
Pessimism Meets Invariance: Provably Efficient Offline Mean-Field Multi-Agent RL

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Mean-Field Multi-Agent Reinforcement Learning (MF-MARL) is attractive in the
2 applications involving a large population of homogeneous agents, as it exploits
3 the permutation invariance of agents and avoids the curse of many agents. Most
4 existing results only focus on online settings, in which agents can interact with
5 the environment during training. In some applications such as social welfare
6 optimization, however, the interaction during training can be prohibitive or even
7 unethical in the societal systems. To bridge such a gap, we propose a SAFARI
8 (pessimistic mean-field value iteration) algorithm for off-line MF-MARL, which
9 only requires a handful of pre-collected experience data. Theoretically, under a
10 weak coverage assumption that the experience dataset contains enough information
11 about the optimal policy, we prove that for an episodic mean-field MDP with
12 a horizon H and N training trajectories, SAFARI attains a sub-optimality gap
13 of $\mathcal{O}(H^2 d_{\text{eff}}/\sqrt{N})$, where d_{eff} is the effective dimension of the function class
14 for parameterizing the value function, but independent on the number of agents.
15 Moreover, we also propose a variant of our SAFARI algorithm for the online
16 setting, which is of independent interest.

17 1 Introduction

18 Significant progress has been made towards multi-agent reinforcement learning (MARL) for many
19 prominent sequential decision making problems, such as social welfare optimization [1], fleet control
20 of autonomous vehicles [2] and playing multiplayer online battle arena (MOBA) games [3]. As the
21 joint state and action space scales exponentially with the number of agents, however, MARL becomes
22 computationally expensive. One remedy is the mean-field regime when an extremely large number
23 of homogenous agents are involved, e.g., social welfare optimization. The effect of each agent on
24 the overall multi-agent system can become infinitesimal, and therefore all agents can be considered
25 interchangeable/indistinguishable [4–6]. Accordingly, the interaction among agents can be captured
26 by some mean-field quantity such as the empirical distribution of states, and therefore each agent
27 only needs to find the best response to the so-called “mean-field state”, which avoids the curse of
28 many agents.

29 Most existing results on mean-field MARL (MF-MARL) are for the online setting [4, 7], where the
30 agents can interact with the environment during training. However, such interaction during training
31 can be prohibitive for some important applications [1, 8–10]. Taking social welfare optimization as
32 an example, repeatedly conducting social experiments on human being can be unaffordable or even
33 unethical in the societal systems. Therefore, we can only consider the offline settings, i.e., we learn
34 the optimal policy based on some pre-collected experience data [10]. Unfortunately, existing offline
35 reinforcement learning (RL) algorithms and theories all focus on the single agent settings, and no
36 algorithms and theories have been developed for MARL under the offline settings, regardless of the
37 mean-field regime or not.

38 To bridge such a critical gap, we propose the first pessimistic algorithm – named SAFARI (peSsimistic
 39 meAn-Field vAlue iteRatIon) for mean-field MARL, which can provably achieve sample efficiency
 40 under the offline setting. Our proposed algorithm contains two important components: (1) To
 41 incorporate the permutation invariance of the homogenous agents, we adopt a RKHS (Reproducing
 42 Kernel Hilbert Space) mean-embedding approach for approximating value functions, which avoids the
 43 exponential blowup of the agents’ state and action spaces; (2) We develop an uncertainty quantifier,
 44 and integrate it into the value iteration procedure as the penalty function. Such a penalty function
 45 can effectively screen the “spuriously correlated trajectories”, i.e., which possibly happen to appear
 46 in the experience data, but are actually unrelated to the optimal policy, but by chance induce large
 47 cumulative rewards and hence may potentially mislead the learned policy.

48 Theoretically, we establish a data-dependent upper bound on the suboptimality of SAFARI for MF-
 49 MARL without the stringent assumptions on the sufficient coverage of the experience data (e.g., finite
 50 concentrability coefficients [11] or uniformly lower bounded densities of visitation measure [12]).
 51 More specifically, we only assume that the experience data of N training trajectories contains enough
 52 information about the optimal policy. Then we prove that for an episodic MF-MARL problem with
 53 a horizon H , SAFARI attains a sub-optimality gap of $\mathcal{O}(H^2 d_{\text{eff}}/\sqrt{N})$, where d_{eff} is the effective
 54 dimension of the function class (RKHS) for parameterizing the value function and independent on
 55 the number of agents.

56 In addition to the offline settings, our SAFARI algorithm can also be extended to MF-MARL under
 57 the online setting, which is of independent interest. Specifically, we only need to flip the sign of our
 58 uncertainty quantifier, and convert it to a bonus function for promoting the exploration [13]. We then
 59 prove that such an optimistic variant attains a sample efficient regret bound, which is also independent
 60 on the number of agents under the online setting.

61 2 Related Work

62 • **Mean-Field MARL.** Existing literature has proposed various mean-field approximation approaches
 63 to model the population behavior of the agents for MARL with a large number, even infinitely many
 64 homogenous agents. [14] investigate a mean-field game with deterministic linear state transitions, and
 65 reformulate it as a mean-field MDP, where the mean-field state lies in finite-dimensional probability
 66 simplex. [4] propose a mean-field approximation approach over actions, which approximates the
 67 interaction between any given agent and the population by the interaction between the agent’s action
 68 and the averaged actions of its neighboring agents. Such an averaging approach over the local actions,
 69 however, is only applicable when a sparse graph over agents is given, which requires extensive prior
 70 knowledge. [5] investigate a mean-field MDP from the perspective of mean-field control. As the
 71 mean-field state lies in a probability simplex and continuous in nature, they propose to discretize
 72 the joint state-action space such that conventional RL algorithms can be applied. [15] investigate a
 73 mean-field MDP motivated by permutation invariance. They require a central controller managing
 74 the actions of all the agents, and therefore is restricted to handling the curse of many agents from
 75 the exponential blowup of joint state space. More recently, [6] investigate a similar mean-field MDP,
 76 which allows agents to make their own local actions without resorting to a centralized controller. All
 77 these methods focus the online settings. In comparison, our proposed SAFARI algorithm and theory
 78 focus on the offline settings.

79 • **RL for Mean-Field Game.** Our work is also related to the literature that studies RL methods
 80 for mean-field games [16–19]. Such a game can be viewed as the infinite-agent limit of general-
 81 sum Markov game with homogeneous agents, and the aggregated effect of the other agents is also
 82 summarized as a mean-field state. In contrast to mean-field MARL, the solution concept of mean-field
 83 game is the Nash equilibrium, which corresponds to a pair of a local policy π^* of the representative
 84 agent and a mean-field state d^* satisfying the following two properties: (i) when the mean-field state
 85 is set to d^* , π^* is the optimal policy of the representative agent; and (ii) when all agents adopt π^* ,
 86 the resulting mean-field state is d^* . Recently, there are many recent works developing RL methods
 87 for solving mean-field games. See, e.g., [20–30] and the references therein. Most of these methods
 88 adopts a double-loop structure, where the inner loop finds the optimal local policy given the current
 89 mean field state and the outer loop updates the mean-field states. Moreover, these works often assume
 90 the data distribution is well-explored with either a generative model [31] or bounded concentrability
 91 coefficients [32]. Our mean-field MARL problem is similar to the inner-loop problem of finding
 92 the optimal local policy in mean-field games. In contrast to these existing works, our algorithm and

theory can be applied to datasets that are possibly not well-explored. Moreover, as mean-field MARL and mean-field games are different models, our work is not directly comparable to these works.

• **Offline Single-Agent RL.** Our work is also closely related to the literature on offline single-agent RL, which often focuses on either policy evaluation or policy optimization. In particular, in policy evaluation, the goal is to estimate the value function of a target policy, whereas in policy optimization, we aim to learn the optimal policy, which can be achieved via estimating the optimal value function. For both these tasks, in the offline setting, due to the lack of continuing exploration [33], the distribution shift [10] is a fundamental challenge. That is, the trajectories in the dataset and those induced by the target policy or the optimal policy might have diverse distributions. Such a challenge is further exacerbated when function approximators are adopted to represent the desired value functions. To overcome such a challenge, most of the existing theoretical works impose certain well-exploration assumptions on the dataset. Some of commonly made assumptions include uniformly lower bounded visitation measure of the behavior policy, uniformly upper bounded importance sampling ratio, and bounded concentrability coefficients. See, e.g., [34–48, 11, 49–54, 12, 55–63] and the references therein.

However, in practice, such assumptions on the dataset often fail to hold [64–67]. In light of this, there is a line of recent works that proposes various pessimism-based offline single-agent RL algorithms with empirical evidence or theoretical guarantees [68–74]. In particular, [71] propose a regularized variant of fitted Q-iteration [34–36], which is shown to attain the optimal policy within a restricted policy class without assuming the dataset is well-explored. Moreover, with an arbitrary dataset, [72–74] identify the critical role of pessimism in achieving offline sample efficiency. Among these works, our work is particularly related to [73], which develops a pessimistic variant of the value iteration algorithm with finite-dimensional linear function approximation. In comparison, our SAFARI algorithm extends such an algorithm to mean-field MARL and we propose to employ RKHS mean embedding for handling the difference between finite-agent empirical mean-field state and its infinite-agent counterpart. Moreover, our algorithm and analysis involve infinite-dimensional RKHS, which strictly generalizes those in [73].

Notation: Given a space \mathcal{X} , we denote $\mathcal{M}(\mathcal{X})$ as the collection of probability distributions supported on \mathcal{X} . Let $u, v, w \in \mathcal{H}$ be elements in a Hilbert space, we denote $\langle u, v \rangle$ as the inner product, and $u \otimes v$ as the outer product satisfying $(u \otimes v)w = u \langle v, w \rangle$. For a scalar a , we denote $\{a\}^+ = \max\{0, a\}$. We use $\mathcal{O}(\cdot)$ to hide absolute constants and log factors.

3 Mean-Field Multi-Agent RL

We consider a Multi-Agent Reinforcement Learning (MARL) problem with $m + 1$ agents and time horizon H . For the i -th agent (also known as the Representative Agent (RA)), at step h , we denote $s_{i,h} \in \mathcal{S}$ and $a_{i,h} \in \mathcal{A}$ as its state and action, respectively. We assume \mathcal{S} and \mathcal{A} are compact.

Different from single agent RL problem, the transition kernel, reward function, and policy of a representative agent in MARL depend not only on its individual state, but the states of m other agents. Furthermore, we assume that the interaction of the representative agent to the other agents is permutation invariant, i.e., the influence of all the other agents is modeled using the empirical distribution of states $\hat{d}_{s,h} = \frac{1}{m} \sum_{j \neq i}^m \delta_{s_{j,h}} \in \mathcal{M}(\mathcal{S})$. To this end, we define the transition kernel $p_h : \mathcal{S} \times \mathcal{M}(\mathcal{S}) \times \mathcal{A} \mapsto \mathcal{M}(\mathcal{S})$, the (deterministic) reward function $r_h : \mathcal{S} \times \mathcal{M}(\mathcal{S}) \times \mathcal{A} \mapsto \mathbb{R}$, and the policy $\pi_h : \mathcal{S} \times \mathcal{M}(\mathcal{S}) \mapsto \mathcal{M}(\mathcal{A})$ all depending on a “meta state” denoted as $\hat{\omega}_h = (s_{i,h}, \hat{d}_{s,h}) \in \mathcal{S} \times \mathcal{M}(\mathcal{S})$. For simplicity, we denote $\Omega = \mathcal{S} \times \mathcal{M}(\mathcal{S})$ as the meta state space.

Remark 1. The empirical distribution of states $\hat{d}_{s,h}$ is naturally permutation invariant and evolves according to the transition kernel p_h and policy π_h . To see this, suppose each agent takes the same policy π_h at step h . Then at step $h + 1$, the state $s_{h+1,j}$ of the j -th agent is sampled from the distribution $p_h(\cdot | s_{h,j} \times \hat{d}_{s,h}, a_{h,j})$, where $a_{h,j}$ is determined by policy $\pi_h(\cdot | s_{h,j} \times \hat{d}_{s,h})$. Collecting m states $d_{h+1,j}$ for $j \neq i$ induces the empirical distribution of states $\hat{d}_{s,h+1}$.

We now define several important notions in MARL. Given a policy π , the value function $V_h^\pi : \Omega \mapsto \mathbb{R}$ at step $h \leq H$ for a representative agent is

$$V_h^\pi(\omega) = \mathbb{E}_\pi \left[\sum_{i=h}^H r_i(\omega_i, a_i) \mid \omega_h = \omega \right], \quad (1)$$

143 where \mathbb{E}_π denotes the expectation over the randomness in trajectories induced by policy π . The
 144 action-value function (Q -function) $Q_h^\pi : \Omega \times \mathcal{A} \mapsto \mathbb{R}$ is defined as

$$Q_h^\pi(\omega, a) = \mathbb{E}_\pi \left[\sum_{i=h}^H r_i(\omega_i, a_i) \mid \omega_h = \omega, a_h = a \right].$$

145 By definition, V_h^π and Q_h^π are related via $V_h^\pi(\omega) = \int_{\mathcal{A}} Q_h^\pi(\omega, a) \pi(a|\omega) da \triangleq \langle Q_h^\pi, \pi \rangle_{\mathcal{A}}$. Next, we
 146 define the Bellman operator and conditional transition operator. At each step $h \leq H$, the Bellman
 147 operator denoted as \mathbb{B}_h is

$$\begin{aligned} (\mathbb{B}_h g)(\omega, a) &= \mathbb{E} [r_h(\omega_h, a_h) + g(\omega_{h+1}) \mid \omega_h = \omega, a_h = a] \\ &= r_h(\omega, a) + (\mathbb{P}_h g)(\omega, a), \end{aligned} \quad (2)$$

148 where g is a function defined on Ω , and \mathbb{P}_h is referred to as the conditional transition operator.

149 **Mean-Field MARL** As the number of agents goes to infinity, the empirical distribution of states
 150 \widehat{d}_s converges to a (continuous) limit d_s . Then the mean-field MARL problem for a representative
 151 agent is defined as a tuple $(\Omega, \mathcal{A}, H, P, r)$, where Ω and \mathcal{A} are the meta state space and action space,
 152 respectively, H is the horizon, $P = \{p_h\}_{h=1}^H : \Omega \times \mathcal{A} \mapsto \mathcal{M}(\mathcal{S})$ is the transition kernel, and
 153 $r = \{r_h\}_{h=1}^H$ is the reward function defined on $\Omega \times \mathcal{A}$. Following Remark 1, the transition of d_s is
 154 also determined by $P = \{p_h\}_{h=1}^H$.

155 To tackle the infinite-dimensional joint distribution of states, we embed the meta state-action space
 156 $\Omega \times \mathcal{A}$ into a reproducing kernel Hilbert space (RKHS). Specifically, denote $\Xi = \mathcal{S} \times \mathcal{S} \times \mathcal{A}$
 157 and let $K : \Xi \times \Xi \mapsto \mathbb{R}$ be a symmetric positive kernel. The corresponding feature mapping of
 158 kernel k is denoted as ψ , which verifies $\langle \psi(\cdot), \psi(\cdot) \rangle = K(\cdot, \cdot)$ can be infinite dimensional. For any
 159 $(\omega, a) \in \Omega \times \mathcal{A}$, we define mean embedding as

$$\mu(\omega, a) = \mathbb{E}_{s' \sim d_s} [\psi(s, s', a)]. \quad (3)$$

160 Based on the embedding, we parameterize the reward r_h and Markov transition p_h as linear functionals
 161 of $\mu(\omega, a)$ in RKHS \mathcal{H}_K induced by kernel K , i.e.,

$$r_h(\omega, a) = \langle \mu(\omega, a), \theta_h \rangle, \quad p_h(\omega' \mid \omega, a) = \langle \mu(\omega, a), v_h(\omega') \rangle, \quad (4)$$

162 where θ_h, v_h are understood as “weights” and have bounded Hilbert norm (see Assumption 3). Such
 163 a parameterization encodes a rich family of functions, once the kernel is universal [15]. By the
 164 definition of Q -function and value function, we can show that the Bellman operator can also be
 165 parameterized in \mathcal{H}_K .

166 **Proposition 1.** Suppose the reward function r_h and the transition kernel p_h is parameterized in \mathcal{H}_K
 167 by (4) for $h = 1, \dots, H$. Then for any $g : \Omega \mapsto \mathbb{R}$, the Bellman operator $(\mathbb{B}_h g)$ and conditional
 168 transition operator $(\mathbb{P}_h g)$ defined in (2) can be written as

$$(\mathbb{B}_h g)(\omega, a) = \langle \mu(\omega, a), w_g \rangle, \quad (\mathbb{P}_h g)(\omega, a) = \langle \mu(\omega, a), w_g + \theta_h \rangle,$$

169 where w_g depends on the function g .

170 The proof is provided in Appendix C.1, which follows from pure algebraic manipulation. From
 171 the perspective of policy learning in mean-field MARL, Proposition 1 motivates us to estimate the
 172 Bellman operator \mathbb{B}_h in \mathcal{H}_K , and then optimize the estimated Q -function to obtain a policy. We
 173 introduce the detailed learning procedure in Section 4 (Algorithm 1).

174 4 Offline Pessimistic Value Iteration

175 In this section, we introduce our dataset and learning algorithm. We collect multiple trajectories of
 176 a representative agent in a mean-field MARL problem. Here the mean-field state distribution d_s is
 177 prohibitive to trace. Instead, we only independently observe the states of a finite number of agents.
 178 Accordingly, the batched dataset $\mathcal{D}_{N,H}$ consists of N trajectories of length H , within which the n -th
 179 sequence is $\tau_n = \{(s_h^n \in \mathcal{S}^{m+1}, a_h^n \in \mathcal{A}, r_h^n \in \mathbb{R})\}_{h=1}^H$. Without loss of generality, we assume $s_{h,0}$
 180 is the state of the representative agent, and the reward function is bounded by 1, i.e., $|r_h(\omega, a)| \leq 1$
 181 for any $\omega \in \Omega, a \in \mathcal{A}$. The collected trajectories are generated by some unknown behavior policy.

182 Recall $\widehat{d}_{s_h^n} = \frac{1}{m} \sum_{j=1}^m s_{h,j}^n$ is the empirical state distribution induced by s_h^n . (We slightly alter the
 183 notation to emphasize the empirical distribution is generated by the collection of m states $s_{h,1:m}^n$,

184 while in the previous context, we use a general purpose notation $(\widehat{d}_{s,h,\cdot})$. We denote $\widehat{\omega}_h^n = s_{h,0}^n \times \widehat{d}_{s_h^n}$,
 185 and compute the empirical mean embedding of $(\widehat{\omega}_h^n, a_h^n)$ as

$$\mu(\widehat{\omega}_h^n, a_h^n) = \mathbb{E}_{s' \sim \widehat{d}_{s_h^n}} [\psi(s_{h,0}^n, s', a_h^n)] = \frac{1}{m} \sum_{j=1}^m \psi(s_{h,0}^n, s_{h,j}^n, a_h^n).$$

186 Under mild conditions, the empirical mean embedding $\mu(\widehat{\omega}_h^n, a_h^n)$ concentrates around the infinite
 187 agent mean embedding $\mu(\omega_h^n, a_h^n)$ defined in (3), where ω_h^n is the infinite agent meta state. See a
 188 detailed error quantification in Lemma 3.

189 **Pessimistic Value Iteration** Our goal is to learn an optimal policy to be deployed for all the agents
 190 based on the experience data of the representative agent. The idea is to estimate the Q -function at
 191 each time step in the RKHS \mathcal{H}_K , and then optimize the Q -function to obtain an optimal policy. In
 192 more detail, at step $h \leq H$, we estimate Bellman operator by optimizing the empirical mean squared
 193 Bellman error

$$(\widehat{\mathbb{B}}_h \widehat{V}_{h+1}) = \underset{f(\cdot) = \langle \cdot, \alpha \rangle \in \mathcal{H}_K}{\operatorname{argmin}} \sum_{n=1}^N \left(f(\mu(\widehat{\omega}_h^n, a_h^n)) - r_h^n - \widehat{V}_{h+1}(\widehat{\omega}_{h+1}^n) \right)^2 + \lambda \|f\|_{\mathcal{H}}^2, \quad (5)$$

194 where $\lambda \geq 1$ controls the regularization strength, \widehat{V} is the estimated value function, and $\|\cdot\|_{\mathcal{H}}$ denotes
 195 the Hilbert norm.

The solution to (5) can be written in a closed form. For notational simplicity, we define

$$K((\omega, a), \cdot) = \mathbb{E}_{s' \sim d_s} [K((s, s', a), \cdot)] \quad \text{with} \quad \omega = s \times d_s.$$

196 Then we denote the gram matrix $K_h \in \mathbb{R}^{N \times N}$ as

$$[K_h]_{\ell, \ell'} = K((\widehat{\omega}_h^\ell, a_h^\ell), (\widehat{\omega}_h^{\ell'}, a_h^{\ell'})) \triangleq \mathbb{E}_{s_1 \sim \widehat{d}_h^\ell, s_2 \sim \widehat{d}_h^{\ell'}} \langle \psi(s_{h,0}^\ell, s_1, a_h^\ell), \psi(s_{h,0}^{\ell'}, s_2, a_h^{\ell'}) \rangle$$

197 for $\ell, \ell' = 1, \dots, N$. Meanwhile, for any (ω, a) , we denote feature vector $\phi_h(\omega, a) =$
 198 $[K((\widehat{\omega}_h^1, a_h^1), (\omega, a)), \dots, K((\widehat{\omega}_h^N, a_h^N), (\omega, a))]^\top \in \mathbb{R}^N$. Then the estimated Bellman operator
 199 $\widehat{\mathbb{B}}_h \widehat{V}_{h+1}$ can be written as

$$\begin{aligned} (\widehat{\mathbb{B}}_h \widehat{V}_{h+1})(\omega, a) &= \phi_h(\omega, a)^\top \widehat{\alpha}_h \\ \text{with } \widehat{\alpha}_h &= (K_h + \lambda I)^{-1} [r_h^1 + \widehat{V}_{h+1}(\widehat{\omega}_{h+1}^1), \dots, r_h^N + \widehat{V}_{h+1}(\widehat{\omega}_{h+1}^N)]^\top, \end{aligned} \quad (6)$$

200 We summarize the proposed SAFARI algorithm in Algorithm 1.

Algorithm 1 Pessimistic Mean-Field Value Iteration (SAFARI)

Input: Dataset $\mathcal{D}_{N,H}$, coefficient β , regularization coefficient λ .

Initialize: Set $\widehat{V}_{H+1} = 0$.

for $h = H, H-1, \dots, 1$ **do**

 Compute $\Lambda_h = K_h + \lambda I$.

 Estimate $\widetilde{Q}_h(\omega, a) \triangleq (\widehat{\mathbb{B}}_h \widehat{V}_{h+1})(\omega, a) = \phi_h(\omega, a)^\top \widehat{\alpha}_h$ as in (6).

 Set $\Gamma_h(\omega, a) = \beta \cdot \lambda^{-1/2} (K((\omega, a), (\omega, a)) - \phi_h(\omega, a)^\top \Lambda_h^{-1} \phi_h(\omega, a))^{1/2}$.

 Let $\widehat{Q}_h(\omega, a) = \min\{\widetilde{Q}_h(\omega, a) - \Gamma_h(\omega, a), H - h + 1\}^+$.

 Optimal policy $\widehat{\pi}_h = \operatorname{argmax}_{\pi} \langle \widehat{Q}_h(\omega, \cdot), \pi(\cdot | \omega) \rangle_{\mathcal{A}}$.

 Set $\widehat{V}_h(\omega) = \langle \widehat{Q}_h(\omega, \cdot), \widehat{\pi}_h(\cdot | \omega) \rangle_{\mathcal{A}}$.

end for

Output: Estimated Q -function \widehat{Q}_h , value function \widehat{V}_h , and optimal policy $\widehat{\pi}_h$ for $h = 1, \dots, H$.

201 The quantity Γ_h quantifies the uncertainty in estimating the Bellman operator $\widehat{\mathbb{B}}_h \widehat{V}_{h+1}$ using kernel
 202 ridge regression. We subtract Γ_h for estimating the Bellman operator to account for the spurious
 203 correlation in the experience data (see Technical Overview following Theorem 1 for a detailed
 204 explanation). We truncate \widehat{Q}_h at $H - h + 1$, since the reward function is bounded by 1.

205 5 Suboptimality of Policy Learned by SAFARI

206 We investigate the performance of the optimal policy $\hat{\pi}$ learned by Algorithm 1. Before we proceed,
207 we state the following assumptions.

208 **Assumption 1** (Boundedness of Kernel). Kernel $K(\cdot, \cdot)$ is bounded, i.e., without loss of generality,
209 we assume $\sup_{\xi \in \Xi} |K(\xi, \xi)| \leq 1$.

210 By Cauchy-Schwarz inequality, Assumption 1 implies for any $\xi_1, \xi_2 \in \Xi$, $K(\xi_1, \xi_2) \leq$
211 $\sqrt{K(\xi_1, \xi_1)K(\xi_2, \xi_2)} \leq 1$. Such an assumption holds for a rich family of commonly used ker-
212 nels, e.g., RBF kernel and Laplacian kernel, and is a standard assumption in literature [75, 76].

213 The second assumption characterizes the spectrum of kernel K . We first introduce the integral
214 operator induced by kernel K . Let $f : \Xi \mapsto \mathbb{R}$ be a square-integrable function. Then we define the
215 integral operator \mathcal{T}_K as

$$(\mathcal{T}_K f)(\xi) = \int K(\xi, x)f(x)dx \quad \text{for } \xi \in \Xi.$$

216 By Mercer's theorem [77], \mathcal{T}_K has corresponding positive eigenvalues σ_i and eigenfunctions ν_i . Then
217 the kernel K admits a decomposition

$$K(\xi_1, \xi_2) = \sum_{i=1}^{\infty} \sigma_i \nu_i(\xi_1) \nu_i(\xi_2).$$

218 **Assumption 2** (Spectrum of Kernel). The eigenvalue σ_i satisfies one of the following three condi-
219 tions:

- 220 1. (*Finite Spectrum*). There exists a positive integer γ , such that $\sigma_i = 0$ for all $i > \gamma$.
- 221 2. (*Exponential Decay*). There exist positive constants C_1, C_2 and exponent $\gamma > 0$ such that
222 $\sigma_i \leq C_1 \exp(-C_2 i^\gamma)$.
- 223 3. (*Polynomial Decay*). There exists a positive constant C and exponent $\gamma \geq 3 + \mathcal{O}(\frac{1}{d})$ such that
224 $\sigma_i \leq C i^{-\gamma}$, where d is the dimension of $\mathcal{S} \times \mathcal{S} \times \mathcal{A}$.

225 Furthermore, in (*Exponential Decay*) and (*Polynomial Decay*), we assume the eigenfunction ν_i is
226 uniformly bounded, i.e., $\sup_i \|\nu_i\|_\infty \leq 1$.

227 As we will show in our theory, the decay rate of the spectrum significantly influences the performance
228 of the proposed SAFARI algorithm. We give examples to better interpret the three categories above.
229 In (*Finite Spectrum*) case, by (4), the reward function and transition kernel is a linear function of
230 a finite dimensional feature map. Such a parameterization is satisfied by linear MDP [78, 79]. In
231 (*Exponential Decay*) and (*Polynomial Decay*) cases, the feature map is infinite dimensional. For
232 example, RBF kernel belongs to (*Exponential Decay*) case, while Laplacian kernel and neural tangent
233 kernel belong to (*Polynomial Decay*) case. We assume $\gamma \geq 3 + \mathcal{O}(1/d)$ in (*Polynomial Decay*) for
234 technical simplicity, yet it is not restrictive: Laplacian kernel and neural tangent kernel both have a
235 polynomial decay rate of $\gamma = d$ [80].

236 The last assumption imposes some regularity on the reward function and transition probabilities.

237 **Assumption 3** (Boundedness). The weights θ_h and v_h in reward function r_h and Markov transition
238 kernel p_h are bounded for any $h = 1, \dots, H$, respectively, i.e., $\|\theta_h\|_{\mathcal{H}} \leq 1$ and $\int_{\Omega} \|v_h(x)\|_{\mathcal{H}} dx \leq$
239 $\sqrt{d_{\text{eff}}}$, where $d_{\text{eff}} = \sup_{K_h} \log \det(I + K_h/\lambda)$ is the effective dimension of \mathcal{H}_K with supremum
240 over all Gram matrix $K_h \in \mathbb{R}^{N \times N}$.

241 The effective dimension describes the complexity of \mathcal{H}_K for parameterizing the MDP [81], whose
242 scale is closely related to the spectrum of kernel K . In the special case of K having a γ -finite
243 spectrum as in Assumption 2, we have $d_{\text{eff}} = \mathcal{O}(\gamma)$, which resembles the dimensionality of a finite
244 dimensional Euclidean space.

245 We measure the pointwise suboptimality of the learned policy $\hat{\pi}$. We define the global optimal policy
246 by the recursion,

$$\pi_h^* = \operatorname{argmax}_{\pi} \langle Q_h^*, \pi \rangle_{\mathcal{A}}, \quad \text{with } Q_h^* = \mathbb{B}_h V_{h+1}^*, V_h^* = \langle Q_h^*, \pi_h^* \rangle_{\mathcal{A}}, \text{ and } V_{H+1}^* = 0.$$

247 Then the suboptimality of $\hat{\pi}$ is given as

$$\text{SubOpt}(\hat{\pi}; \omega) = V_1^{\pi^*}(\omega) - V_1^{\hat{\pi}}(\omega).$$

248 Our main result is provided in the following theorem, which upper bounds $\text{SubOpt}(\hat{\pi}; \omega)$.

249 **Theorem 1.** Suppose Assumption 1 – 3 hold. For any $\delta \in (0, 1)$, let $\hat{\pi}_h$ be the policy returned by
250 Algorithm 1 with

$$m \geq \log(2/\delta), \quad \lambda = 1, \quad \beta = \begin{cases} c \max\{d, \gamma\} H \sqrt{\log(\max\{d, \gamma\} H N / \delta)} & (\text{Finite Spectrum}) \\ c H \sqrt{d (\log(H N / \delta))^{1+2/\gamma}} & (\text{Exponential Decay}) \\ c N^{\frac{d+1}{d+\gamma}} H \sqrt{d \log(H N / \delta)} & (\text{Polynomial Decay}) \end{cases}$$

251 where d is the dimension of $\Xi = \mathcal{S} \times \mathcal{S} \times \mathcal{A}$ and c is some constant depending on C, C_1, C_2 and
252 Lebesgue measure of Ξ . Then for any meta state ω , with probability at least $1 - \delta$ over the randomness
253 of the dataset $\mathcal{D}_{N,H}$, we have

$$\text{SubOpt}(\hat{\pi}; \omega) \leq 2 \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(\omega_h, a_h) \mid \omega_1 = \omega].$$

254 Theorem 1 indicates that the suboptimality of learned policy depends on the uncertainty quantifier
255 Γ_h . The scale of Γ_h depends on how well the collected data explore the state-action space. Moreover,
256 from a Bayesian learning perspective, Γ_h measures the eliminated uncertainty in estimating the
257 Bellman operator given dataset $\mathcal{D}_{N,H}$ [73]. To better understand the convergence of SubOpt , we
258 specialize Theorem 1 under a weak data coverage assumption.

259 **Assumption 4 (Weak Coverage).** Suppose the dataset is collected under some behavior policy $\bar{\pi}$
260 such that there exists a constant $c_{\min} > 0$ satisfying

$$\inf_{\|f\|_{\mathcal{H}}=1} \langle f, \mathbb{E}_{\bar{\pi}} [\mu(\omega_h, a_h) \otimes \mu(\omega_h, a_h)] f \rangle \geq c_{\min} \quad \text{for any } h = 1, \dots, H.$$

261 Recall that μ is the mean embedding in \mathcal{H}_K .

262 Assumption 4 says that the operator $\mathbb{E}_{\bar{\pi}} [\mu(\omega_h, a_h) \otimes \mu(\omega_h, a_h)]$ is positive definite. Intuitively, this
263 requires that the collected data relatively well spread over the state-action space. We present the
264 following Corollary providing a concrete convergence rate of SubOpt .

265 **Corollary 1.** Under the setting in Theorem 1, we additionally assume Assumption 4 holds. Then for
266 $N \geq \Omega(\log(d_{\text{eff}} H / \delta))$ sufficiently large, with probability $1 - \delta$, we have

$$\text{SubOpt}(\hat{\pi}; \omega) = \mathcal{O} \left(H^2 d_{\text{eff}} \sqrt{\frac{\log(d_{\text{eff}} H N / \delta)}{N}} \right).$$

267 Here d_{eff} is the effective dimension of RKHS \mathcal{H}_K , which takes value

$$d_{\text{eff}} = \begin{cases} \max\{d, \gamma\} \log N & (\text{Finite Spectrum}) \\ d (\log N)^{1+1/\gamma} & (\text{Exponential Decay}) \\ d N^{\frac{d+1}{d+\gamma}} \log N & (\text{Polynomial Decay}) \end{cases}$$

268 **Impact of Kernel Spectrum** The spectrum of kernel K significantly influences the performance of
269 the learned policy. In *(Finite Spectrum)* case, the effective dimension scales linearly with dimension
270 d and γ , and SubOpt converges at a rate of $\mathcal{O}(H^2 \max\{d, \gamma\} / \sqrt{N})$, which recovers the result
271 of Corollary 4.5 in [73] on linear MDP. In *(Exponential Decay)* case, the convergence rate is
272 $\mathcal{O}(H^2 d (\log N)^{1+1/\gamma} / \sqrt{N})$, which is similar to *(Finite Spectrum)* case with additional logarithmic
273 dependence on N . However, in *(Polynomial Decay)* case, the convergence rate is considerably slower,
274 and relies heavily on the decay rate γ . Consider, for instance, Laplacian kernel and NTK, whose
275 spectrum decays with $\gamma = d$. Then SubOpt converges at a rate of $\mathcal{O}(H^2 d N^{-\frac{1}{2d}} \log N)$, which
276 suffers from the curse of dimensionality without further assumptions on data.

277 **No Curse of Many Agents** The convergence of SubOpt does not suffer from the curse of many
278 agents. In particular, both Theorem 1 and Corollary 1 only impose a mild requirement on the number
279 m of neighboring agents to be sampled. This is due to the permutation invariance in mean-field
280 MARL, since the interactive influence of neighboring agents are captured by the distribution of states.

281 **Technical Overview** We briefly discuss the proof of Theorem 1 and Corollary 1. The full proof is
 282 deferred to Appendix A and B. We first decompose SubOpt into three terms (see Lemma 1):

$$\text{SubOpt}(\pi; \omega) = \mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3.$$

283 Here $\mathcal{E}_1 = \sum_{h=1}^H \mathbb{E}_{\hat{\pi}}[\widehat{Q}_h(\omega_h, a_h) - (\mathbb{B}_h \widehat{V}_{h+1})(\omega_h, a_h) \mid \omega_1 = \omega]$ reflects the uncertainty
 284 in estimating the Bellman operator. Note that the evaluating trajectory is generated by the
 285 learned policy $\hat{\pi}$, which has spurious correlation with the estimated Bellman operator; $\mathcal{E}_2 =$
 286 $\sum_{h=1}^H \mathbb{E}_{\pi^*}[(\mathbb{B}_h \widehat{V}_{h+1})(\omega_h, a_h) - \widehat{Q}_h(\omega_h, a_h) \mid \omega_1 = \omega]$ is the estimation error of Bellman oper-
 287 ator again, yet it is evaluated by a trajectory generated by π^* . Compared to \mathcal{E}_1 , \mathcal{E}_2 does not suffer
 288 from the spurious correlation with the learned policy and the estimated Bellman operator. Lastly,
 289 $\mathcal{E}_3 = \sum_{h=1}^H \mathbb{E}_{\pi}[(\widehat{Q}_h(\omega_h, \cdot), \pi_h^*(\cdot \mid \omega_h)) - \widehat{\pi}_h(\cdot \mid \omega_h)]_{\mathcal{A}} \mid \omega_1 = \omega]$ is the optimization error. By the
 290 optimality of $\hat{\pi}$, we immediately have $\mathcal{E}_3 \leq 0$.

291 In order to tackle \mathcal{E}_1 and \mathcal{E}_2 , we properly choose Γ_h so that the event $E = \{|\mathbb{B}_h \widehat{V}_{h+1} - \widehat{\mathbb{B}}_h \widehat{V}_{h+1}| \leq$
 292 $\Gamma_h\}$ happens with high probability. In fact, Γ_h is understood as the uncertainty quantifier of estimating
 293 $\mathbb{B}_h \widehat{V}_{h+1}$ with high confidence $1 - \delta$. Then we can show $\mathcal{E}_1 \leq 0$ conditioned on event E , meanwhile
 294 $\mathcal{E}_2 \leq 2 \sum_{h=1}^H \mathbb{E}_{\pi^*}[\Gamma_h(\omega_h, a_h) \mid \omega_1 = \omega]$. To this end, we reduce the upper bound of SubOpt to
 295 bounding the uncertainty quantifier Γ_h , which allows us to leverage statistical tools. In particular, Γ_h
 296 consists of two types of statistical error: 1) covariate concentration error on mean embedding, i.e.,
 297 finite agent empirical embedding $\mu(\widehat{\omega}_h^n, a_h^n)$ concentration error with respect to population counterpart
 298 $\mu(\omega_h^n, a_h^n)$; 2) regression error in Bellman operator estimation. We bound 1) by concentration of
 299 empirical means in Hilbert spaces (see Lemma 3). In bounding 2), we exploit the closed form solution
 300 of kernel ridge regression and concentration of self-normalizing processes (see Lemma 5).

301 6 Online Policy Learning

302 We extend our pessimistic policy learning to online mean-field settings. Different from performing
 303 backward value iteration, online mean-field MARL aims to learn a policy forward in time. Specifically,
 304 we consider the mean-field MARL problem $(\Omega, \mathcal{A}, H, P, r)$ again, where we denote s_0 as the
 305 representative agent and sample m neighboring agents from the mean-field distribution of states. We
 306 simulate N episodes of length H . In the n -th episode, we sample a starting meta state ω_1^n from a
 307 fixed initial distribution on Ω . However, we can only observe the states of m agents to approximate
 308 the mean-field state distribution, and denote $\widehat{\omega}_1^n$ as the received empirical counterpart of ω_1^n . Then the
 309 representative agent determines a policy $\pi^n = \{\pi_h^n\}_{h=1}^H$, to be deployed for all the agents. At step h ,
 310 the representative agent takes an action a_h^n sampled from $\pi_h^n(\cdot \mid \widehat{\omega}_h^n)$, and transits to the next state
 311 following the Markov transition p_h . The reward function r_h^n is revealed to the representative agent
 312 after taking the action. Note that the reward function can be chosen by the environment adversarially.
 313 The n -th episode ends, when the representative agent receives the terminal reward $r_H^n(\widehat{\omega}_H^n, a_H^n)$. We
 314 define regret of the learned policy, competing with the globally optimal policy:

$$\text{Regret}(\{\pi^n\}_{n=1}^N; H) = \sup_{\pi} \sum_{n=1}^N \left(V_1^{\pi, n}(\omega_1^n) - V_1^{\pi^*, n}(\omega_1^n) \right),$$

315 where the value function V is defined in (1).

316 Similar to the offline setting, we parameterize the transition kernel $P = \{p_h\}_{h=1}^H$ in an RKHS.
 317 Specifically, we slightly abuse the notation to overload $\Xi = \mathcal{S} \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{S}$. We consider a
 318 kernel $K : \Xi \times \Xi \mapsto \mathbb{R}$ with the corresponding feature map ψ . Then we parameterize the transition
 319 kernel p_h as

$$p_h(\omega' \mid \omega, a) = \langle \mu(\omega, a, \omega'), v_h \rangle \quad \text{with} \quad \mu(\omega, a, \omega') = \mathbb{E}_{s \sim d_s, s' \sim d'_s} [\psi(s_0, s, a, s'_0, s')],$$

320 where v_h is the weight and we denote $\omega = (s_0, d_s), \omega' = (s'_0, d'_s)$. Here μ is the mean embedding
 321 depending additionally on the next meta state. We also note the slight variant of parameterizing the
 322 transition kernel compared to the off-line setting in (4). Furthermore, we assume the received reward
 323 $r_h^n(\omega, a) = \mathbb{E}_{s \sim d_s} [g_h^n(s_0, s, a)]$, where g_h^n represents the pairwise reward between the representative
 324 agent and its neighbors. In the sequel, we assume

$$\|v_h\|_{\mathcal{H}} \leq 1 \quad \text{and} \quad \|g_h^n\|_{\infty} \leq 1 \quad \text{for} \quad h = 1, \dots, H, \quad n = 1, \dots, N.$$

325 We propose Online Mean-field Proximal Policy Optimization (OMPPPO) in Algorithm 2 for solving
 326 the online mean-field MARL problem.

Algorithm 2 Online Mean-field Proximal Policy Optimization (OMPPO)

Input: Coefficient β , constant α .

Initialize: Set $\{\pi_0^n\}_{h=1}^H$ to be uniform on \mathcal{A} ; $Q_h^0 = 0$ for $h = 1, \dots, H$.

for $n = 1, \dots, N$ **do**

 Sample starting state $\widehat{\omega}_1^n$.

for $h = 1, \dots, H$ **do**

 Set $\pi_h^n(a | \omega) \propto \pi_h^{n-1}(a | \omega) \exp(\alpha Q_h^{n-1}(\omega, a))$.

 Take action $a_h^n \sim \pi_h^n(\cdot | \widehat{\omega}_h^n)$; receive reward function r_h^n ; transit to next state $\widehat{\omega}_{h+1}^n$.

end for

 Set $V_{H+1}^n = 0$.

for $h = H, H-1, \dots, 1$ **do**

 Compute Gram matrix K_h^{n-1} and vector $\phi_h^n(\omega, a)$: For $\ell, \ell' = 1, \dots, n-1$

$$[K_h^{n-1}]_{\ell, \ell'} = \int_{\Omega \times \Omega} K((\widehat{\omega}_h^\ell, a_h^\ell, \omega'), (\widehat{\omega}_h^{\ell'}, a_h^{\ell'}, \omega'')) V_{h+1}^\ell(\omega') V_{h+1}^{\ell'}(\omega'') d\omega' d\omega''$$

$$[\phi_h^n(\omega, a)]_\ell = \int_{\Omega \times \Omega} K((\widehat{\omega}_h^\ell, a_h^\ell, \omega'), (\omega, a, \omega'')) V_{h+1}^\ell(\omega') V_{h+1}^n(\omega'') d\omega' d\omega''.$$

 Compute $\Lambda_h^{n-1} = K_h^{n-1} + \lambda I$.

 Estimate $\widetilde{Q}_h^n(\omega, a) = \phi_h^{n-1}(\omega, a)^\top (\Lambda_h^{n-1})^{-1} [V_{h+1}^1(\widehat{\omega}_{h+1}^1), \dots, V_{h+1}^{n-1}(\widehat{\omega}_{h+1}^{n-1})]^\top$.

 Set $\Gamma_h^n(\omega, a) = \beta \cdot \lambda^{-1/2} (K((\omega, a), (\omega, a)) - \phi_h^{n-1}(\omega, a)^\top (\Lambda_h^{n-1})^{-1} \phi_h^{n-1}(\omega, a))^{1/2}$.

 Let $Q_h^n(\omega, a) = \min \left\{ \widetilde{Q}_h^n(\omega, a) + r_h^n(\omega, a) + \Gamma_h^n(\omega, a), H - h + 1 \right\}^+$.

 Compute $V_h^n(\omega) = \langle Q_h^n(\omega, \cdot), \pi_h^n(\cdot | \omega) \rangle_{\mathcal{A}}$.

end for

end for

Output: Value function V_h^n and optimal policy π_h^n for $h = 1, \dots, H$ and $n = 1, \dots, N$.

327 Note that different from the off-line setting, OMPPO iteratively improves its policy while accumulating
 328 experience data. During the $(n-1)$ -th episode, OMPPO executes the policy π^{n-1} and collects
 329 the trajectory data of the episode. At the end of the $(n-1)$ -th episode, OMPPO uses the collected
 330 samples to update the estimated Q -functions $\{Q_h^{n-1}\}_{h=1}^H$, and obtains the improved policy π^n using
 331 the updated Q -functions. Further, we flip the sign of Γ_h to promote exploration.

332 Recall that d_{eff} is the effective dimension of \mathcal{H}_K and takes values as in Corollary 1 based on the
 333 spectrum of K . We then present the following regret bound.

334 **Theorem 2.** Suppose Assumption 1 and 2 hold. Given $\delta \in (0, 1)$, we further assume the cardinality
 335 of the action space is bounded by $\log |\mathcal{A}| = \mathcal{O}(d_{\text{eff}}^2 \log^2(d_{\text{eff}} H N / \delta))$. Let π^n be the policy learned
 336 by Algorithm 2 with $\alpha = \sqrt{2 \log |\mathcal{A}| / (H^2 N)}$, $m \geq 2H^4 N^2 \log(2HN/\delta)$, $\lambda = 1$, and

$$\beta = \begin{cases} c \max\{d, \gamma\} H \sqrt{\log(\max\{d, \gamma\} H N / \delta)} & (\text{Finite Spectrum}) \\ cH \sqrt{d (\log(HN/\delta))^{1+2/\gamma}} & (\text{Exponential Decay}) \\ cN^{\frac{d+1}{d+\gamma}} H \sqrt{d \log(HN/\delta)} & (\text{Polynomial Decay}) \end{cases},$$

337 where d is the dimension of Ξ and c is some constant depending on C, C_1, C_2 , and Lebesgue measure
 338 of Ξ . Then with probability $1 - \delta$, we have

$$\text{Regret}(\{\pi^n\}_{n=1}^N; H) = \mathcal{O}\left(H^2 d_{\text{eff}} \sqrt{N} \log(HN/\delta)\right).$$

339 The regret in Theorem 2 does not suffer from the curse of many agents, thanks to the permutation
 340 invariance in mean-field MARL. We also observe that the spectrum of kernel has significant influence
 341 on the regret of learned policy. The proof is provided in Appendix D.

342 7 Conclusion

343 This paper proposes a SAFARI (Pessimistic Mean-Field Value Iteration) algorithm in mean-field
 344 MARL, with an extension to online settings (OMPPO algorithm). We prove a suboptimality bound
 345 $\mathcal{O}(H^2 d_{\text{eff}} / \sqrt{N})$ in the off-line setting, under a weak data coverage assumption. The suboptimality
 346 bound is free of the curse of many agents due to the permutation invariance in mean-field formulation.
 347 In online settings, we provide a regret bound of the order $\mathcal{O}(H^2 d_{\text{eff}} \sqrt{N})$.

348 **References**

- 349 [1] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-
350 agent reinforcement learning in sequential social dilemmas. *arXiv preprint arXiv:1702.03037*,
351 2017.
- 352 [2] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement
353 learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.
- 354 [3] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy
355 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large
356 scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 357 [4] Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field
358 multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages
359 5571–5580. PMLR, 2018.
- 360 [5] René Carmona, Mathieu Laurière, and Zongjun Tan. Model-free mean-field reinforcement
361 learning: mean-field mdp and mean-field q-learning. *arXiv preprint arXiv:1910.12802*, 2019.
- 362 [6] Yan Li, Lingxiao Wang, Jiachen Yang, Ethan Wang, Zhaoran Wang, Tuo Zhao, and Hongyuan
363 Zha. Permutation invariant policy optimization for mean-field multi-agent reinforcement
364 learning: A principled approach. *arXiv preprint arXiv:2105.08268*, 2021.
- 365 [7] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A
366 selective overview of theories and algorithms. *arXiv preprint arXiv:1911.10635*, 2019.
- 367 [8] Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy
368 evaluation across representations with applications to educational games. In *AAMAS*, pages
369 1077–1084, 2014.
- 370 [9] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza,
371 Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement
372 learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- 373 [10] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:
374 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 375 [11] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement
376 learning. In *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- 377 [12] Ming Yin, Yu Bai, and Yu-Xiang Wang. Near optimal provable uniform convergence in
378 off-policy evaluation for reinforcement learning. *arXiv preprint arXiv:2007.03760*, 2020.
- 379 [13] Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy
380 optimization. In *International Conference on Machine Learning*, pages 1283–1294. PMLR,
381 2020.
- 382 [14] Jiachen Yang, Xiaojing Ye, Rakshit Trivedi, Huan Xu, and Hongyuan Zha. Learning deep mean
383 field games for modeling large population behavior. *arXiv preprint arXiv:1711.03156*, 2017.
- 384 [15] Lingxiao Wang, Zhuoran Yang, and Zhaoran Wang. Breaking the curse of many agents:
385 Provable mean embedding q-iteration for mean-field reinforcement learning. In *International
386 Conference on Machine Learning*, pages 10092–10103. PMLR, 2020.
- 387 [16] Minyi Huang, Peter E Caines, and Roland P Malhamé. Individual and mass behaviour in
388 large population stochastic wireless power control problems: centralized and nash equilibrium
389 solutions. In *42nd IEEE International Conference on Decision and Control (IEEE Cat. No.
390 03CH37475)*, volume 1, pages 98–103. IEEE, 2003.
- 391 [17] Jean-Michel Lasry and Pierre-Louis Lions. Jeux à champ moyen. i–le cas stationnaire. *Comptes
392 Rendus Mathématique*, 343(9):619–625, 2006.
- 393 [18] Jean-Michel Lasry and Pierre-Louis Lions. Jeux à champ moyen. ii–Horizon fini et contrôle
394 optimal. *Comptes Rendus Mathématique*, 343(10):679–684, 2006.
- 395 [19] Minyi Huang, Peter E Caines, and Roland P Malhamé. Large-population cost-coupled LQG
396 problems with nonuniform agents: individual-mass behavior and decentralized *var* ϵ -
397 Nash equilibria. *IEEE transactions on automatic control*, 52(9):1560–1571, 2007.
- 398 [20] Xin Guo, Anran Hu, Renyuan Xu, and Junzi Zhang. Learning mean-field games. *arXiv preprint
399 arXiv:1901.09585*, 2019.

- 400 [21] Xin Guo, Renyuan Xu, and Thaleia Zariphopoulou. Entropy regularization for mean field games
401 with learning. *arXiv preprint arXiv:2010.00145*, 2020.
- 402 [22] Xin Guo, Anran Hu, Renyuan Xu, and Junzi Zhang. A general framework for learning mean-
403 field games. *arXiv preprint arXiv:2003.06069*, 2020.
- 404 [23] Zuyue Fu, Zhuoran Yang, Yongxin Chen, and Zhaoran Wang. Actor-critic provably finds
405 nash equilibria of linear-quadratic mean-field games. In *International Conference on Learning*
406 *Representations*, 2019.
- 407 [24] Berkay Anahtarçı, Can Deha Karıksız, and Naci Saldi. Fitted q-learning in mean-field games.
408 *arXiv preprint arXiv:1912.13309*, 2019.
- 409 [25] Berkay Anahtarçı, Can Deha Karıksız, and Naci Saldi. Q-learning in regularized mean-field
410 games. *arXiv preprint arXiv:2003.12151*, 2020.
- 411 [26] Berkay Anahtarçı, Can Deha Karıksız, and Naci Saldi. Value iteration algorithm for mean-field
412 games. *Systems & Control Letters*, 143:104744, 2020.
- 413 [27] Sarah Perrin, Julien Pérolat, Mathieu Laurière, Matthieu Geist, Romuald Elie, and Olivier
414 Pietquin. Fictitious play for mean field games: Continuous time analysis and applications. *arXiv*
415 *preprint arXiv:2007.03458*, 2020.
- 416 [28] Romuald Elie, Julien Pérolat, Mathieu Laurière, Matthieu Geist, and Olivier Pietquin. On the
417 convergence of model free learning in mean field games. In *Proceedings of the AAAI Conference*
418 *on Artificial Intelligence*, volume 34, pages 7143–7150, 2020.
- 419 [29] Muhammad Aneeq uz Zaman, Kaiqing Zhang, Erik Miehling, and Tamer Başar. Reinforcement
420 learning in non-stationary discrete-time linear-quadratic mean-field games. In *IEEE Conference*
421 *on Decision and Control (CDC)*, pages 2278–2284. IEEE, 2020.
- 422 [30] Kai Cui and Heinz Koepl. Approximately solving mean field games via entropy-regularized
423 deep reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*,
424 pages 1909–1917. PMLR, 2021.
- 425 [31] Mohammad Gheshlaghi Azar, Rémi Munos, and Bert Kappen. On the sample complexity of
426 reinforcement learning with a generative model. *arXiv preprint arXiv:1206.6461*, 2012.
- 427 [32] Rémi Munos. Performance bounds in ℓ_p -norm for approximate value iteration. *SIAM journal*
428 *on control and optimization*, 46(2):541–561, 2007.
- 429 [33] Csaba Szepesvári. *Algorithms for reinforcement learning*. Morgan & Claypool, 2010.
- 430 [34] András Antos, Csaba Szepesvári, and Rémi Munos. Fitted Q-iteration in continuous action-space
431 MDPs. In *Advances in Neural Information Processing Systems*, 2007.
- 432 [35] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with
433 Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine*
434 *Learning*, 71(1):89–129, 2008.
- 435 [36] Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of*
436 *Machine Learning Research*, 9(May):815–857, 2008.
- 437 [37] Amir-massoud Farahmand, Csaba Szepesvári, and Rémi Munos. Error propagation for ap-
438 proximate policy and value iteration. In *Advances in Neural Information Processing Systems*,
439 2010.
- 440 [38] Amir-massoud Farahmand, Mohammad Ghavamzadeh, Csaba Szepesvári, and Shie Mannor.
441 Regularized policy iteration with nonparametric function spaces. *Journal of Machine Learning*
442 *Research*, 17(1):4809–4874, 2016.
- 443 [39] Bruno Scherrer, Mohammad Ghavamzadeh, Victor Gabillon, Boris Lesner, and Matthieu Geist.
444 Approximate modified policy iteration and its application to the game of Tetris. *Journal of*
445 *Machine Learning Research*, 16:1629–1676, 2015.
- 446 [40] Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning.
447 In *International Conference on Machine Learning*, pages 652–661, 2016.
- 448 [41] Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforce-
449 ment learning. In *International Conference on Machine Learning*, pages 2139–2148, 2016.

- 450 [42] Mehrdad Farajtabar, Yinlam Chow, and Mohammad Ghavamzadeh. More robust doubly robust
451 off-policy evaluation. In *International Conference on Machine Learning*, pages 1447–1456,
452 2018.
- 453 [43] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon:
454 Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*,
455 2018.
- 456 [44] Tengyang Xie, Yifei Ma, and Yu-Xiang Wang. Towards optimal off-policy evaluation for rein-
457 forcement learning with marginalized importance sampling. In *Advances in Neural Information*
458 *Processing Systems*, pages 9668–9678, 2019.
- 459 [45] Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. DualDICE: Behavior-agnostic estimation
460 of discounted stationary distribution corrections. In *Advances in Neural Information Processing*
461 *Systems*, 2019.
- 462 [46] Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Al-
463 gaeDICE: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- 464 [47] Ziyang Tang, Yihao Feng, Lihong Li, Dengyong Zhou, and Qiang Liu. Doubly robust bias
465 reduction in infinite horizon off-policy estimation. *arXiv preprint arXiv:1910.07186*, 2019.
- 466 [48] Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. GenDICE: Generalized offline
467 estimation of stationary values. In *International Conference on Learning Representations*, 2020.
- 468 [49] Nathan Kallus and Masatoshi Uehara. Efficiently breaking the curse of horizon in off-policy
469 evaluation with double reinforcement learning. *arXiv preprint arXiv:1909.05850*, 2019.
- 470 [50] Nathan Kallus and Masatoshi Uehara. Doubly robust off-policy value and gradient estimation
471 for deterministic policies. *arXiv preprint arXiv:2006.03900*, 2020.
- 472 [51] Nan Jiang and Jiawei Huang. Minimax value interval for off-policy evaluation and policy
473 optimization. In *Advances in Neural Information Processing Systems*, 2020.
- 474 [52] Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and Q-function learning for
475 off-policy evaluation. In *International Conference on Machine Learning*, pages 9659–9668,
476 2020.
- 477 [53] Yaqi Duan, Zeyu Jia, and Mengdi Wang. Minimax-optimal off-policy evaluation with linear
478 function approximation. In *International Conference on Machine Learning*, pages 2701–2709,
479 2020.
- 480 [54] Ming Yin and Yu-Xiang Wang. Asymptotically efficient off-policy evaluation for tabular
481 reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*,
482 pages 3948–3958, 2020.
- 483 [55] Ofir Nachum and Bo Dai. Reinforcement learning via Fenchel-Rockafellar duality. *arXiv*
484 *preprint arXiv:2001.01866*, 2020.
- 485 [56] Mengjiao Yang, Ofir Nachum, Bo Dai, Lihong Li, and Dale Schuurmans. Off-policy evaluation
486 via the regularized Lagrangian. In *Advances in Neural Information Processing Systems*, 2020.
- 487 [57] Zuyue Fu, Zhuoran Yang, and Zhaoran Wang. Single-timescale actor-critic provably finds
488 globally optimal policy. *arXiv preprint arXiv:2008.00483*, 2020.
- 489 [58] Jianqing Fan, Zhaoran Wang, Yuchen Xie, and Zhuoran Yang. A theoretical analysis of deep
490 Q-learning. In *Learning for Dynamics and Control*, pages 486–489, 2020.
- 491 [59] Tengyang Xie and Nan Jiang. Batch value-function approximation with only realizability. *arXiv*
492 *preprint arXiv:2008.04990*, 2020.
- 493 [60] Tengyang Xie and Nan Jiang. Q^* -approximation schemes for batch reinforcement learning: A
494 theoretical comparison. *arXiv preprint arXiv:2003.03924*, 2020.
- 495 [61] Peng Liao, Zhengling Qi, and Susan Murphy. Batch policy learning in average reward Markov
496 decision processes. *arXiv preprint arXiv:2007.11771*, 2020.
- 497 [62] Junyu Zhang, Alec Koppel, Amrit Singh Bedi, Csaba Szepesvári, and Mengdi Wang. Variational
498 policy gradient method for reinforcement learning with general utilities. In *Advances in Neural*
499 *Information Processing Systems*, 2020.
- 500 [63] Tongzheng Ren, Jialian Li, Bo Dai, Simon S Du, and Sujay Sanghavi. Nearly horizon-free
501 offline reinforcement learning. *arXiv preprint arXiv:2103.14077*, 2021.

- 502 [64] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning
503 without exploration. In *International Conference on Machine Learning*, pages 2052–2062,
504 2019.
- 505 [65] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective
506 on offline reinforcement learning. In *International Conference on Machine Learning*, pages
507 104–114, 2020.
- 508 [66] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: Datasets for
509 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 510 [67] Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Tom Le Paine, Sergio Gómez Colmenarejo,
511 Konrad Zolna, Rishabh Agarwal, Josh Merel, Daniel Mankowitz, Cosmin Paduraru, et al. RL
512 Unplugged: Benchmarks for offline reinforcement learning. *arXiv preprint arXiv:2006.13888*,
513 2020.
- 514 [68] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea
515 Finn, and Tengyu Ma. MOPO: Model-based offline policy optimization. *arXiv preprint*
516 *arXiv:2005.13239*, 2020.
- 517 [69] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. MOREL:
518 Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- 519 [70] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-learning for
520 offline reinforcement learning. *arXiv preprint arXiv:2006.04779*, 2020.
- 521 [71] Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably good batch
522 reinforcement learning without great exploration. *arXiv preprint arXiv:2007.08202*, 2020.
- 523 [72] Jacob Buckman, Carles Gelada, and Marc G Bellemare. The importance of pessimism in
524 fixed-dataset policy optimization. *arXiv preprint arXiv:2009.06799*, 2020.
- 525 [73] Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline RL?
526 *arXiv preprint arXiv:2012.15085*, 2020.
- 527 [74] Chenjun Xiao, Yifan Wu, Tor Lattimore, Bo Dai, Jincheng Mei, Lihong Li, Csaba Szepesvari,
528 and Dale Schuurmans. On the optimality of batch policy optimization algorithms. *arXiv preprint*
529 *arXiv:2104.02293*, 2021.
- 530 [75] Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares
531 algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- 532 [76] Krikamol Muandet, Kenji Fukumizu, Francesco Dinuzzo, and Bernhard Schölkopf. Learning
533 from distributions via support measure machines. *arXiv preprint arXiv:1202.6504*, 2012.
- 534 [77] Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. Support
535 vector machines. *IEEE Intelligent Systems and their applications*, 13(4):18–28, 1998.
- 536 [78] Lin Yang and Mengdi Wang. Sample-optimal parametric Q-learning using linearly additive
537 features. In *International Conference on Machine Learning*, pages 6995–7004, 2019.
- 538 [79] Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement
539 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–
540 2143, 2020.
- 541 [80] Alberto Bietti and Francis Bach. Deep equals shallow for relu networks in kernel regimes. *arXiv*
542 *preprint arXiv:2009.14397*, 2020.
- 543 [81] Zhuoran Yang, Chi Jin, Zhaoran Wang, Mengdi Wang, and Michael I. Jordan. On function
544 approximation in reinforcement learning: Optimism in the face of large state spaces, 2020.
- 545 [82] Yasemin Altun and Alex Smola. Unifying divergence minimization and statistical inference
546 via convex duality. In *International Conference on Computational Learning Theory*, pages
547 139–153. Springer, 2006.
- 548 [83] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear
549 stochastic bandits. In *NIPS*, volume 11, pages 2312–2320, 2011.
- 550 [84] Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International*
551 *Conference on Machine Learning*, pages 844–853. PMLR, 2017.

552 **Checklist**

- 553 1. For all authors...
- 554 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
555 contributions and scope? [Yes]
- 556 (b) Did you describe the limitations of your work? [No]
- 557 (c) Did you discuss any potential negative societal impacts of your work? [No] *We do not*
558 *foresee direct negative societal impacts of this work.*
- 559 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
560 them? [Yes]
- 561 2. If you are including theoretical results...
- 562 (a) Did you state the full set of assumptions of all theoretical results? [Yes] *See Assumption*
563 *1 – 4 in Section 5.*
- 564 (b) Did you include complete proofs of all theoretical results? [Yes] *The full proofs of the-*
565 *oretical results are provided in the supplementary, organized in sections corresponding*
566 *to the main text of the paper.*
- 567 3. If you ran experiments...
- 568 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
569 perimental results (either in the supplemental material or as a URL)? [N/A] *This*
570 *submission does not provide empirical results.*
- 571 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
572 were chosen)? [N/A]
- 573 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
574 ments multiple times)? [N/A]
- 575 (d) Did you include the total amount of compute and the type of resources used (e.g., type
576 of GPUs, internal cluster, or cloud provider)? [N/A]
- 577 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 578 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 579 (b) Did you mention the license of the assets? [N/A]
- 580 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 581
- 582 (d) Did you discuss whether and how consent was obtained from people whose data you’re
583 using/curating? [N/A]
- 584 (e) Did you discuss whether the data you are using/curating contains personally identifiable
585 information or offensive content? [N/A]
- 586 5. If you used crowdsourcing or conducted research with human subjects...
- 587 (a) Did you include the full text of instructions given to participants and screenshots, if
588 applicable? [N/A]
- 589 (b) Did you describe any potential participant risks, with links to Institutional Review
590 Board (IRB) approvals, if applicable? [N/A]
- 591 (c) Did you include the estimated hourly wage paid to participants and the total amount
592 spent on participant compensation? [N/A]