/4}} \right) + \eta^2 \mathcal{O}(R_\beta^2) \right]. \end{aligned}$$

Now take $\eta = \mathcal{O}(T^{-1/2})$, we have

$$\mathbb{E} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\bar{\beta}_T}^0(s, d_S, \bar{a}) \right]^2 \leq \mathcal{O} \left(\frac{R_\beta^2}{T^{1/2}} + \frac{R_\beta^{5/2}}{m^{1/4}} + \frac{R_\beta^2}{N^{1/2}} \right).$$

We can now apply Lemma E.1 and conclude that:

$$\begin{aligned} \mathbb{E} \left[F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\bar{\beta}_T}^0(s, d_S, \bar{a}) \right]^2 & \leq 2 \mathbb{E} \left[\left(F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\bar{\beta}_T}^0(s, d_S, \bar{a}) \right)^2 + \left(F_{\beta_*}^0(s, d_S, \bar{a}) - F_{\bar{\beta}_T}^0(s, d_S, \bar{a}) \right)^2 \right] \\ & \leq \mathcal{O} \left(\frac{R_\beta^2}{T^{1/2}} + \frac{R_\beta^{5/2}}{m^{1/4}} + \frac{R_\beta^2}{N^{1/2}} + \frac{R_\beta^3}{m^{1/2}} \right). \end{aligned}$$

□

We can now specialize Theorem E.1 to policy optimization and policy evaluation.

Proof of Lemma 4.1. Note that F_{θ_*} satisfies $F_{\theta_*}^0 = \Pi_{\mathcal{F}_{R_\theta, m_Q}^{\mathcal{P}}}(\zeta_{F_{\theta_*}^0})$. For policy evaluation, by definition (E.2), the operator ζ equals to the Bellman evaluation operator: $\zeta_{F_\theta}(s, d_S, \bar{a}) = [\mathcal{T}^{\pi_k} F_\theta](s, d_S, \bar{a})$. Hence we have $F_{\theta_*}^0 = \Pi_{\mathcal{F}_{R_\theta, m_Q}^{\mathcal{P}}}([\mathcal{T}^{\pi_k} F_{\theta_*}^0])$. Now by Assumption 1, we know that $[\mathcal{T}^{\pi_k} F_{\theta_*}^0] \in \mathcal{F}_{R_\theta, m_Q}^{\mathcal{P}}$, together with the fact that $F_{\theta_*}^0 \in \mathcal{F}_{R_\theta, m_Q}^{\mathcal{P}}$. By uniqueness of the projection, we have $F_{\theta_*}^0 = [\mathcal{T}^{\pi_k} F_{\theta_*}^0]$, which implies $F_{\theta_*}^0 = Q^{\pi_k}$. The claim follows immediately by applying Theorem E.1. □

Proof of Lemma 4.2. Note that here we have $F_{\alpha_*}^0 = \Pi_{F_{R_\alpha, m_A}}(\zeta_{F_{\alpha_*}^0})$. For policy optimization the operator ζ is defined by: $\zeta_{F_{\alpha_*}^0}(s, d_S, \bar{a}) = \tau_{k+1}(v_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a}))$. Then we

have

$$\begin{aligned}
& \mathbb{E}_{\tilde{\sigma}_k} \left[F_{\alpha(T)}^A(s, d_S, \bar{a}) - \tau_{k+1}(v_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a})) \right]^2 \\
& \leq 2\mathbb{E}_{\tilde{\sigma}_k} \left[F_{\alpha(T)}^A(s, d_S, \bar{a}) - F_{\alpha_*}^0(s, d_S, \bar{a}) \right]^2 + 2\mathbb{E}_{\tilde{\sigma}_k} \left[F_{\alpha_*}^0(s, d_S, \bar{a}) - \tau_{k+1}(v_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a})) \right]^2 \\
& \leq 2\mathbb{E}_{\tilde{\sigma}_k} \left[\tau_{k+1}(v_k^{-1} F_{\theta_k}^0(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^0(s, d_S, \bar{a})) - \tau_{k+1}(v_k^{-1} F_{\theta_k}^Q(s, d_S, \bar{a}) + \tau_k^{-1} F_{\alpha_k}^A(s, d_S, \bar{a})) \right]^2 \\
& \quad + 2\mathbb{E}_{\tilde{\sigma}_k} \left[F_{\alpha(T)}^A(s, d_S, \bar{a}) - F_{\alpha_*}^0(s, d_S, \bar{a}) \right]^2 \\
& \leq \mathcal{O} \left(\frac{R_\alpha^2}{T^{1/2}} + \frac{R_\alpha^{5/2}}{m_A^{1/4}} + \frac{R_\alpha^2}{N^{1/2}} + \frac{R_\alpha^3}{m_A^{1/2}} \right),
\end{aligned}$$

where the second inequality comes from the definition of projection and $F_{\alpha_*}^0 = \Pi_{F_{R_\alpha}, m_A}(\zeta_{F_{\alpha_*}^0})$, and the last inequality comes from direct application of Theorem E.1 and Lemma E.1, together with the fact that F_α^A and F_θ^Q share the same initialization. \square

Proof of Theorem 4.1. From Lemma 4.3, we know that

$$\begin{aligned}
& (1 - \gamma) (\mathcal{L}(\pi^*) - \mathcal{L}(\pi_k)) \\
& \leq \mathbb{E}_{\nu^*} \left\{ v_k [\text{KL}(\pi^*(\cdot|s, d_S) \|\pi_k(\cdot|s, d_S)) - \text{KL}(\pi^*(\cdot|s, d_S) \|\pi_{k+1}(\cdot|s, d_S)) + \varepsilon_k + \varepsilon'_k] + v_k^{-1} M \right\}.
\end{aligned} \tag{E.23}$$

where ε_k and ε'_{k+1} are defined as in Lemma 4.1 and Lemma 4.2, respectively. In addition, $\varepsilon_k = \tau_{k+1}^{-1} \epsilon_{k+1} \phi_{k+1}^* + v_k^{-1} \epsilon'_k \psi_k^*$, $\varepsilon'_k = |\mathcal{A}| \tau_{k+1}^{-2} \epsilon_{k+1}^2$, and $M = \mathbb{E}_{\nu^*} \left[\max_{a \in \mathcal{A}} (F_{\theta_0}^Q(s, d_S, \bar{a}))^2 \right] + 2R_\alpha^2$. ϕ_k^* and ψ_k^* are defined by:

$$\begin{aligned}
\phi_k^* &= \mathbb{E}_{\tilde{\sigma}_k} [|\text{d}\pi^*/\text{d}\pi_0 - \text{d}\pi_k/\text{d}\pi_0|^2]^{1/2}, \\
\psi_k^* &= \mathbb{E}_{\sigma_k} [|\text{d}\sigma^*/\text{d}\sigma_k - \text{d}(\nu^* \times \pi_k)/\text{d}\sigma_k|^2]^{1/2}.
\end{aligned}$$

Now sum up (E.23) from $k = 0$ to $K - 1$, with $v_k = v\sqrt{K}$, we have

$$\begin{aligned}
& (1 - \gamma) K \min_{0 \leq k \leq K-1} (\mathcal{L}(\pi^*) - \mathcal{L}(\pi_k)) \\
& \leq \mathbb{E}_{\nu^*} \left\{ v\sqrt{K} \left[\text{KL}(\pi^*(\cdot|s, d_S) \|\pi_0(\cdot|s, d_S)) - \text{KL}(\pi^*(\cdot|s, d_S) \|\pi_K(\cdot|s, d_S)) + \sum_{k=0}^{K-1} (\varepsilon_k + \varepsilon'_k) \right] + v^{-1} \sqrt{K} M \right\}.
\end{aligned}$$

Since we initialize policy π_0 to be uniform policy, we have $\text{KL}(\pi^*(\cdot|s, d_S) \|\pi_0(\cdot|s, d_S)) \leq \log|\mathcal{A}|$. Rearrange, we obtain

$$\min_{0 \leq k \leq K} \{\mathcal{L}(\pi^*) - \mathcal{L}(\pi_k)\} \leq \frac{v^2 \left(\log|\mathcal{A}| + \sum_{k=1}^{K-1} (\varepsilon_k + \varepsilon'_k) \right) + M}{(1 - \gamma)v\sqrt{K}}, \tag{E.24}$$

From Lemma 4.1 and Lemma 4.2, we have

$$\epsilon_k = \mathcal{O} \left(\frac{R_\theta^2}{T^{1/2}} + \frac{R_\theta^{5/2}}{m_Q^{1/4}} + \frac{R_\theta^2}{N^{1/2}} + \frac{R_\theta^3}{m_Q^{1/2}} \right), \epsilon'_{k+1} = \mathcal{O} \left(\frac{R_\alpha^2}{T^{1/2}} + \frac{R_\alpha^{5/2}}{m_A^{1/4}} + \frac{R_\alpha^2}{N^{1/2}} + \frac{R_\alpha^3}{m_A^{1/2}} \right).$$

To control ε_k and ε'_k , we have

$$\begin{aligned}
\tau_{k+1}^{-1} \epsilon'_{k+1} \phi_{k+1}^* &= \mathcal{O} \left(kK^{-1} \phi_k^* (R_\alpha^2 T^{-1/2} + R_\alpha^{5/2} m_A^{-1/4} + R_\alpha^2 N^{-1/2}) \right) \\
|\mathcal{A}| \tau_{k+1}^{-2} (\epsilon'_{k+1})^2 &= \mathcal{O} \left(k^2 K^{-1} |\mathcal{A}| (R_\alpha^2 T^{-1/2} + R_\alpha^{5/2} m_A^{-1/4} + R_\alpha^2 N^{-1/2})^2 \right) \\
v_k^{-1} \epsilon_k \psi_k^* &= \mathcal{O} \left(K^{-1/2} \psi_k^* (R_\theta^2 T^{-1/2} + R_\theta^2 N^{-1/2} + R_\theta^{5/2} m^{-1/4}) \right),
\end{aligned}$$

if $m_A = \Omega(R_\alpha^2)$ and $m_Q = \Omega(R_\theta^2)$. One can verify that with the choice of N, T, m_A, m_Q in Theorem 4.1, we have $\varepsilon_k = \varepsilon'_k = \mathcal{O}(K^{-1})$. Plug into (E.24), we obtain the second part of Theorem 4.1. \square