

Offline Multi-Agent Reinforcement Learning via In-Sample Sequential Policy Optimization

Zongkai Liu^{1,3}, Qian Lin¹, Chao Yu^{1,2*}, Xiawei Wu¹, Yile Liang⁴, Donghui Li⁴, Xuetao Ding⁴

¹School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou, China

²Pengcheng Laboratory, Shenzhen, China

³Shanghai Innovation Institute, Shanghai, China

⁴Meituan, Beijing, China

{liuzk, linq67}@mail2.sysu.edu.cn, yuchao3@mail.sysu.edu.cn

Abstract

Offline Multi-Agent Reinforcement Learning (MARL) is an emerging field that aims to learn optimal multi-agent policies from pre-collected datasets. Compared to single-agent case, multi-agent setting involves a large joint state-action space and coupled behaviors of multiple agents, which bring extra complexity to offline policy optimization. In this work, we revisit the existing offline MARL methods and show that in certain scenarios they can be problematic, leading to uncoordinated behaviors and out-of-distribution (OOD) joint actions. To address these issues, we propose a new offline MARL algorithm, named In-Sample Sequential Policy Optimization (InSPO). InSPO sequentially updates each agent's policy in an in-sample manner, which not only avoids selecting OOD joint actions but also carefully considers teammates' updated policies to enhance coordination. Additionally, by thoroughly exploring low-probability actions in the behavior policy, InSPO can well address the issue of premature convergence to sub-optimal solutions. Theoretically, we prove InSPO guarantees monotonic policy improvement and converges to quantal response equilibrium (QRE). Experimental results demonstrate the effectiveness of our method compared to current state-of-the-art offline MARL methods.

Code — <https://github.com/kkkaiaiai/InSPO/>

Introduction

Offline Reinforcement Learning (RL) is a rapidly evolving field that aims to learn optimal policies from pre-collected datasets without interacting directly with the environment (Figueiredo Prudencio, Maximo, and Colombini 2024). The primary challenge in offline RL is the issue of distributional shift (Yang et al. 2021), which occurs when policy evaluation on out-of-distribution (OOD) samples leads to the accumulation of extrapolation errors. Existing research usually tackles this problem by employing conservatism principles, compelling the learning policy to remain close to the data manifold through various data-related regularization techniques (Yang et al. 2021; Pan et al. 2022; Matsunaga et al. 2023; Shao et al. 2023; Wang et al. 2023b).

*Corresponding author.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

In comparison to the single-agent counterpart, offline Multi-Agent Reinforcement Learning (MARL) has received relatively less attention. Under the multi-agent setting, it not only faces the challenges inherent to offline RL but also encounters common MARL issues, such as difficulties in coordination and large joint-action spaces (Zhang, Yang, and Başar 2021). These issues cannot be simply resolved by combining state-of-the-art offline RL solutions with modern multi-agent techniques (Yang et al. 2021). In fact, due to the increased number of agents and the offline nature of the problem, these issues become even more challenging. For example, under the offline setting, even if each agent selects an in-sample action, the resulting joint action may still be OOD. Additionally, in cooperative MARL, agents need to consider both their own actions and the actions of other agents in order to determine their contributions to the global return for high overall performance. Thus, under offline settings, discovering and learning cooperative joint policies from the dataset poses a unique challenge for offline MARL.

To address the aforementioned issues, recent works have developed specific offline MARL algorithms. These approaches generally integrate the conservatism principle into the Centralized Training with Decentralized Execution (CTDE) framework, such as value decomposition structures (Yang et al. 2021; Pan et al. 2022; Matsunaga et al. 2023; Shao et al. 2023; Wang et al. 2023b), which is developed under the Individual-Global-Max (IGM) assumption. Although these approaches have demonstrated successes in certain offline multi-agent tasks, they still exhibit several limitations. For example, due to the inherent limitations of the IGM principle, algorithms that utilize value decomposition structures may struggle to find optimal solutions because of constraints in their representation capabilities, and can even lead to the selection of OOD joint actions, as we show in the Proposed Method section.

In this work, we propose a principled approach to tackle OOD joint actions issue. By introducing a behavior regularization into the policy learning objective and derive the closed-form solution of the optimal policy, we develop a sequential policy optimization method in an entirely in-sample learning manner without generating potentially OOD actions. Besides, the sequential update scheme used in this method enhances both the representation capabil-

ity of the joint policy and the coordination among agents. Then, to prevent premature convergence to local optima, we encourage sufficient exploration of low-probability actions in the behavior policy through the use of policy entropy. The proposed novel algorithm, named the In-Sample Sequential Policy Optimization (InSPO), enjoys the properties of monotonic improvement and convergence to quantal response equilibrium (QRE) (McKelvey and Palfrey 1995), a solution concept in game theory. We evaluate InSPO in the XOR game, Multi-NE game, and Bridge to demonstrate its effectiveness in addressing OOD joint action and local optimum convergence issues. Additionally, we test it on various types of offline datasets in the StarCraft II micromanagement benchmark to showcase its competitiveness with current state-of-the-art offline MARL algorithms.

Related Work

MARL. The CTDE framework dominates current MARL research, facilitating agent cooperation. In CTDE, agents are centrally trained using global information but rely only on local observations to make decisions. Value decomposition is a notable method, representing the joint Q-function as a combination of individual agents' Q-functions (Wang et al. 2021; Son et al. 2019; Rashid et al. 2018). These methods typically depend on the IGM principle, assuming the optimal joint action corresponds to each agent's greedy actions. However, environments with multi-modal reward landscapes frequently violate the IGM assumption, limiting the effectiveness of value decomposition in learning optimal policies (Fu et al. 2022).

Another influential class of methods is Multi-Agent Policy Gradient (MAPG), with notable algorithms such as MAPPO(Yu et al. 2022), CoPPO(Wu et al. 2021), and HAPPO(Kuba et al. 2022). However, on-policy learning approaches like these struggle in offline settings due to OOD action issues, leading to extrapolation errors.

Offline MARL. OMAR (Pan et al. 2022) combines Independent Learning and zeroth-order optimization to adapt CQL (Kumar et al. 2020) for multi-agent scenarios. However, OMAR fundamentally follows a single-agent learning paradigm, which treats other agents as part of the environment, and does not handle cooperative behavior learning and OOD joint actions insufficiently.

To enhance cooperation and efficiency in complex environments such as StarCraft II, some existing works employ value decomposition as a foundation for algorithm design. For instance, ICQ (Yang et al. 2021) introduces conservatism to prevent optimization on unseen state-action pairs, mitigating extrapolation errors. OMIGA (Wang et al. 2023b) and CFCQL (Shao et al. 2023) are the latest offline MARL methods, both integrating value decomposition structures. OMIGA applies implicit local value regularization to enable in-sample learning, while CFCQL calculates counterfactual regularization per agent, avoiding the excessive conservatism caused by direct value decomposition-CQL integration. Nonetheless, the IGM principle has been shown to fail in identifying optimal policies in multi-modal reward landscapes (Fu et al. 2022), due to the limited expressive-

ness of the Q-value network, which poses a potential risk of encountering the OOD joint actions issue in offline settings.

An alternative research direction in offline RL applies constraints on state-action distributions, called DIstribution Correction Estimation (DICE) methods (Figueiredo Prudencio, Maximo, and Colombini 2024). AlberDICE (Matsunaga et al. 2023) is a pioneering DICE-based method in offline MARL, which is proved to converge to NEs. However, when multiple NEs exist, its convergence results heavily depends on the dataset distribution. If the behavior policy is near a sub-optimal NE, AlberDICE will converge directly to that sub-optimal solution rather than the global optimum. This is primarily because AlberDICE lacks sufficient exploration of low-probability state-action pairs in dataset, leading to premature convergence to a deterministic policy. Additionally, AlberDICE employs an out-of-sample learning during policy extraction, i.e., it uses actions produced by the policy rather the actions in datasets, which could lead to OOD joint actions (Xu et al. 2023a; Kostrikov, Nair, and Levine 2022).

Additionally, some works consider using model-based method (Barde et al. 2024) or using diffusion models (Li, Pan, and Huang 2023; Zhu et al. 2023) to solve the OOD action issue. More discussion about the related work is given in Appendix E.

Background

Cooperative Markov Game

The cooperative MARL problem is usually modeled as a cooperative Markov game (Littman 1994) $\mathcal{G} = \langle \mathcal{N}, \mathcal{S}, \mathcal{A}, P, r, \gamma, d \rangle$, where $\mathcal{N} = \{1, \dots, N\}$ is the set of agent indices, \mathcal{S} is the finite state space, $\mathcal{A} = \prod_{i \in \mathcal{N}} \mathcal{A}_i$ is the joint action space, with \mathcal{A}_i denoting the finite action space of agent i , $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the common reward function shared with all agents, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the transition probability function, $\gamma \in [0, 1]$ is the discount factor, and $d \in \Delta(\mathcal{S})$ is the initial state distribution. At time step $t \in \{1, \dots, T\}$, each agent $i \in \mathcal{N}$ at state $s_t \in \mathcal{S}$ selects an action $a_t^i \sim \pi^i(\cdot | s_t)$ and moves to the next state $s_{t+1} \sim P(\cdot | s_t, a_t)$. It then receives a reward $r_t = r(s_t, a_t)$ according to the joint action $a_t = \{a_t^1, \dots, a_t^N\}$. We denote the joint policy as $\pi(\cdot | s) = \prod_{i \in \mathcal{N}} \pi^i(\cdot | s)$, and the joint policy except the i -th player as π^{-i} . In a cooperative Markov game, all agents aim to learn a optimal joint policy π that jointly maximizes the expected discount returns $\mathbb{E}_{s \in \mathcal{S}, a \sim \pi} [\sum_{t=0}^T \gamma^t r(s_t, a_t)]$. Under the offline setting, only a pre-collected dataset $\mathcal{D} = \{(s, a, r, s')_k\}_{k=1}^{|\mathcal{D}|}$ collected by an unknown behavior policy $\mu = \prod_{i \in \mathcal{N}} \mu^i$ is given and the environment interactions are not allowed.

IGM Principle and Value Decomposition

Value-based methods aim to learn a joint Q-function $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ to estimate the future expected return given the current state s and joint action a . However, directly computing the joint Q-function is challenging due to the huge state-action space in MARL. To address this issue, value decomposition decomposes the joint Q-function Q into individual Q-functions Q^i for each agent: $Q(s, a) =$

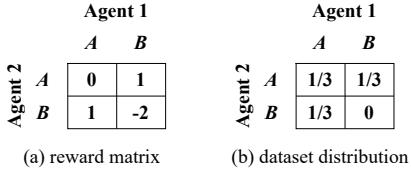


Figure 1: XOR game. (a) is the reward matrix of joint actions. (b) is the distribution of dataset.

$f_{\text{mix}}(Q^1(s, a^1), \dots, Q^N(s, a^N); s)$, where f_{mix} represents the mixing function conditioned on the state (Fu et al. 2022). The mixing function f_{mix} must satisfy the IGM principle that any optimal joint action a_* should satisfy

$$a_* = \arg \max_{a \in \mathcal{A}} Q(s, a) = \bigcup_{i \in \mathcal{N}} \{\arg \max_{a^i \in \mathcal{A}^i} Q^i(s, a^i)\}. \quad (1)$$

Under the IGM assumption, value decomposition enables the identification of the optimal joint action through the greedy actions of each agent.

Behavior-Regularized Markov Game in Offline MARL

Behavior-Regularized Markov Game is a useful framework to avoid distribution shift by incorporating a data-related regularization term on rewards (Wang et al. 2023b; Xu et al. 2023a). In this framework, the goal is to optimize policy by

$$\max_{\pi} \mathbb{E} \left[\sum_{t=1}^T \gamma^t \left(r(s_t, a_t) - \alpha f(\pi(\cdot|s_t), \mu(\cdot|s_t)) \right) \right], \quad (2)$$

where $f(\cdot, \cdot)$ is a regularization function, and $\alpha \geq 0$ is a temperature constant. The unknown behavior policy μ here can usually be approximated by using Behavior Cloning (Wang et al. 2023b; Xu et al. 2023a). Policy evaluation operator in this framework is given by

$$\mathcal{T}_{\pi} Q_{\pi}(s, a) \triangleq r(s, a) + \gamma \mathbb{E}_{s'|s, a} [V_{\pi}(s')], \quad (3)$$

$$\text{where } V_{\pi}(s) = \mathbb{E}_{a \sim \pi} [Q_{\pi}(s, a) - \alpha f(\pi(a|s), \mu(a|s))].$$

Thus, the objective (2) can be represented as

$$\max_{\pi} \mathbb{E}_{a \sim \pi} [Q_{\pi}(s, a) - \alpha f(\pi(\cdot|s), \mu(\cdot|s))]. \quad (4)$$

The Proposed Method

OOD Joint Action in Offline MARL

In offline MARL, value decomposition methods are more prone to encountering OOD joint actions due to the constraints of the IGM principle in certain scenarios. We use the XOR game, shown in Figure 1, to illustrate this phenomenon. Figure 1(a) shows the reward matrix of the XOR game, while Figure 1(b) depicts the dataset considered in the offline setting. Since it is necessary to minimize temporal difference (TD) error $\mathbb{E}_{\mathcal{D}}[(f_{\text{mix}}(Q^1(a^1), Q^2(a^2)) - r(a^1, a^2))^2]$ while satisfying the IGM principle, the local Q-functions for both agents are forced to satisfy $Q^i(B) >$

$Q^i(A)$, $i = 1, 2$ (See Appendix D for a detailed derivation). As a result, both agents tend to choose action B , resulting in the OOD joint action (B, B) .

Another line in offline MARL research combines MAPG methods and data-related regularization (Pan et al. 2022). However, they can still encounter the OOD joint actions issue although not constrained by the IGM principle. Considering again the above offline task, both learned agents are likely to choose (A, A) due to the data-related regularization. For agent 1, given that its teammate selects action A , choosing action B would yield a higher payoff. The same is true for agent 2, resulting in the OOD joint action (B, B) . This situation arises because these methods do not fully consider the change of teammates' policies, leading to conflicting directions in policy updates.

MAPG methods employing sequential update scheme can effectively address this issue, as they fully consider the direction of teammates' policy updates, thereby avoiding conflicts (Matsunaga et al. 2023; Kuba et al. 2022). In the same scenario as above, but with sequential updates, where agent 1 updates first followed by agent 2, agent 1 would still choose action B for a higher payoff. Then, when agent 2 updates, knowing that agent 1 chose B , it would find that sticking with action A is best. Consequently, sequential-update MAPG methods converge to the optimal policy.

In-Sample Sequential Policy Optimization

Inspired by the above discussions, we introduce an in-sample sequential policy optimization method under the behavior-regularized Markov game framework, i.e., Eq.(4). Here we consider the reverse KL divergence as the regularization, which means $f(x, y) = \log(\frac{x}{y})$. The benefit of choosing reverse KL divergence is that the global regularization can be decomposed naturally as $\log(\frac{\pi}{\mu}) = \sum_{i \in \mathcal{N}} \log(\frac{\pi^i}{\mu^i})$, making the simplified computation of sequential-update possible. Denoting $i_{1:n}$ as an ordered subset $\{i_1, \dots, i_n\}$ of \mathcal{N} , and $-i_{1:n}$ as its complement, where i_k is the k -th agent in the ordered subset and $i_{1:0} = \emptyset$, the sequential-update objectives are given by:

$$\pi_{\text{new}}^{i_n} = \arg \max_{\pi^{i_n}} \mathbb{E}_{a^{i_n} \sim \pi^{i_n}} \left[Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n}) - \alpha \log(\frac{\pi^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)}) \right], \quad (5)$$

where

$$Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n}) \triangleq \mathbb{E}_{\pi_{\text{new}}^{i_{1:n}-1}, \pi_{\text{old}}^{-i_{1:n}}} [Q_{\pi_{\text{old}}}(s, a^{-i_n}, a^{i_n})].$$

However, the optimization objective (5) requires actions produced by the policy, which is in a out-of-sample learning manner, potentially leading to OOD actions. In order to achieve in-sample learning using only the dataset actions, we derive the closed-form solution of objectives (5) by the Karush-Kuhn-Tucker (KKT) conditions

$$\pi_{\text{new}}^{i_n}(a^{i_n}|s) \propto \mu^{i_n}(a^{i_n}|s) \cdot \exp \left(\frac{Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n})}{\alpha} \right), \quad (6)$$

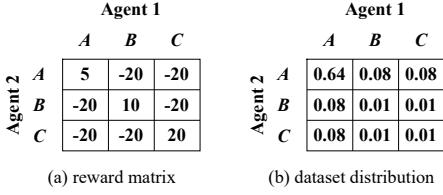


Figure 2: M-NE game. (a) is the reward matrix of joint actions. (b) is the distribution of dataset.

and thus obtain the in-sample optimization objectives for parametric policy $\pi_{\theta^{in}}$ by minimizing the KL divergence:

$$\begin{aligned} \theta_{\text{new}}^{in} &= \arg \min_{\theta^{in}} D_{\text{KL}}(\pi_{\text{new}}^{in}(\cdot|s), \pi_{\theta^{in}}(\cdot|s)) \\ &= \arg \min_{\theta^{in}} \mathbb{E}_{(s, a^{in}) \sim \mathcal{D}} \left[-\exp \left(\frac{A_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{in})}{\alpha} \right) \right. \\ &\quad \left. \cdot \log \pi_{\theta^{in}}(a^{in}|s) \right], \end{aligned} \quad (7)$$

where $A_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{in}) \triangleq Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{in}) - \mathbb{E}_{\pi_{\text{new}}^{in}}[Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{in})]$.

A potential problem of this method is that it may lead to premature convergence to local optima due to exploitation of vested interests. This concern is especially pronounced when the behavior policy is a local optimum, as we will show in the next subsection.

Maximum-Entropy Behavior-Regularized Markov Game

The existence of multiple local optima is a common phenomenon in many multi-agent tasks, where finding the global optimum is often extremely challenging. Therefore, near-optimal (or expert) behavior policies can easily fall into or stay near local optima. In such cases, because the data-related regularization enforces the learned policy to remain close to the behavior policy, optimizing the objective in Eq.(5) is more likely to cause the sequential policy optimization method to converge towards a deterministic policy that exploits this local optimum. Moreover, escaping this local optimum becomes challenging, as when one of the agents attempts to deviate unilaterally, the optimization objective (5) impedes this since it hurts the overall benefits.

We examine this issue using the M-NE game depicted in Figure 2, with Figure 2(a) showing the reward matrix and Figure 2(b) illustrating the offline dataset. In this game, there are three NEs: (A, A) , (B, B) , and (C, C) , with rewards of 5, 10, and 20, respectively, where (C, C) represents the global optimal NE and other NEs are local optima. On the considered dataset, data-related regularization enforces agents to select A with a high probability. As a result, agents confidently converge to the local optimum (A, A) based on the observed high probability of their teammates choosing A , failing to recognize the optimal joint action (C, C) .

One way to address this issue is to introduce perturbations to the rewards, preventing sequential policy optimization method from deterministically converging to a local optimum and thereby encouraging it to escape the local optimum and identify the global optimal solution. From a game-

theoretic perspective, the optimal solution of the perturbed game aligns with the solution concept of quantal response equilibrium (QRE) (McKelvey and Palfrey 1995).

Definition 1. For a Behavior-Regularized Markov Game \mathcal{G} with a reward function r , denote the perturbed reward as \tilde{r} . Then, a joint policy π_* is a QRE if it holds

$$\tilde{J}(\pi_*) \geq \tilde{J}(\pi_*^{-i}, \pi^i), \quad \forall i \in \mathcal{N}, \pi^i, \quad (8)$$

where $\tilde{J}(\pi) \triangleq \mathbb{E}_\pi[\sum_t \gamma^t (\tilde{r}(s_t, \mathbf{a}_t) - \alpha f(\pi(\cdot|s_t), \mu(\cdot|s_t)))]$.

Therefore, our goal is to design an in-sample sequential policy optimization method with QRE convergence guarantees. One simple and effective way to introduce disturbances is to add policy entropy into the rewards, which is also a commonly used regularization in online RL to improve exploration (Liu et al. 2024; Haarnoja et al. 2018). Therefore, we introduce the following Maximum-Entropy Behavior-Regularized Markov Game (MEBR-MG) problem, which is a generalization of Behavior-Regularized Markov Game (2).

$$\begin{aligned} \max_{\pi} \mathbb{E} \left[\sum_{t=1}^T \gamma^t \left(r(s_t, \mathbf{a}_t) - \alpha D_{\text{KL}}(\pi(\cdot|s_t), \mu(\cdot|s_t)) \right. \right. \\ \left. \left. + \beta \mathcal{H}(\pi(\cdot|s_t)) \right) \right], \end{aligned} \quad (9)$$

where $\mathcal{H}(\pi(\cdot|s_t))$ is policy entropy, and $\beta \geq 0$ is a temperature constant. In the following context, we first give some facts about MEBR-MG, and then give the in-sample sequential policy optimization method under MEBR-MG.

In MEBR-MG, we have the following modified policy evaluation operator given by:

$$\mathcal{T}_\pi Q_\pi(s, \mathbf{a}) \triangleq r(s, \mathbf{a}) + \gamma \mathbb{E}_{s'|s, \mathbf{a}}[V_\pi(s')], \quad (10)$$

where

$$V_\pi(s) = \mathbb{E}_{\mathbf{a} \sim \pi} \left[Q_\pi(s, \mathbf{a}) - \sum_{i \in \mathcal{N}} \left(\alpha \log \frac{\pi^i(a^i|s)}{\mu^i(a^i|s)} + \beta \log \pi^i(a^i|s) \right) \right].$$

Lemma 2. Given a policy π , consider the modified policy evaluation operator \mathcal{T}_π in Eq.(10) and a initial Q-function $Q_0 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, and define $Q_{k+1} = \mathcal{T}_\pi Q_k$. Then the sequence Q_k will converge to the Q-function Q_π of policy π as $k \rightarrow \infty$.

Proof can be found in Appendix A. This lemma indicates Q-function will converge to the Q-value under the joint policy π by repeatedly applying the policy evaluation operator.

Moreover, the additional smoothness introduced by the regularization term allows the QRE of MEBR-MG to be expressed in the form of Boltzmann distributions, as demonstrated by the following Proposition 3.

Proposition 3. In a MEBR-MG, a joint policy π_* is a QRE if it holds

$$V_{\pi_*}(s) \geq V_{\pi_*^{-i}, \pi_*^i}(s), \quad \forall i \in \mathcal{N}, \pi^i, s \in \mathcal{S}. \quad (11)$$

Then the QRE policies for each agent i are given by

$$\begin{aligned} \pi_*^i(a^i|s) &\propto \mu^i(a^i|s) \\ &\cdot \exp \left(\frac{\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_*^{-i}}[Q_{\pi_*}(s, a^i, \mathbf{a}^{-i})] - \beta \log \mu^i(a^i|s)}{\alpha + \beta} \right) \end{aligned} \quad (12)$$

Algorithm 1: InSPO

Input: Offline dataset \mathcal{D} , initial policy π_0 and Q-function Q_0
Output: π_K

- 1: Compute behavior policy μ by simple Behavior Cloning
- 2: **for** $k = 1, \dots, K$ **do**
- 3: Compute Q_k by Eq.(10)
- 4: Draw a permutation $i_{1:N}$ of agents at random
- 5: **for** $n = 1, \dots, N$ **do**
- 6: Update $\pi_k^{i_n}$ by Eq.(13)
- 7: **end for**
- 8: **end for**

Proof can be found in Appendix A. Eq.(12) for QRE demonstrates that incorporating the reverse KL divergence term ensures that the learned policy shares the same support set as the behavior policy, thereby avoiding OOD actions; and the addition of an entropy term allows the policy to place greater emphasis on actions with lower probabilities in the behavior policy, preventing premature convergence to local optima.

Similar to Eq.(7), the in-sample sequential policy optimization procedure under MEBR-MG is given by:

$$\begin{aligned} \theta_{\text{new}}^{i_n} = \arg \min_{\theta^{i_n}} & \mathbb{E}_{(s, a^{i_n}) \sim \mathcal{D}} \left[\right. \\ & \left. - \exp \left(\frac{A_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n}) - \beta \log \mu^{i_n}(a^{i_n}|s)}{\alpha + \beta} \right) \cdot \log \pi_{\theta^{i_n}}(a^{i_n}|s) \right]. \end{aligned} \quad (13)$$

Proposition 4. *The sequential policy optimization procedure under MEBR-MG guarantees policy improvement, i.e., $\forall s \in \mathcal{S}, a \in \mathcal{A}$,*

$$Q_{\pi_{\text{new}}}(s, a) \geq Q_{\pi_{\text{old}}}(s, a), V_{\pi_{\text{new}}}(s) \geq V_{\pi_{\text{old}}}(s).$$

Proof can be found in Appendix A. Proposition 4 demonstrates that the policy improvement step defined in Eq.(13) ensures a monotonic increase in performance at each iteration. By alternating between the policy evaluation step and the policy improvement step, we derive InSPO, as shown in Algorithm 1, and we furthermore prove that InSPO converges to QRE as follows.

Theorem 5. *Joint policy π updated by Algorithm 1 converges to QRE.*

Proof can be found in Appendix A.

The Practical Implementation of InSPO

In this section, we design a practical implementation of InSPO to handle the issue of large state-action space, making it more suitable for offline MARL. More details can be found in Appendix B.

Policy Evaluation. According to Eq.(10), we need to train a global Q-function to estimate the expected future return based on the current state and joint action. However, in MARL, the joint action space grows exponentially with the

number of agents. To circumvent this exponential complexity, we instead maintain a local Q-function $Q_{\phi^{i_n}}$ for each agent $i_n \in \mathcal{N}$ to approximate $Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n})$. Besides, $Q_{\phi^{i_n}}$ should be updated sequentially in conjunction with the policy in order to incorporate the information of updated teammates (i.e., $\pi_{\text{new}}^{i_{1:n-1}}$) into the local Q-function. Thus, we optimize the following objective for each local Q-function ϕ^{i_n} :

$$\min_{\phi^{i_n}} \mathbb{E}_{(s, a, s', r) \sim \mathcal{D}} \left[\rho^{i_n} \cdot (Q_{\phi^{i_n}}(s, a^{i_n}) - y)^2 \right], \quad (14)$$

where

$$\begin{aligned} \rho^{i_n} &= \rho^{i_n}(s, a) \triangleq \frac{(\pi_{\text{new}}^{i_{1:n-1}} \cdot \pi_{\text{old}}^{-i_{1:n}})(a^{-i_n}|s)}{\mu^{-i_n}(a^{-i_n}|s)}, \\ y &= y(s, a, s', r) \triangleq r + \gamma \mathbb{E}_{a^{i_{n'}} \sim \pi_{\text{old}}^{i_n}} [Q_{\phi^{i_n}}(s', a^{i_{n'}})] \\ &\quad - \alpha D_{\text{KL}}(\pi_{\text{old}}^{i_n}(\cdot|s'), \mu^{i_n}(\cdot|s')) + \beta \mathcal{H}(\pi_{\text{old}}^{i_n}(\cdot|s')). \end{aligned}$$

Here we omit the regularization terms for other agents to simplify the computation. Furthermore, to reduce the high variance of importance sampling ratio ρ^{i_n} , InSPO adopts importance resampling (Schlegel et al. 2019) in practice, which resamples experience with probability proportional to ρ^{i_n} to construct a resampled dataset $\mathcal{D}_{\rho^{i_n}}$, stabilizing the algorithm training effectively. Thus, Eq.(14) is replaced with

$$\min_{\phi^{i_n}} \mathbb{E}_{(s, a, s', r) \sim \mathcal{D}_{\rho^{i_n}}} \left[(Q_{\phi^{i_n}}(s, a^{i_n}) - y)^2 \right]. \quad (15)$$

Policy Improvement. After obtaining the optimal local value functions, we can adopt the in-sample sequential policy optimization method in Eq.(13) to learn the local policy for each agent:

$$\begin{aligned} \theta_{\text{new}}^{i_n} = \arg \min_{\theta^{i_n}} & \mathbb{E}_{(s, a^{i_n}) \sim \mathcal{D}_{\rho^{i_n}}} \left[\right. \\ & \left. - \exp \left(\frac{A_{\phi^{i_n}}(s, a^{i_n}) - \beta \log \mu^{i_n}(a^{i_n}|s)}{\alpha + \beta} \right) \log \pi_{\theta^{i_n}}(a^{i_n}|s) \right], \end{aligned}$$

where $A_{\phi^{i_n}}(s, a^{i_n}) \triangleq Q_{\phi^{i_n}}(s, a^{i_n}) - \mathbb{E}_{\pi_{\theta_{\text{old}}^{i_n}}} [Q_{\phi^{i_n}}(s, a^{i_n})]$.

Experiments

We conduct a series of experiments to evaluate InSPO on XOR game, M-NE game, Bridge (Fu et al. 2022) and StarCraft II Micromanagement (Xu et al. 2023b). In addition to Behavior Cloning (BC), our baselines also include the current state-of-the-art offline MARL algorithms: OMAR (Pan et al. 2022), CFCQL (Shao et al. 2023), OMIGA (Wang et al. 2023b) and AlberDICE (Matsunaga et al. 2023). Each algorithm is run for five random seeds, and we report the mean performance with standard deviation. For the final results, we indicate the algorithm with the best mean performance in bold, and an asterisk (*) denotes that the metric is not significantly different from the top-performing metric in that case, based on a heteroscedastic two-sided t-test with a 5% significance level. See Appendix C for experimental details.

Dataset	BC	OMAR	AlberDICE	CFCQL	OMIGA	InSPO
(a)	0.00 ± 0.01	1.00 ± 0.00	1.00 ± 0.00	-0.64 ± 0.71	0.00 ± 0.01	1.00 ± 0.00
(b)	0.23 ± 0.01	-2.00 ± 0.00	1.00 ± 0.00	-0.38 ± 0.11	0.21 ± 0.03	1.00 ± 0.00
(c)	0.00 ± 0.01	0.00 ± 0.00	1.00 ± 0.00	-0.73 ± 0.48	0.05 ± 0.00	1.00 ± 0.00

Table 1: Averaged test return on XOR game.

Dataset	BC	OMAR	AlberDICE	CFCQL	OMIGA	InSPO
balanced	-9.79 ± 0.41	20.00 ± 0.00	20.00 ± 0.00	20.00 ± 0.00	20.00 ± 0.00	20.00 ± 0.00
imbalanced	-3.47 ± 0.17	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	5.00 ± 0.00	20.00 ± 0.00

Table 2: Averaged test return on M-NE game.

Dataset	BC	OMAR	AlberDICE	CFCQL	OMIGA	InSPO
Optimal	-1.26	-2.21 ± 0.90	-6.01 ± 0.00	$-1.27 \pm 0.03^*$	-12.75 ± 1.60	-9.87 ± 0.68
Mixed	-4.56	-5.88 ± 0.49	-6.01 ± 0.00	-1.29 ± 0.00	-13.79 ± 1.78	-13.45 ± 0.42

Table 3: Averaged test return on Bridge.

Map	Dataset	BC	OMAR	CFCQL	OMIGA	InSPO
2s3z	medium	0.16 ± 0.07	0.15 ± 0.04	0.40 ± 0.10	0.23 ± 0.01	0.23 ± 0.06
	medium-replay	0.33 ± 0.04	0.24 ± 0.09	$0.55 \pm 0.07^*$	0.42 ± 0.02	0.58 ± 0.09
	expert	0.97 ± 0.02	0.95 ± 0.04	0.99 ± 0.01	0.98 ± 0.02	0.99 ± 0.01
	mixed	0.44 ± 0.06	0.60 ± 0.04	$0.84 \pm 0.09^*$	0.62 ± 0.03	0.85 ± 0.04
3s_vs_5z	medium	0.08 ± 0.02	0.00 ± 0.00	0.28 ± 0.03	0.02 ± 0.02	0.17 ± 0.05
	medium-replay	0.01 ± 0.01	0.00 ± 0.00	0.12 ± 0.04	0.02 ± 0.01	$0.10 \pm 0.05^*$
	expert	0.98 ± 0.02	0.64 ± 0.08	0.99 ± 0.01	0.98 ± 0.02	0.99 ± 0.01
	mixed	0.21 ± 0.04	0.00 ± 0.00	0.60 ± 0.14	0.20 ± 0.06	0.78 ± 0.09
5m_vs_6m	medium	$0.28 \pm 0.37^*$	0.19 ± 0.06	0.29 ± 0.05	0.25 ± 0.08	$0.28 \pm 0.06^*$
	medium-replay	0.18 ± 0.06	0.03 ± 0.02	$0.22 \pm 0.06^*$	0.16 ± 0.05	0.24 ± 0.07
	expert	0.82 ± 0.04	0.33 ± 0.06	0.84 ± 0.03	0.74 ± 0.05	0.79 ± 0.12
	mixed	0.21 ± 0.12	0.10 ± 0.10	$0.76 \pm 0.07^*$	0.38 ± 0.23	0.78 ± 0.06
6h_vs_8z	medium	0.40 ± 0.03	0.04 ± 0.03	$0.41 \pm 0.04^*$	0.34 ± 0.01	0.43 ± 0.06
	medium-replay	0.11 ± 0.04	0.00 ± 0.00	0.21 ± 0.05	0.11 ± 0.04	0.23 ± 0.02
	expert	0.60 ± 0.04	0.01 ± 0.01	$0.70 \pm 0.06^*$	0.54 ± 0.04	0.74 ± 0.11
	mixed	0.27 ± 0.06	0.00 ± 0.00	0.49 ± 0.08	0.36 ± 0.06	0.60 ± 0.12
average performance		0.38	0.21	0.54	0.39	0.55

Table 4: Averaged test winning rate on StarCraft II Micromanagement.

Comparative Evaluation

Matrix Game. We evaluate whether InSPO can address the two issues highlighted in previous section using the XOR game and M-NE game shown in Figure 1(a) and Figure 2(a). First, we evaluate the ability of InSPO to handle

the OOD joint actions issue using the XOR game. Table 1 compares the performance of all algorithms on four datasets, each comprising an equal mix of different joint actions: (a) $\{(A, B), (B, A)\}$, (b) $\{(A, A), (A, B), (B, A)\}$, and (c) $\{(A, A), (A, B), (B, A), (B, B)\}$. Figure 3 illustrates the converged joint policy of OMAR, OMIGA, and InSPO on dataset (c), representing decentralized training, value decomposition, and sequential-update methods, respectively.

As observed, only InSPO and AlberDICE, two sequential-update MAPG methods, successfully converge to the optimal policy, while the other algorithms fail, even opting for OOD joint actions. Specifically, when a more relaxed behavior policy constraint is applied, the value decomposition method (i.e., OMIGA (relaxed) in Figure 3) converges to an OOD joint action (B, B) , consistent with our analysis in previous section. These results suggest that decentralized training and value decomposition methods have limitations in environments that demand high levels of coordination.

Next, we conduct InSPO on M-NE game to evaluate its ability to alleviate the local optimum convergence issue. Ta-

Agent 1		Agent 1		Agent 1		Agent 1		
A B		A B		A B		A B		
Agent 2 A	0.0	0.0	0.0	0.0	0.36	0.20	0.0	1.0
	0.0	1.0	0.0	1.0	0.29	0.15	0.0	0.0
OMAR	0.0	0.0	0.0	0.0	0.36	0.20	0.0	1.0
	0.0	1.0	0.0	1.0	0.29	0.15	0.0	0.0
OMIGA (relaxed)		OMIGA		InSPO				

Figure 3: Final joint policy on XOR game for dataset (b).

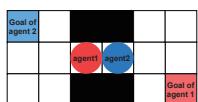


Figure 4: Bridge at the beginning.

ble 2 shows the results on datasets: (a) a balanced dataset collected by a uniform policy $\mu^i(A) = \mu^i(B) = \mu^i(C) = 1/3$, and (b) a imbalanced dataset collected by a near local optimum $\mu^i(A) = 0.8, \mu^i(B) = \mu^i(C) = 0.1$, for $i = 1, 2$. The results show that on the balanced dataset (a), most algorithms find the global optimal NE, while on the imbalanced dataset (b), only InSPO correctly identifies the global optimal NE. This indicates that in environments with multiple local optima, the convergence of most algorithms can be heavily influenced by the dataset distribution. Specifically, when the dataset is biased toward a local optimum, algorithms are prone to converging on this sub-optimal solution. In contrast, InSPO converges to the global optimal solution through comprehensive exploration of the dataset. The results demonstrate in offline scenarios, algorithms must make full use of all available dataset information to prevent being heavily influenced by behavior policies.

Bridge. Bridge, illustrated in Figure 4, is a grid-world Markov game resembling a temporal version of the XOR game. Two agents must alternately cross a one-person bridge as fast as possible. Starting side by side on the bridge, they must move together to allow one agent to cross first. For this experiment, we use two datasets provided by Matsunaga et al. (2023): optimal and mixed. The optimal dataset contains 500 trajectories, generated by combining two optimal deterministic policies: either agent 1 steps back to let agent 2 cross first, or vice versa. The mixed dataset includes the optimal dataset plus an additional 500 trajectories generated using a uniform random policy.

The performance is shown in Table 3, where the results of BC, OMAR and AlberDICE are from the report in Matsunaga et al. (2023). The performance is similar to that of the XOR game: only InSPO and AlberDICE, both using sequential-update, achieve near-optimal performance on both datasets. In contrast, both value decomposition methods fail to converge and produced undesirable outcomes.

StarCraft II. We further extend our study to the StarCraft II micromanagement benchmarks, a high-dimensional and complex environment widely used in both online and offline MARL. In this environment, we consider four representative maps: 2 easy maps (2s3z, 3s_vs_5z), 1 hard map (5m_vs_6m) and 1 super-hard map (6h_vs_8z). We use four datasets provided by Shao et al. (2023): medium, expert, medium-replay and mixed. The medium-replay dataset is a replay buffer collected during training until the policy achieves medium performance, while the mixed dataset is the equal mixture of the medium and expert datasets. The results are shown in Table 4, with the performance of CFCQL, OMAR, and BC taken from Shao et al.’s report.

In contrast to the previous benchmarks, StarCraft II does not exhibit a highly multi-modal reward landscape. Additionally, the agents share nearly identical local objectives, making this environment suitable for the IGM principle. Therefore, value decomposition methods have achieved state-of-the-art performance in this environment both in offline and online settings. Even so, as shown in Table 4, InSPO still demonstrates competitive performance and achieves state-of-the-art results in most tasks.

			Agent 1						Agent 1			
			A	B	C	A	B		A	B		
Agent 2			A	1.0	0.0	0.0	A	0.44	0.23	A	0.44	0.23
			B	0.0	0.0	0.0	B	0.21	0.12	B	0.21	0.12
			C	0.0	0.0	0.0						

(a) InSPO w/o entropy

(b) InSPO w/o SPO

Figure 5: Ablation on entropy and sequential update scheme. (a) is InSPO without entropy on M-NE game for the imbalanced dataset. (b) is simultaneous-update version of InSPO on XOR game for dataset (b).

dataset	$\alpha = 0.1$	$\alpha = 5$	$\alpha = 10$	$\alpha = 50$	auto- α
expert	0.51 ± 0.17	0.54 ± 0.06	0.69 ± 0.05	0.62 ± 0.03	0.74
mixed	0.17 ± 0.07	0.54 ± 0.06	0.53 ± 0.12	0.48 ± 0.07	0.60

Table 5: Ablation results for α on 6h_vs_8z.

Ablation Study. Here we present the impact of different components on performance of InSPO. Figure 5(a) shows the converged policy of InSPO without entropy in the M-NE game on the imbalanced dataset. Without the perturbations of entropy in the optimization objective, InSPO w/o entropy cannot escape the local optimum. Figure 5(b) shows the policy of InSPO using the simultaneous update scheme instead of sequential policy optimization (denoted as InSPO w/o SPO) on dataset (b) of the XOR game. Due to conflicting update directions, InSPO w/o SPO fails to learn the optimal policy and faces the OOD joint actions issue.

Temperature α is used to control the degree of conservatism. A too large α will result in an overly conservative policy, while a too small one will easily causes distribution shift. Thus, to obtain a suitable α , we implement both fixed and auto-tuned α in practice (see Appendix B for details), where the auto-tuned α is adjusted by $\min_{\alpha} \mathbb{E}_{\mathcal{D}}[\alpha D_{\text{KL}}(\pi, \mu) - \alpha \bar{D}_{\text{KL}}]$, where \bar{D}_{KL} is the target value. Table 5 gives ablation results for α , which shows that the auto-tuned α can find an appropriate α to further improve performance.

Furthermore, we explore the impact of update order on performance and the training efficiency of sequential updates. These results are provided in Appendix C.

Conclusion

In this paper, we study the offline MARL problem, a topic of significant practical importance and challenges that has not received adequate attention. We begin with two simple yet highly illustrative matrix games, highlighting some limitations of current offline MARL algorithms in addressing OOD joint actions and sub-optimal convergence issues. To overcome these challenges, we propose a novel algorithm called InSPO, which utilizes sequential-update insample learning to avoid OOD joint actions, and introduces policy entropy to ensure comprehensive exploration of the dataset, thus avoiding the influence of local optimum behavior policies. Furthermore, we theoretically demonstrate that InSPO possesses monotonic improvement and QRE convergence properties, and then empirically validate its superior

performance on various MARL benchmarks. For future research, integrating sequential-update in-sample learning and enhanced dataset utilization with other offline MARL algorithms presents an intriguing direction.

Acknowledgments

We gratefully acknowledge the support from the National Natural Science Foundation of China (No. 62076259, 62402252), the Fundamental and Applicational Research Funds of Guangdong Province (No. 2023A1515012946), the Fundamental Research Funds for the Central Universities Sun Yat-sen University, and the Pengcheng Laboratory Project (PCL2023A08, PCL2024Y02). This research is also supported by Meituan.

References

- Barde, P.; Foerster, J.; Nowrouzezahrai, D.; and Zhang, A. 2024. A Model-Based Solution to the Offline Multi-Agent Reinforcement Learning Coordination Problem. In *International Conference on Autonomous Agents and Multiagent Systems*. International Foundation for Autonomous Agents and Multiagent Systems / ACM.
- Ding, Z.; Su, K.; Hong, W.; Zhu, L.; Huang, T.; and Lu, Z. 2022. Multi-Agent Sequential Decision-Making via Communication. arXiv:2209.12713.
- Figueiredo Prudencio, R.; Maximo, M. R. O. A.; and Colombini, E. L. 2024. A Survey on Offline Reinforcement Learning: Taxonomy, Review, and Open Problems. *IEEE Transactions on Neural Networks and Learning Systems*, 35(8): 10237–10257.
- Fu, W.; Yu, C.; Xu, Z.; Yang, J.; and Wu, Y. 2022. Revisiting Some Common Practices in Cooperative Multi-Agent Reinforcement Learning. In *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 6863–6877. PMLR.
- Haarnoja, T.; Zhou, A.; Abbeel, P.; and Levine, S. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, 1861–1870. PMLR.
- Haarnoja, T.; Zhou, A.; Hartikainen, K.; Tucker, G.; Ha, S.; Tan, J.; Kumar, V.; Zhu, H.; Gupta, A.; Abbeel, P.; and Levine, S. 2019. Soft Actor-Critic Algorithms and Applications. arXiv:1812.05905.
- Kostrikov, I.; Nair, A.; and Levine, S. 2022. Offline Reinforcement Learning with Implicit Q-Learning. In *International Conference on Learning Representations*. OpenReview.net.
- Kuba, J. G.; Chen, R.; Wen, M.; Wen, Y.; Sun, F.; Wang, J.; and Yang, Y. 2022. Trust Region Policy Optimisation in Multi-Agent Reinforcement Learning. In *The Tenth International Conference on Learning Representations*. OpenReview.net.
- Kumar, A.; Zhou, A.; Tucker, G.; and Levine, S. 2020. Conservative Q-Learning for Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems*.
- Li, Z.; Pan, L.; and Huang, L. 2023. Beyond Conservatism: Diffusion Policies in Offline Multi-agent Reinforcement Learning. arXiv:2307.01472.
- Li, Z.; Zhao, W.; Wu, L.; and Pajarinen, J. 2024. Backpropagation Through Agents. In *Annual AAAI Conference on Artificial Intelligence*. AAAI Press.
- Littman, M. L. 1994. Markov Games as a Framework for Multi-Agent Reinforcement Learning. In *Machine Learning, Proceedings of the Eleventh International Conference*, 157–163. Morgan Kaufmann.
- Liu, J.; Zhong, Y.; Hu, S.; Fu, H.; Fu, Q.; Chang, X.; and Yang, Y. 2024. Maximum Entropy Heterogeneous-Agent Reinforcement Learning. In *International Conference on Learning Representations*. OpenReview.net.
- Matsunaga, D. E.; Lee, J.; Yoon, J.; Leonards, S.; Abbeel, P.; and Kim, K. 2023. AlberDICE: Addressing Out-Of-Distribution Joint Actions in Offline Multi-Agent RL via Alternating Stationary Distribution Correction Estimation. In *Advances in Neural Information Processing Systems*.
- McKelvey, R. D.; and Palfrey, T. R. 1995. Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1): 6–38.
- Pan, L.; Huang, L.; Ma, T.; and Xu, H. 2022. Plan Better Amid Conservatism: Offline Multi-Agent Reinforcement Learning with Actor Rectification. In *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 17221–17237. PMLR.
- Rashid, T.; Samvelyan, M.; de Witt, C. S.; Farquhar, G.; Foerster, J. N.; and Whiteson, S. 2018. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, 4292–4301. PMLR.
- Schlegel, M.; Chung, W.; Graves, D.; Qian, J.; and White, M. 2019. Importance Resampling for Off-policy Prediction. In *Advances in Neural Information Processing Systems*, 1797–1807.
- Shao, J.; Qu, Y.; Chen, C.; Zhang, H.; and Ji, X. 2023. Counterfactual Conservative Q Learning for Offline Multi-agent Reinforcement Learning. In *Advances in Neural Information Processing Systems*.
- Son, K.; Kim, D.; Kang, W. J.; Hostallero, D.; and Yi, Y. 2019. QTRAN: Learning to Factorize with Transformation for Cooperative Multi-Agent Reinforcement Learning. In *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, 5887–5896. PMLR.
- Sutton, R. S. 2018. Reinforcement Learning: An Introduction.
- Wang, J.; Ren, Z.; Liu, T.; Yu, Y.; and Zhang, C. 2021. QPLEX: Duplex Dueling Multi-Agent Q-Learning. In *International Conference on Learning Representations*. OpenReview.net.
- Wang, J.; Ye, D.; and Lu, Z. 2023. More Centralized Training, Still Decentralized Execution: Multi-Agent Conditional Policy Factorization. In *International Conference on Learning Representations*. OpenReview.net.

- Wang, X.; Tian, Z.; Wan, Z.; Wen, Y.; Wang, J.; and Zhang, W. 2023a. Order Matters: Agent-by-agent Policy Optimization. In *International Conference on Learning Representations*. OpenReview.net.
- Wang, X.; Xu, H.; Zheng, Y.; and Zhan, X. 2023b. Offline Multi-Agent Reinforcement Learning with Implicit Global-to-Local Value Regularization. In *Advances in Neural Information Processing Systems*.
- Wu, Z.; Yu, C.; Ye, D.; Zhang, J.; Piao, H.; and Zhuo, H. H. 2021. Coordinated Proximal Policy Optimization. In *Advances in Neural Information Processing Systems*, 26437–26448.
- Xu, H.; Jiang, L.; Li, J.; Yang, Z.; Wang, Z.; Chan, W. K. V.; and Zhan, X. 2023a. Offline RL with No OOD Actions: In-Sample Learning via Implicit Value Regularization. In *International Conference on Learning Representations*. OpenReview.net.
- Xu, H.; Jiang, L.; Li, J.; Yang, Z.; Wang, Z.; Chan, W. K. V.; and Zhan, X. 2023b. Offline RL with No OOD Actions: In-Sample Learning via Implicit Value Regularization. In *International Conference on Learning Representations*. OpenReview.net.
- Yang, Y.; Ma, X.; Li, C.; Zheng, Z.; Zhang, Q.; Huang, G.; Yang, J.; and Zhao, Q. 2021. Believe What You See: Implicit Constraint Approach for Offline Multi-Agent Reinforcement Learning. In *Advances in Neural Information Processing Systems*.
- Yu, C.; Velu, A.; Vinitsky, E.; Gao, J.; Wang, Y.; Bayen, A. M.; and Wu, Y. 2022. The Surprising Effectiveness of PPO in Cooperative Multi-Agent Games. In *Advances in Neural Information Processing Systems*.
- Zhang, K.; Yang, Z.; and Başar, T. 2021. Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms. arXiv:1911.10635.
- Zhu, Z.; Liu, M.; Mao, L.; Kang, B.; Xu, M.; Yu, Y.; Ermon, S.; and Zhang, W. 2023. Madiff: Offline multi-agent learning with diffusion models. arXiv:2305.17330.

A Proofs

A.1 Proof of Policy Evaluation

Lemma 6. Given a policy π , consider the modified policy evaluation operator \mathcal{T}_π under MEBR-MGs and a initial Q-function $\mathbf{Q}_0 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, and define $\mathbf{Q}_{k+1} = \mathcal{T}_\pi \mathbf{Q}_k$. Then the sequence \mathbf{Q}_k will converge to the Q-function \mathbf{Q}_π of policy π as $k \rightarrow \infty$.

Proof. With a pseudo-reward $r_\pi(s, \mathbf{a}) \triangleq r(s, \mathbf{a}) - \gamma \mathbb{E}_{s' \mid s, \mathbf{a}} [\alpha D_{\text{KL}}(\pi(\cdot|s'), \mu(\cdot|s')) - \beta \mathcal{H}(\pi(\cdot|s'))]$, the update rule of Q-function can be represented as:

$$\mathbf{Q}(s, \mathbf{a}) \leftarrow r_\pi(s, \mathbf{a}) + \gamma \mathbb{E}_{s' \mid s, \mathbf{a}, \mathbf{a}' \sim \pi} [\mathbf{Q}(s', \mathbf{a}')].$$

Then, we can apply the standard convergence results for policy evaluation (Sutton 2018). \square

A.2 Proof of QRE

Proposition 7. In a MEBR-MG, a joint policy π_* is a QRE if it holds

$$V_{\pi_*}(s) \geq V_{\pi^i, \pi_*^{-i}}(s), \quad \forall i \in \mathcal{N}, \pi^i, s \in \mathcal{S}. \quad (16)$$

Then the QRE policies for each agent i are given by

$$\begin{aligned} \pi_*^i(a^i | s) &\propto \mu^i(a^i | s) \\ &\cdot \exp\left(\frac{\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_*^{-i}}[\mathbf{Q}_{\pi_*}(s, a^i, \mathbf{a}^{-i})] - \beta \log \mu^i(a^i | s)}{\alpha + \beta}\right) \end{aligned} \quad (17)$$

Proof. We consider the following constrained policy optimization problem to agent i :

$$\begin{aligned} \max_{\pi^i} & \mathbb{E}_{a^i \sim \pi^i, \mathbf{a}^{-i} \sim \pi^{-i}} [\mathbf{Q}_\pi(s, \mathbf{a})] \\ & - \alpha \sum_{j=1}^N \sum_{a^j} \pi^j(a^j | s) \log \frac{\pi^j(a^j | s)}{\mu^j(a^j | s)} \\ & - \beta \sum_{j=1}^N \sum_{a^j} \pi^j(a^j | s) \log \pi^j(a^j | s), \\ \text{s.t. } & \sum_{a^i} \pi^i(a^i | s) = 1, \quad \forall s \in \mathcal{S}. \end{aligned}$$

Its associated Lagrangian function is

$$\begin{aligned} \mathcal{L}(\pi^i, \lambda) = & \mathbb{E}_{a^i \sim \pi^i, \mathbf{a}^{-i} \sim \pi^{-i}} [\mathbf{Q}_\pi(s, \mathbf{a})] \\ & - \alpha \sum_{j=1}^N \sum_{a^j} \pi^j(a^j | s) \log \frac{\pi^j(a^j | s)}{\mu^j(a^j | s)} \\ & - \beta \sum_{j=1}^N \sum_{a^j} \pi^j(a^j | s) \log \pi^j(a^j | s) \\ & + \lambda \left(\sum_{a^i} \pi^i(a^i | s) - 1 \right). \end{aligned}$$

Therefore, we have:

$$\begin{aligned} \frac{\partial \mathcal{L}(\pi^i, \lambda)}{\partial \pi^i(a^i | s)} = & \mathbb{E}_{\mathbf{a}^{-i} \sim \pi^{-i}} [\mathbf{Q}_\pi(s, a^i, \mathbf{a}^{-i})] \\ & - \alpha \log \frac{\pi^i(a^i | s)}{\mu^i(a^i | s)} - \beta \log \pi^i(a^i | s) - \alpha - \beta + \lambda. \end{aligned}$$

According to the Karush-Kuhn-Tucker (KKT) conditions, we know the optimal policy (i.e., QRE policy) π_*^i satisfies

$$\begin{aligned} \pi_*^i(a^i | s) &= \exp\left(\frac{\lambda_*}{\alpha + \beta} - 1\right) \cdot \mu^i(a^i | s) \cdot \exp\left(\right. \\ &\left. \frac{\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_*^{-i}}[\mathbf{Q}_{\pi_*}(s, a^i, \mathbf{a}^{-i})] - \beta \log \mu^i(a^i | s)}{\alpha + \beta}\right), \end{aligned}$$

where λ_* is used to ensure $\sum_{a^i} \pi_*^i(a^i | s) = 1$, i.e.,

$$\begin{aligned} \exp\left(1 - \frac{\lambda_*}{\alpha + \beta}\right) &= \sum_{a^i} \left[\mu^i(a^i | s) \right. \\ &\left. \cdot \exp\left(\frac{\mathbb{E}_{\mathbf{a}^{-i} \sim \pi_*^{-i}}[\mathbf{Q}_{\pi_*}(s, a^i, \mathbf{a}^{-i})] - \beta \log \mu^i(a^i | s)}{\alpha + \beta}\right) \right]. \end{aligned}$$

Thus, the proof is completed. \square

A.3 Proof of Policy Improvement

Proposition 8. *The sequential policy optimization procedure under MEBR-MG guarantees policy improvement, i.e., $\forall s \in \mathcal{S}, a \in \mathcal{A}$,*

$$Q_{\pi_{\text{new}}}(s, a) \geq Q_{\pi_{\text{old}}}(s, a), V_{\pi_{\text{new}}}(s) \geq V_{\pi_{\text{old}}}(s).$$

Proof. The policy improvement step of sequential policy optimization is given by

$$\begin{aligned} \pi_{\text{new}}^{i_n} &= \arg \max_{\pi^{i_n}} \mathbb{E}_{a^{i_n} \sim \pi^{i_n}} \left[Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n}) \right. \\ &\quad \left. - \alpha \log \frac{\pi_{\text{new}}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \pi_{\text{new}}^{i_n}(a^{i_n}|s) \right], \end{aligned} \quad (18)$$

where

$$Q_{\pi_{\text{old}}}^{i_{1:n}}(s, a^{i_n}) \triangleq \mathbb{E}_{\pi_{\text{new}}^{i_{1:n-1}}, \pi_{\text{old}}^{-i_{1:n}}} \left[Q_{\pi_{\text{old}}}(s, a^{-i_n}, a^{i_n}) \right].$$

We first show that the resulting joint policy π_{new} satisfies:

$$\begin{aligned} \pi_{\text{new}} &= \arg \max_{\pi} \mathbb{E}_{\mathbf{a} \sim \pi} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) \right. \\ &\quad \left. - \alpha \log \frac{\pi(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi(\mathbf{a}|s) \right] \\ &= \arg \max_{\pi} L_{\pi_{\text{old}}}(\pi). \end{aligned}$$

Otherwise, suppose that there exists a policy $\bar{\pi} \neq \pi_{\text{new}}$ such that $L_{\pi_{\text{old}}}(\bar{\pi}) > L_{\pi_{\text{old}}}(\pi_{\text{new}})$. From the policy improvement step, we have

$$\begin{aligned} &\mathbb{E}_{a^{i_n} \sim \pi_{\text{new}}^{i_n}} \left[\mathbb{E}_{\pi_{\text{new}}^{i_{1:n-1}}, \pi_{\text{old}}^{-i_{1:n}}} \left[Q_{\pi_{\text{old}}}(s, a^{-i_n}, a^{i_n}) \right] \right. \\ &\quad \left. - \alpha \log \frac{\pi_{\text{new}}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \pi_{\text{new}}^{i_n}(a^{i_n}|s) \right] \\ &\geq \\ &\mathbb{E}_{a^{i_n} \sim \bar{\pi}^{i_n}} \left[\mathbb{E}_{\pi_{\text{new}}^{i_{1:n-1}}, \pi_{\text{old}}^{-i_{1:n}}} \left[Q_{\pi_{\text{old}}}(s, a^{-i_n}, a^{i_n}) \right] \right. \\ &\quad \left. - \alpha \log \frac{\bar{\pi}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \bar{\pi}^{i_n}(a^{i_n}|s) \right] \end{aligned}$$

Subtracting both sides of the inequality by $\mathbb{E}_{\pi_{\text{new}}^{i_{1:n-1}}, \pi_{\text{old}}^{-i_{1:n-1}}} \left[Q_{\pi_{\text{old}}}(s, a^{-i_n}, a^{i_n}) \right]$ gives

$$\begin{aligned} &\mathbb{E}_{a^{i_n} \sim \pi_{\text{new}}^{i_n}} \left[\mathbb{E}_{\pi_{\text{new}}^{i_{1:n-1}}} \left[A_{\pi_{\text{old}}}^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n}) \right] \right. \\ &\quad \left. - \alpha \log \frac{\pi_{\text{new}}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \pi_{\text{new}}^{i_n}(a^{i_n}|s) \right] \\ &\geq \\ &\mathbb{E}_{a^{i_n} \sim \bar{\pi}^{i_n}} \left[\mathbb{E}_{\pi_{\text{new}}^{i_{1:n-1}}} \left[A_{\pi_{\text{old}}}^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n}) \right] \right. \\ &\quad \left. - \alpha \log \frac{\bar{\pi}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \bar{\pi}^{i_n}(a^{i_n}|s) \right], \end{aligned} \quad (19)$$

where

$$\begin{aligned} A_{\pi_{\text{old}}}^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n}) &\triangleq \mathbb{E}_{\pi_{\text{old}}^{-i_{1:n}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}^{-i_n}, a^{i_n}) \right] \\ &\quad - \mathbb{E}_{\pi_{\text{old}}^{-i_{1:n-1}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}^{-i_n}, a^{i_n}) \right]. \end{aligned}$$

Combining this inequality (19) with Lemma 9, we have:

$$\begin{aligned} &\mathbb{E}_{\mathbf{a} \sim \pi_{\text{new}}} \left[A_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{new}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{new}}(\mathbf{a}|s) \right] \\ &= \sum_{n=1}^N \mathbb{E}_{\mathbf{a}^{i_{1:n-1}} \sim \pi_{\text{new}}^{i_{1:n-1}}, a^{i_n} \sim \pi_{\text{new}}^{i_n}} \left[A_{\pi_{\text{old}}}^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n}) \right. \\ &\quad \left. - \alpha \log \frac{\pi_{\text{new}}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \pi_{\text{new}}^{i_n}(a^{i_n}|s) \right] \\ &\geq \sum_{n=1}^N \mathbb{E}_{\mathbf{a}^{i_{1:n-1}} \sim \pi_{\text{new}}^{i_{1:n-1}}, a^{i_n} \sim \bar{\pi}^{i_n}} \left[A_{\pi_{\text{old}}}^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n}) \right. \\ &\quad \left. - \alpha \log \frac{\bar{\pi}^{i_n}(a^{i_n}|s)}{\mu^{i_n}(a^{i_n}|s)} - \beta \log \bar{\pi}^{i_n}(a^{i_n}|s) \right] \\ &= \mathbb{E}_{\mathbf{a} \sim \bar{\pi}} \left[A_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\bar{\pi}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \bar{\pi}(\mathbf{a}|s) \right]. \end{aligned}$$

The resulting inequality can be equivalently rewritten as

$$\begin{aligned} &\mathbb{E}_{\mathbf{a} \sim \pi_{\text{new}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{new}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{new}}(\mathbf{a}|s) \right] \\ &\geq \mathbb{E}_{\mathbf{a} \sim \bar{\pi}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\bar{\pi}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \bar{\pi}(\mathbf{a}|s) \right], \end{aligned}$$

which contradicts the claim $L_{\pi_{\text{old}}}(\bar{\pi}) > L_{\pi_{\text{old}}}(\pi_{\text{new}})$. Hence, we have $\pi_{\text{new}} = \arg \max_{\pi} L_{\pi_{\text{old}}}(\pi)$, which gives

$$\begin{aligned} &\mathbb{E}_{\mathbf{a} \sim \pi_{\text{new}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{new}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{new}}(\mathbf{a}|s) \right] \\ &\geq \mathbb{E}_{\mathbf{a} \sim \pi_{\text{old}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{old}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{old}}(\mathbf{a}|s) \right] \\ &= V_{\pi_{\text{old}}}(s) \end{aligned}$$

Therefore, we have

$$\begin{aligned} Q_{\pi_{\text{old}}}(s, \mathbf{a}) &= r(s, \mathbf{a}) + \gamma \mathbb{E}_{s' | s, \mathbf{a}} [V_{\pi_{\text{old}}}(s')] \\ &\leq r(s, \mathbf{a}) + \gamma \mathbb{E}_{s' | s, \mathbf{a}} \left[\mathbb{E}_{\mathbf{a} \sim \pi_{\text{new}}} \left[Q_{\pi_{\text{old}}}(s', \mathbf{a}) - \alpha \log \frac{\pi_{\text{new}}(\mathbf{a}|s')}{\mu(\mathbf{a}|s')} - \beta \log \pi_{\text{new}}(\mathbf{a}|s') \right] \right] \\ &\quad \vdots \\ &\leq Q_{\pi_{\text{new}}}(s, \mathbf{a}), \quad \forall s, \mathbf{a} \end{aligned}$$

And then,

$$\begin{aligned} V_{\pi_{\text{old}}}(s) &= \mathbb{E}_{\mathbf{a} \sim \pi_{\text{old}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{old}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{old}}(\mathbf{a}|s) \right] \\ &\leq \mathbb{E}_{\mathbf{a} \sim \pi_{\text{new}}} \left[Q_{\pi_{\text{old}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{new}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{new}}(\mathbf{a}|s) \right] \\ &\leq \mathbb{E}_{\mathbf{a} \sim \pi_{\text{new}}} \left[Q_{\pi_{\text{new}}}(s, \mathbf{a}) - \alpha \log \frac{\pi_{\text{new}}(\mathbf{a}|s)}{\mu(\mathbf{a}|s)} - \beta \log \pi_{\text{new}}(\mathbf{a}|s) \right] \\ &= V_{\pi_{\text{new}}}(s), \quad \forall s. \end{aligned}$$

Thus, the proof is completed. \square

Lemma 9 (Multi-Agent Advantage Decomposition from Kuba et al.). *For any state s and joint action, the joint advantage function can be decomposed as:*

$$\mathbf{A}_\pi(s, \mathbf{a}) = \sum_{n=1}^N \mathbf{A}_\pi^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n})$$

Proof. By the definition of multi-agent advantage function, we have:

$$\begin{aligned} \mathbf{A}_\pi(s, \mathbf{a}) &= \mathbf{Q}_\pi(s, \mathbf{a}) - \mathbf{V}_\pi(s) \\ &= \sum_{n=1}^N \left[\mathbb{E}_{\pi^{-i_{1:n}}} \left[\mathbf{Q}_\pi(s, \mathbf{a}^{-i_n}, a^{i_n}) \right] \right. \\ &\quad \left. - \mathbb{E}_{\pi^{-i_{1:n-1}}} \left[\mathbf{Q}_\pi(s, \mathbf{a}^{-i_n}, a^{i_n}) \right] \right] \\ &= \sum_{n=1}^N \mathbf{A}_\pi^{i_n}(s, \mathbf{a}^{i_{1:n-1}}, a^{i_n}) \end{aligned}$$

□

A.4 Proof of QRE convergence

Theorem 10. *In the tabular setting, the joint policy π updated by InSPO converges to QRE.*

Proof. First, we have that $\mathbf{Q}_{\pi_{k+1}}(s, \mathbf{a}) \geq \mathbf{Q}_{\pi_k}(s, \mathbf{a})$ by Proposition 8 and that the Q-function is upper-bounded since reward and regularizations are bounded. Hence, the sequence of policies converges to some limit point $\bar{\pi}$.

Then, considering this limit point joint policy $\bar{\pi}$, it must be the case that $\forall i, \pi^i$,

$$\begin{aligned} &\mathbb{E}_{a^i \sim \bar{\pi}^i} \left[\mathbb{E}_{\bar{\pi}^{-i}} \left[\mathbf{Q}_{\bar{\pi}}(s, \mathbf{a}^{-i}, a^i) \right] \right. \\ &\quad \left. - \alpha \log \frac{\bar{\pi}^i(a^i|s)}{\mu^i(a^i|s)} - \beta \log \bar{\pi}^i(a^i|s) \right] \\ &\geq \\ &\mathbb{E}_{a^i \sim \pi^i} \left[\mathbb{E}_{\bar{\pi}^{-i}} \left[\mathbf{Q}_{\bar{\pi}}(s, \mathbf{a}^{-i}, a^i) \right] \right. \\ &\quad \left. - \alpha \log \frac{\pi^i(a^i|s)}{\mu^i(a^i|s)} - \beta \log \pi^i(a^i|s) \right], \end{aligned}$$

which is equivalent with $V_{\bar{\pi}}(s) \geq V_{\pi^i, \bar{\pi}^{-i}}(s)$. Thus, $\bar{\pi}$ is a quantal response equilibrium, which finishes the proof. □

B Details of Practical Algorithm

In the practical implementation of InSPO, we train the policy network θ^{i_n} by loss function

$$\begin{aligned} J(\theta^{i_n}) &\triangleq \mathbb{E}_{(s, a^{i_n}) \sim \mathcal{D}_{\rho^{i_n}}} \left[\right. \\ &\quad \left. - \exp \left(\frac{A_{\bar{\phi}^{i_n}}(s, a^{i_n}) - \beta \log \mu^{i_n}(a^{i_n}|s)}{\alpha + \beta} \right) \log \pi_{\theta^{i_n}}(a^{i_n}|s) \right], \end{aligned} \quad (20)$$

where

$$A_{\bar{\phi}^{i_n}}(s, a^{i_n}) \triangleq Q_{\bar{\phi}^{i_n}}(s, a^{i_n}) - \mathbb{E}_{\pi_{\theta^{i_n}}} [Q_{\bar{\phi}^{i_n}}(s, a^{i_n})],$$

and $\bar{\phi}^{i_n}$ is the soft-target network of ϕ^{i_n} .

The behavior policy μ^{i_n} in $J(\theta^{i_n})$ is pre-trained by the Behavior Cloning

$$\min_{\mu^{i_n}} \mathbb{E}_{(s, a^{i_n}) \sim \mathcal{D}} \left[-\log \mu^{i_n}(a^{i_n}|s) \right]. \quad (21)$$

Instead of using the μ^{i_n} trained by Eq.(21) to calculate μ^{-i_n} in the computation of importance resampling ratio ρ^{i_n} , we pre-train an MLP-based autoregressive behavior policy, same with the one in Matsunaga et al., for numerically stability. The optimization objective of the autoregressive behavior policy is given by

$$\min_{\mu^{-i}} \mathbb{E}_{(s, \mathbf{a}) \sim \mathcal{D}} \left[- \sum_{j=1, j \neq i}^N \log \mu^{-i}(a^j|s, a^i, \mathbf{a}^{1:i-1}) \right]. \quad (22)$$

Then, in the computation of ρ^{i_n} , we use the geometric mean to prevent the collapse of ρ^{i_n} due to the growth of the number of agents:

$$\rho^{i_n} = \left(\frac{\pi^{-i_n}(\mathbf{a}^{-i_n}|s)}{\mu^{-i_n}(\mathbf{a}^{-i_n}|s)} \right)^{\frac{1}{N-1}}. \quad (23)$$

In order to estimate the future return based on state-action pair, we train the local Q-function network by minimizing the temporal difference (TD) error with a CQL regularization term to further penalize OOD action values:

$$\begin{aligned} J(\phi^{i_n}) &\triangleq \mathbb{E}_{(s, \mathbf{a}, s', r) \sim \mathcal{D}_{\rho^{i_n}}} \left[\left(Q_{\phi^{i_n}}(s, a^{i_n}) - y \right)^2 \right] \\ &\quad + \alpha_{\text{CQL}} \mathbb{E}_{s \sim \mathcal{D}_{\rho^{i_n}}} \left[\log \sum_{a^{i_n}} \exp(Q_{\phi^{i_n}}(s, a^{i_n})) \right. \\ &\quad \left. - \mathbb{E}_{a^{i_n} \sim \mu^{i_n}} [Q_{\phi^{i_n}}(s, a^{i_n})] \right], \end{aligned} \quad (24)$$

where

$$\begin{aligned} y &= y(s, \mathbf{a}, s', r) \triangleq r + \gamma \mathbb{E}_{a^{i_n} \sim \pi_{\text{old}}^{i_n}} [Q_{\bar{\phi}^{i_n}}(s', a^{i_n})] \\ &\quad - \alpha D_{\text{KL}}(\pi_{\theta^{i_n}}(\cdot|s'), \mu^{i_n}(\cdot|s')) + \beta \mathcal{H}(\pi_{\theta^{i_n}}(\cdot|s')). \end{aligned}$$

Here we use $\alpha_{\text{CQL}} = 0.1$ as the default value.

Lastly, we give the loss function of auto-tuned α , which is inspired by the auto-tuned temperature extension of SAC (Haarnoja et al. 2019):

$$J(\alpha) \triangleq \mathbb{E}_{s \sim \mathcal{D}} \left[\sum_{i \in \mathcal{N}} \left(\alpha D_{\text{KL}}(\pi^i(\cdot|s), \mu^i(\cdot|s)) - \alpha \bar{D}_{\text{KL}} \right) \right], \quad (25)$$

where \bar{D}_{KL} is the target value with the default value 0.18.

C Experimental Details

C.1 Baselines

BC, OMAR and CFCQL: We use the open-source implementation ¹ provided by Shao et al.. **AlberDICE:** We use the open-source implementation ² provided by Matsunaga et al.. **OMIGA:** We use the open-source implementation ³ provided by Wang et al..

¹<https://github.com/thu-rlab/CFCQL>

²<https://github.com/dematsunaga/alberdice>

³<https://github.com/ZhengYinan-AIR/OMIGA>

Map	Dataset	random	fixed	semi-greedy
5m_vs_6m	medium	0.28 ± 0.06	0.29 ± 0.06	0.28 ± 0.03
	medium-replay	0.24 ± 0.07	0.25 ± 0.10	0.27 ± 0.04
	expert	0.79 ± 0.12	0.84 ± 0.05	0.81 ± 0.05
	mixed	0.78 ± 0.06	0.78 ± 0.05	0.78 ± 0.08
6h_vs_8z	medium	0.43 ± 0.06	0.39 ± 0.07	0.42 ± 0.11
	medium-replay	0.23 ± 0.02	0.23 ± 0.04	0.17 ± 0.08
	expert	0.74 ± 0.11	0.72 ± 0.07	0.71 ± 0.09
	mixed	0.60 ± 0.12	0.62 ± 0.10	0.65 ± 0.11

Table 6: Impact of update order.

C.2 Dataset Details

Bridge. We use the datasets ² provided Matsunaga et al.. The optimal dataset (500 trajectories) is collected by a hand-crafted (multi-modal) optimal policy, and the mixed dataset (500 trajectories) is the equal mixture of the optimal dataset and 500 trajectories collected by a uniform random policy.

StarCraft II. We use the datasets ¹ provided Shao et al., which are collected by QMIX (Rashid et al. 2018). The medium dataset (5000 trajectories) is collected by a partially trained model of QMIX, and the expert dataset (5000 trajectories) is collected by a fully trained model. The mixed dataset (5000 trajectories) is the equal mixture of medium and expert datasets. The medium-replay dataset (5000 trajectories) is the replay buffer during training until the policy reaches the medium performance.

C.3 Resources

We run all the experiments on 4*NVIDIA GeForce RTX 3090 GPUs and 4*NVIDIA A30. Each setting is repeated for 5 seeds. For one seed, it takes about 1.5 hours for StarCraft II, 1 hour for Bridge, and 15 minutes for matrix games.

C.4 Reproducibility

The local Q-function network is represented by 3-layer ReLU activated MLPs with 256 units for each hidden layer. The policy network is implemented in two ways: MLP and RNN. For the Matrix Game and Bridge, the policy network is represented by 3-layer ReLU activated MLPs with 256 units for each hidden layer. For StarCraft II, the policy network consists of two linear layers of 64 units and one GRU-Cell layer, referring to the CFCQL implementation ¹. All the networks are optimized by Adam optimizer.

For all datasets and algorithms, we run all the experiments with $1e7$ training steps on 5 seeds: 0, 1, 2, 3, 4. For evaluation we rolled out policies for 32 episodes and computed the mean episode return (or winning rate). For OMAR, we tune a best CQL weight from $\{0, 0.1, 0.5, 1, 2, 3, 4, 5, 10\}$. For CFCQL, we tune a best CQL weight from $\{0, 0.1, 0.5, 1, 5, 10, 50\}$, softmax temperature from $\{0, 0.1, 0.5, 1, 100\}$. For OMIGA, we tune a best regularization temperature from $\{0.1, 0.5, 1, 3, 5, 7, 10\}$. For InSPO, we tune a best α from $\{0.1, 0.5, 1, 3, 5\}$, \bar{D}_{KL} from $\{0.08, 0.18, 0.3\}$ and an exponentially decaying β from $\{0, 5, 10\}$.

C.5 Impact of Update Order

Update order might impact performance (Ding et al. 2022; Wang et al. 2023a). Ding et al. used a world model to determine the order by comparing the value of intentions. Wang et al. introduced a semi-greedy agent selection rule, which prioritizes updating agents with a higher expected advantage. However, they are heuristic and may not guarantee optimal results. To investigate this, we add an ablation study, comparing ‘random’, ‘fixed’, and semi-greedy update rule in terms of their impact on InSPO’s performance. Table 6 shows that the effect of update order on performance is not significant, but determining the optimal update sequence is still an interesting direction for future work.

C.6 Training Efficiency

The sequential updates can increase training times, which is a common issue with this method. Here we add a comparison of training times between InSPO and CFCQL in Table 7, which shows that sequential updates do increase the training time. However, some work (Wang et al. 2023a) aim to address this issue by grouping agents into blocks for simultaneous updates within blocks, with sequential updates between blocks.

map	CFCQL	InSPO
2s3z	1h	1.75h
3s_vs_5z	1h	1.2h
5m_vs_6m	0.5h	1.25h
6h_vs_8z	0.75h	2.25h

Table 7: Comparison of training times between InSPO and CFCQL.

D Value Decomposition in XOR

In XOR game, the minimization of TD error $\mathbb{E}_{\mathcal{D}}[(\mathbf{Q}(a^1, a^2) - r(a^1, a^2))^2]$ motivates the global Q-network to satisfy

$$\mathbf{Q}(A, B) = \mathbf{Q}(B, A) > \mathbf{Q}(A, A). \quad (26)$$

According to IGM principle, if $Q^i(A) > Q^i(B)$, then $\mathbf{Q}(A, A) > \mathbf{Q}(A, B)$, and $\mathbf{Q}(A, A) > \mathbf{Q}(B, B)$, contradicting Eq.(26). If $Q^i(A) = Q^i(B)$, then $\mathbf{Q}(A, A) = \mathbf{Q}(A, B) = \mathbf{Q}(B, B)$, contradicting Eq.(26) again. Therefore, we have $Q^i(A) < Q^i(B)$, which leads to the OOD joint action (B, B) .

E Similar Techniques to the Concept of “sequential”

The idea of “sequential” has been explored in various directions within MARL. Specifically, Ding et al. introduced SeqComm, a communication framework where agents condition their actions based on the *ordered actions* of others, mitigating circular dependencies that arise in *simultaneous* communication. MACPF (Wang, Ye, and Lu 2023) decomposes joint policies into individual ones, incorporating a correction term to model dependencies on preceding agents’ actions in *sequential execution*. BPPO (Li et al. 2024) employs an auto-regressive joint policy with a *fixed execution order*. During training, agents act *sequentially* based on the prior agents’ actions, and update their policies based on the feedback from subsequent agents. This bidirectional mechanism enables each agent adapts to the changing behavior of the team efficiently.

While these works use the concept of “sequential” to improve coordination, applying them directly in offline MARL poses challenges. Both policy and value functions in MACPF and BPPO explicitly condition on prior agents’ actions, which can be challenging in offline settings where required data may be missing, potentially hindering accurate value function updates.