
Graphon Mean Field Games with A Representative Player: Analysis and Learning Algorithm

Anonymous Authors¹

Abstract

We propose a discrete-time graphon game formulation on continuous state and action spaces using a representative player to study stochastic games with heterogeneous interaction among agents. This formulation admits both conceptual and mathematical advantages, compared to a widely adopted formulation using a continuum of players. We prove the existence and uniqueness of the graphon equilibrium with mild assumptions, and show that this equilibrium can be used to construct an approximate solution for the finite player game, which is challenging to analyze and solve due to curse of dimensionality. An online oracle-free learning algorithm is developed to solve the equilibrium numerically, and sample complexity analysis is provided for its convergence.

1. Introduction

Many real-world applications, such as flocking (Perrin et al., 2021), epidemiology (Cui et al., 2022), and autonomous driving (Huang et al., 2020) involve multiagent systems, where agents optimize individual cumulative rewards by selecting sequential actions in an (in)finite horizon, while interacting strategically among one another. In discrete-time, such finite player games form Markov games (Littman, 1994; Solan & Vieille, 2015; Yang et al., 2018b). At a Nash equilibrium (NE), nobody can improve her payoff by unilaterally switching her action policies. The NE is challenging to solve when the agent size gets larger. To address such a challenge, mean-field formulations are proposed to model individuals interacting with others only via an aggregate population.

A school of researchers define a new type of games, namely, mean fields games (MFG) (Lasry & Lions, 2007; Huang

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

et al., 2006) that describe the limiting behavior of its corresponding finite player game as the number of players is large, assuming that the interaction among players are homogeneous.

As a generalization to MFGs, graphon mean field games (GMFG or graphon games) are developed (Caines & Huang, 2021; Gao et al., 2021; Aurell et al., 2022; Tangpi & Zhou, 2023; Cui & Koepll, 2022) to tackle the limiting behavior of finite player games with *heterogeneous* agents who interact *asymmetrically*, deemed as games on networks. GMFGs cover a broader range of models and applications, as it allows infinitely many distinct types of player with flexible heterogeneous interaction, which are normally modeled by graphons, a natural limit of finite graphs when the number of vertices goes to infinity.

Because GMFGs may not be solved explicitly in general, recent years have seen a growing tend of using learning methods for equilibria. Compared to abundant studies on learning MFG (Yang et al., 2018a; Guo et al., 2019; Cardaliaguet & Hadikhanloo, 2015; Elie et al., 2020; Perrin et al., 2020; 2021; 2022; Lauriere et al., 2022; Chen et al., 2023a;b), learning on graphon games (Cui & Koepll, 2022; Zhang et al., 2023) is relatively understudied.

A major roadblock in learning GMFG lies in the fact that there is no consensus on what a mathematically tractable formulation of GMFG should be, since it is not straightforward to describe the limiting behavior of large number of heterogeneous players. To the best of our knowledge, there are two types of formulations. The first type, also the widely adopted one, models a game for uncountably infinite players with distinct types (Caines & Huang, 2021), so-called “continuum-player” games (Carmona et al., 2021).

Unfortunately, this formulation suffers from limitations. Theoretically, the mapping from players’ types to state dynamic is not measurable under the usual σ -algebra, which potentially pose challenges in analytical investigation of solution properties (Appendix B.2); And practically, it is difficult to develop an algorithm that directly solves a system of optimal control problems for a continuum of players. Moreover, these studies could lack consistency between the formulations (that model infinitely many players) and the

| Reference Perspective | Formulation | Assumptions | | Analysis | | Algorithms | | |
|------------------------|----------------|-------------|------------|--------------------|------------|-------------|-------------|---------------------|
| Perspective | Player-type | Time domain | Graphon | State action space | Uniqueness | Approx Eqbm | Oracle-free | Complexity analysis |
| (Cui & Koepll, 2022) | Continuum | Discrete | Continuous | Finite | ✗ | ✓ | ✗ | ✗ |
| (Zhang et al., 2023) | Continuum | Discrete | Continuous | Finite | ✗ | ✗ | ✗ | ✗ |
| (Lacker & Soret, 2022) | Representative | Continuous | General | \mathbb{R}^d | ✓ | ✓ | NA | NA |
| This paper | Representative | Discrete | General | \mathbb{R}^d | ✓ | ✓ | ✓ | ✓ |

algorithms (that only sample a single representative agent).

To tackle the aforementioned challenges, a second kind of formulation refers to a generic player who represents all types of agents while interacting with the aggregate population (Lacker & Soret, 2022). This formulation is amenable to theoretical guarantees and more importantly, ease the algorithmic design and implementation.

In this paper, we study discrete-time graphon games of the second formulation with rigorous analysis and learning methods. We start from finite player games to motivate graphon games, which in turn provide approximate equilibria to finite games in dense interaction networks. Subsequently, GMFG always refer to representative-player graphon game, unless otherwise specified.

Related work. Tab. 1 compares the most relevant studies on learning GMFGs. **Continuum-player formulation:** In discrete time regime, Cui & Koepll (2022) showed existence and approximate equilibrium under Lipschitz transition kernel and graphon, and Zhang et al. (2023) only showed the existence of GMFGs with entropic regularization. Both studies assumed access to an oracle that returns the population dynamics, and the latter further assumes access to an action-value function oracle that returns the optimal policies. Under these assumptions, Zhang et al. (2023) provides a convergence rate of their algorithm, while Cui & Koepll (2022) only shows the asymptotic convergence. In continuous time regime, Caines & Huang (2021) focused on finite networks where each vertex represents a population. Gao et al. (2021); Aurell et al. (2022); Tangpi & Zhou (2023) studied linear quadratic games, and the latter two adopted rich Fubini extension to address the measurability issue. **Representative-player formulation:** As the establisher and the only work to the best of our knowledge, Lacker & Soret (2022) rigorously studied the equilibrium existence uniqueness and approximate equilibrium in continuous-time, with no discussion in algorithm implementation.

Contributions. Our major contributions are:

- As opposed to the widely used formulation of GMFG with a continuum of players, we offer a new formulation with only one representative player, which inherits the spirit of classic MFGs, and more importantly, provide technical advantages (see Appendix B.2).
- Our model framework is general in terms of state-action space and transition kernel. We allow the state and action

spaces to be Euclidean spaces, as opposed to the finite state-action space. The state dynamic transitions and reward functions are allowed to be time-variant.

- We present comprehensive analysis of mathematical properties of our GMFG, namely, equilibrium existence, uniqueness and approximate equilibrium convergence. All rely on weaker assumptions than those used in existing studies, such as continuous graphon and Lipschitz state transition dynamics. See Secs. 4 and D-F.
- We provide the first fully online oracle-free learning scheme for solving the equilibrium, and justify its efficiency with a sample complexity analysis. See Sec. 5.

2. Preliminaries

2.1. Notations

Let E be any Polish space (complete separable metric topological space). We use $\mathcal{P}(E)$ to represent all the probability measures on E equipped with the weak topology, with \Rightarrow being the weak convergence. Let $\mathcal{M}_+(E)$ denote the space of nonnegative Borel measures of finite variation. Denote $\|\cdot\|_{\text{TV}}$ the total variation norm. Given a random element X valued in E , let $\mathcal{L}(X) \in \mathcal{P}(E)$ be the probabilistic law (distribution) of X . For any $\mu \in \mathcal{P}(E)$, we write $X \sim \mu$ if $\mathcal{L}(X) = \mu$. For simplicity, we represent the integral with $\langle \mu, \phi \rangle = \int_E \phi d\mu$ for $\mu \in \mathcal{M}_+(E)$ and measurable ϕ .

Let $\mathcal{P}_{\text{unif}}([0, 1] \times E)$ denote a measure on product space $[0, 1] \times E$ with uniform first marginal. We always consider E to be a regular space, and thus each element μ admits a disintegration $d\mu^u(dx)$ where $\mu^u(dx)$ is a Lebesgue almost every uniquely defined kernel $[0, 1] \rightarrow E$.

2.2. Graphon

2.2.1. DEFINITION

A graphon W is an L_1 integrable function : $[0, 1]^2 \rightarrow \mathbb{R}_+$. It represents a graph with infinitely many vertices taking labels in $[0, 1]$, and the edge weight connecting vertex u and v is given by $W(u, v)$. It is a natural notion for the limit of a sequence of graphs as the size of vertices grows.

Any finite graph can be expressed equivalently as a graphon: given any graph on $n \geq 1$ vertices with non-negative edge weights, it can be equivalently expressed as a matrix $\xi \in \mathbb{R}_+^{n \times n}$, where ξ_{ij} is the edge weight between vertex i and j . We define a *step graphon associated with* ξ , denoted as W_ξ

110 on $[0, 1]^2$ below:

$$111 \quad W_\xi(u, v) := \sum_{i,j=1}^n \xi_{ij} 1_{\{u \in I_i^n, v \in I_j^n\}}$$

112 where the interval of $[0, 1]$ is divided into n bins with the
 113 i^{th} bin as $I_i^n := [(i-1)/n, i/n], \forall i = 1, \dots, n-1; I_n^n := [(n-1)/n, 1]$.

114 2.2.2. GRAPHON OPERATOR

115 Given a Polish space E and any graphon W , the graphon
 116 operator \mathbf{W} , which maps a measure on $\mathcal{P}_{\text{unif}}([0, 1] \times E)$ to
 117 a function $[0, 1] \rightarrow \mathcal{M}_+(E)$, is defined as follows (Lacker
 118 & Soret, 2022): for any $m \in \mathcal{P}_{\text{unif}}([0, 1] \times E)$,

$$119 \quad \mathbf{W}m(u) := \int_{[0,1] \times E} W(u, v) \delta_x m(dv, dx) \quad (1)$$

120 where δ_x is Dirac delta measure at x . Intuitively, assume
 121 m admits disintegration $m(du, dx) = dvm^u(dx)$, and W
 122 represents a graph with infinitely many vertices where each
 123 vertex $u \in [0, 1]$ bears a random value on E with distri-
 124 bution m^u . Then $\mathbf{W}m(u) = \int_{[0,1] \times E} W(u, v) \delta_x m_v(dx) dv$
 125 is an average of the distribution of the random value over
 126 all vertices, weighted by edges with u as one end. Note
 127 that $\mathbf{W}m(u) \in \mathcal{M}_+(E)$ since the weighted average may no
 128 longer be a probability measure.

129 2.2.3. STRONG OPERATOR TOPOLOGY

130 Now we define convergence of graphons in strong operator
 131 topology. We abuse the notation by denoting the usual
 132 integral operator $\mathbf{W} : L_\infty[0, 1] \rightarrow L_1[0, 1]$,

$$133 \quad \mathbf{W}\phi(u) := \int_{[0,1]} W(u, v) \phi(v) dv \quad \forall \phi \in L_\infty[0, 1] \quad (2)$$

134 and it should lead to no ambiguity as graphon operators
 135 and integral operators have different domains. We say a
 136 sequence of graphons W^n converges to a limit graphon W
 137 in the strong operator topology if for any $\phi \in L_\infty[0, 1]$,
 138 $\|\mathbf{W}^n\phi - \mathbf{W}\phi\|_1 \rightarrow 0$, denoted as $W^n \rightarrow W$. Convergence
 139 in strong operator topology is usually weaker than converge
 140 in cut norm, see appendix A.5.

141 3. Finite Player Games

142 Consider a game with $n \in \mathbb{N}_+$ players. Let $\xi \in \mathbb{R}_+^{n \times n}$
 143 be an interaction matrix with nonnegative entries, where
 144 ξ_{ij} is the interaction influence of player j onto player i
 145 for $i, j \in [n]$. Let $T \in \mathbb{N}_+$ be terminal time of the game,
 146 and $\mathbb{T} := \{0, 1, 2, \dots, T-1\}$. At each time t , denote
 147 $\mathbf{X}_t = (X_t^1, \dots, X_t^n) \in (\mathbb{R}^d)^n$ the state dynamics of all the
 148 players, i.e. each player's state takes value in \mathbb{R}^d for some
 149 fixed $d \geq 1$, and let $\mathcal{C} := (\mathbb{R}^d)^{T+1}$ be the space of state
 150 paths. For any $x \in \mathcal{C}$, write x_t the value of path at time
 151 t . The initial states \mathbf{X}_0 follow a vector of initial measures
 152 $\lambda = (\lambda^1, \dots, \lambda^n) \in (\mathcal{P}(\mathbb{R}^d))^n$. At each time every player
 153

154 may choose an action from the action space A , and we
 155 assume that $A \subset \mathbb{R}^d$ is compact. Let \mathcal{A}_n be the collection of
 156 all feedback policies $\mathbb{T} \times (\mathbb{R}^d)^n \rightarrow \mathcal{P}(A)$, and each player's
 157 action follows a policy from this collection. For any policy
 158 $\pi^i \in \mathcal{A}_n$ chosen by player i , the state process of player i
 159 evolves by a transition kernel $P : \mathbb{T} \times \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \times A \rightarrow$
 160 $\mathcal{P}(\mathbb{R}^d)$ as follows

$$161 \quad X_0^i \sim \lambda^i \\ 162 \quad a_t^i \sim \pi_t^i(\mathbf{X}_t) \quad X_{t+1}^i \sim P_t(X_t^i, M_t^i, a_t^i)$$

163 for $i = 1, \dots, n$, where

$$164 \quad M^i := \frac{1}{n} \sum_{j=1}^n \xi_{ij} \delta_{X^j} \in \mathcal{M}_+(\mathcal{C})$$

165 is the empirical weighted neighborhood measure of player
 166 i , and M_t^i is the time t marginal of M^i . “Empirical” means
 167 the measure is an average of the Dirac measures at the
 168 realizations, in particular, M^i is a random measure. M^i
 169 depicts an average of all players' states, weighted by their
 170 influence on player i . At each time, every player chooses an
 171 action according to her policy, and her state process X is
 172 Markov decision process (MDP), which now depends not
 173 only on her current state and action, but also the empirical
 174 weighted neighborhood measure. Note that at each time
 175 t , the policy π^i of player i may depend on each of other
 176 players' state, while the transition law P should only depend
 177 on other players by an aggregation of their states, i.e. the
 178 empirical weighted neighborhood measure.

179 At each time all players receive running reward according
 180 to some $f : \mathbb{T} \times \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \times A \rightarrow \mathbb{R}$ and they receive
 181 a terminal reward at the terminal time T according to some
 182 function $g : \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \rightarrow \mathbb{R}$. The objective of player
 183 i is to maximize her expected accumulated reward

$$184 \quad J^i(\boldsymbol{\pi}) := \mathbb{E} \left[\sum_{t \in \mathbb{T}} f_t(X_t^{\boldsymbol{\pi}, i}, M_t^{\boldsymbol{\pi}, i}, a_t^{\boldsymbol{\pi}, i}) + g(X_T^{\boldsymbol{\pi}, i}, M_T^{\boldsymbol{\pi}, i}) \right]$$

185 which is a function of the policy of all players $\boldsymbol{\pi} = (\pi^1, \dots, \pi^n) \in (\mathcal{A}_n)^n$. We write $X^{\boldsymbol{\pi}, i}$, $M^{\boldsymbol{\pi}, i}$ and $a^{\boldsymbol{\pi}, i}$
 186 to emphasize that the state dynamic of player i depends on
 187 $\boldsymbol{\pi}$.

188 **Definition 3.1.** For any nonnegative vector $\epsilon = (\epsilon^1, \dots, \epsilon^n) \in \mathbb{R}_+^n$, an ϵ -equilibrium of the n -player game
 189 is defined as $\hat{\boldsymbol{\pi}} = (\hat{\pi}^1, \dots, \hat{\pi}^n) \in (\mathcal{A}_n)^n$ such that for an i ,

$$190 \quad J^i(\hat{\boldsymbol{\pi}}) \geq \sup_{\boldsymbol{\pi} \in \mathcal{A}_n} J^i(\hat{\boldsymbol{\pi}}^{-i}, \boldsymbol{\pi}) - \epsilon_i \quad (3)$$

191 where $(\hat{\boldsymbol{\pi}}^{-i}, \boldsymbol{\pi})$ denotes the vector $\hat{\boldsymbol{\pi}}$ with i^{th} coordinate
 192 replaced by $\boldsymbol{\pi}$.

193 **Mapping n -player indices onto a continuous label space.**
 194 This part serves as a transition from finite player game
 195 defined above, to its limiting system in the next section.
 196 In the finite n -player game, we map the index of agent
 197 $i \in \{1, \dots, n\}$ onto a continuous label space $[0, 1]$, by
 198 assigning player i a label $u_i \in I_i^n := [(i-1)/n, i/n], \forall i = 1, \dots, n-1$ and $u_n \in I_n^n := [(n-1)/n, 1]$.

We demonstrate that the empirical weighted neighborhood measure M^i can be expressed in terms of the graphon operator. Let W_ξ be the step graphon associated with interaction matrix ξ , then the interaction between player i and j can be expressed by $\xi_{ij} = W_\xi(u_i, u_j)$. Define the empirical label-state joint measure

$$S := \frac{1}{n} \sum_{i=1}^n \delta_{(u_i, X^i)} \in \mathcal{P}([0, 1] \times \mathcal{C}) \quad (4)$$

which is an empirical measure of the label-state pairs of all players. Then we have for $i = 1, \dots, n$,

$$M^i = \frac{1}{n} \sum_{j=1}^n \xi_{ij} \delta_{X^j} = \int W(u_i, v) \delta_x S(dv, dx) = W_\xi S(u_i) \quad (5)$$

This demonstrates that the graphon operator is a generalization of the weighted neighborhood measure when there are infinitely many players: with W being the interaction among a continuum of players, and μ being their population label-state joint measure, $W\mu(u)$ is the weighted neighborhood measure for player of label $u \in [0, 1]$.

4. Representative-player Graphon Game

4.1. Game formulation

Given a graphon $W \in L_+^1[0, 1]^2$ representing the interactions of a continuum types of players, we define the graphon game associated with W for a single representative player as follows. Let the state and action space be defined as in Section 3. Let $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ be a filtered probability space that supports an \mathcal{F}_0 -measurable random variable U uniform on $[0, 1]$, and an adapted Markov process X valued in \mathbb{R}^d . We understand U as the label for the representative player, and X as her state dynamic. The initial label-state law of the representative player is given by $\lambda := \mathcal{L}(U, X_0) \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d)$. The term “label-state” always refer to the joint measure of a player’s label and state pair (U, X) . Let $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$ be fixed, and we understand it as the label-state joint measure of the population that the representative player reacts to. As in mean-field games, all other players except the representative player are abstracted into μ . Let $\mu_t \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d)$ be the marginal of μ under image $(u, x) \mapsto (u, x_t)$.

Let \mathcal{V}_U be the collections of all the open-loop policies, i.e. all the adapted process valued in $\mathcal{P}(A)$. Let \mathcal{A}_U denotes the collection of all the closed-loop (Markovian) policies, i.e. measurable functions $\mathbb{T} \times [0, 1] \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$. \mathcal{A}_U is usually a proper subset of \mathcal{V}_U , unless the filtration is generated by U and X . For any $\pi \in \mathcal{V}_U$, the label-state pair (U, X) follows the transition dynamic $(U, X_0) \sim \lambda$ and at each $t \in \mathbb{T}$,

$$a_t \sim \pi_t \quad X_{t+1} \sim P_t(X_t, W\mu_t(U), a_t)$$

for the same $\{P_t\}_{t \in \mathbb{T}}$ as in the finite player game introduced in section 3. In words, the representative player is uniformly

assigned a label U at time 0, and her later state transition depends on her current state, action and weighted neighborhood measure $W\mu_t(U) \in \mathcal{P}(\mathbb{R}^d)$. Recall the identity in equation (5), μ is now a generalization of S defined in (4) when there are infinitely many types of players, and $W\mu(u)$ is the distribution of the states of the population (infinitely many other players), reweighted by their interaction with the representative player when her label is $u \in [0, 1]$.

Let $f : \mathbb{T} \times \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \times A \rightarrow \mathbb{R}$ be the running reward and $g : \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \rightarrow \mathbb{R}$ be the terminal reward. The objective of the representative player is to choose a policy $\pi \in \mathcal{V}_U$ to maximize

$$J_W(\mu, \pi) := \mathbb{E} \left[\sum_{t \in \mathbb{T}} f_t(X_t^\pi, W\mu_t(U), a_t^\pi) + g(X_T^\pi, W\mu_T(U)) \right]$$

Note that the expectation is w.r.t. all random elements on \mathcal{F} , i.e. (U, X) and π , and we use X^π , a^π to emphasize that they depends on the policy π .

Definition 4.1. We say that the measure-policy pair $(\hat{\mu}, \hat{\pi}) \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C}) \times \mathcal{V}_U$ is a W -equilibrium if

$$J_W(\hat{\mu}, \hat{\pi}) = \sup_{\pi \in \mathcal{V}_U} J_W(\hat{\mu}, \pi) \quad (6)$$

$$\hat{\mu} = \mathcal{L}(U, X^{\hat{\pi}}) \quad (7)$$

$\hat{\mu}, \hat{\pi}$ are called equilibrium population measure and equilibrium optimal policy respectively.

Intuitively, the game is formulated for a representative player, while all other players are abstracted into a label-state joint measure μ . The representative player interacts with the population only through the weighted neighborhood measure $W\mu(U)$, according to which she takes action to optimize her reward. We give a comprehensive comparison between our formulation with the continuum-player graphon game in Appendix B.

Remark 4.2. We define an infinite horizon version of graphon game with time-invariant dynamics and rewards in Appendix A.1. The analysis in the rest of this section could be easily adapted to the infinite horizon formulation by eliminating the time dependency of functions.

GMFG as MFG with augmented state space. The graphon games defined here could be transformed into classical MFGs with an augmented state space, by imposing the label space $[0, 1]$ as an additional dimension to the state space. However, this does not simplify the analysis or proof, and it is not appropriate to adapt existing results from MFG directly (See Appendix A.2).

4.2. Existence of Equilibrium

Assumption 4.3. 1. The action space A is a compact subspace of \mathbb{R}^d .
 2. The running rewards $f_t, \forall t \in \mathbb{T}$ and terminal reward g are bounded and jointly continuous.

- 220 3. The intial distribution $\lambda \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d)$ admits
 221 disintegration $\lambda(du, dx) = du\lambda^u(dx)$, and the following
 222 collection of measures is tight¹:

$$\{\lambda^u\}_{u \in [0, 1]} \subset \mathcal{P}(\mathbb{R}^d)$$

- 223 4. For each $t \in \mathbb{T}$, the following collection of measures is
 224 tight:

$$\zeta_t := \{P_t(x, m, a)\}_{(x, m, a) \in \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \times A} \subset \mathcal{P}(\mathbb{R}^d)$$

- 225 5. For each $t \in \mathbb{T}$, $P_t(x, m, \cdot)$ is continuous in A for every
 226 $(x, m) \in \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d)$.

An example case where Assumption 4.3(3, 4) are trivially satisfied is that there exists some compact subspace $\mathcal{X} \subset \mathbb{R}^d$ such that the collection ζ_t are uniformly supported on \mathcal{X} , i.e. the Markov process X takes values in the state space \mathcal{X} . Also note that we do not assume the graphon W to be continuous.

Theorem 4.4. Suppose assumption 4.3 holds. Then there exists a W -equilibrium $(\mu, \hat{\pi})$. Moreover, the equilibrium optimal policy $\hat{\pi}$ can be chosen to be a closed-loop policy.

The theorem is proved with probabilistic compactification and Kakutani-Fan-Glicksberg fixed point theorem in Appendix D.

4.3. Uniqueness of Equilibrium

Assumption 4.5. 1. The state transition law P does not depend on the measure argument. Then it reads $P_t : \mathbb{R}^d \times A \rightarrow \mathcal{P}(\mathbb{R}^d)$ for $t \in \mathbb{T}$.

2. For each $t \in \mathbb{T}$, the running reward f_t is separable in the measure and action argument: there exists $f_t^1 : \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \rightarrow \mathbb{R}$ and $f_t^2 : \mathbb{R}^d \times A \rightarrow \mathbb{R}$ such that $f_t(x, m, a) = f_t^1(x, m) + f_t^2(x, a)$.

3. The optimal policy is unique. More specifically, for each $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$, the supremum $\sup_{\pi \in \mathcal{V}_U} J_W(\mu, \pi)$ is attained uniquely.

4. The functions f_t^1 and g satisfy the Larsi-Lions Monotonicity condition, in the following sense: for any $m_1, m_2 \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d)$, and $t \in \mathbb{T}$,

$$\int_{[0, 1] \times \mathbb{R}^d} (g(x, \mathbf{W}m_1(u)) - g(x, \mathbf{W}m_2(u)))$$

$$(m_1 - m_2)(du, dx) \leq 0$$

$$\int_{[0, 1] \times \mathbb{R}^d} (f_t^1(x, \mathbf{W}m_1(u)) - f_t^1(x, \mathbf{W}m_2(u)))$$

$$(m_1 - m_2)(du, dx) \leq 0$$

Assumption 4.5 are the graphon game analogies to classic uniqueness assumptions in mean-field games, see for example Carmona & Delarue (2018, Section 3.4) and Lacker (2018, section 8.6).

¹Recall the definition of tightness: for arbitrary index set I and Polish space E , a collection of probability measures $\{P_i\}_{i \in I} \subset \mathcal{P}(E)$ is tight if for any $\epsilon > 0$, there exists some compact measurable subset $K \subset E$ such that $\inf_{i \in I} P_i(K) > 1 - \epsilon$.

Theorem 4.6. Suppose assumption 4.3 and assumption 4.5 holds. Then the graphon game admits a unique W -equilibrium.

The proof follows a standard argument in MFG, see Appendix E.

4.4. Approximate Equilibrium for Finite Player Game

Let $\hat{\pi} : \mathbb{T} \times [0, 1] \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$ be the equilibrium optimal closed-loop policy of the graphon game associated with graphon W , and we construct an n -player game policy from $\hat{\pi}$ as follows. Assign player i the policy

$$\pi^{n, \mathbf{u}^n, i}(t, x_1, \dots, x_n) := \hat{\pi}(t, u_i^n, x_i) \quad (8)$$

and $\boldsymbol{\pi}^{n, \mathbf{u}^n} := (\pi^{n, \mathbf{u}^n, 1}, \dots, \pi^{n, \mathbf{u}^n, n}) \in (\mathcal{A}_n)^n$. Define

$$\epsilon_i^n(\mathbf{u}^n) := \sup_{\beta \in \mathcal{A}_n} J_i(\boldsymbol{\pi}^{n, \mathbf{u}^n, -i}, \beta) - J_i(\boldsymbol{\pi}^{n, \mathbf{u}^n}) \quad (9)$$

and $\boldsymbol{\epsilon}^n(\mathbf{u}^n) := (\epsilon_1^n(\mathbf{u}^n), \dots, \epsilon_n^n(\mathbf{u}^n))$. $\epsilon_i^n(\mathbf{u}^n)$ is the largest reward improvement player i could achieve by changing her own policy, when all other players follow policies $\boldsymbol{\pi}^{n, \mathbf{u}^n}$. By definition 3.1, $\boldsymbol{\pi}^{n, \mathbf{u}^n}$ is an $\epsilon^n(\mathbf{u}^n)$ -equilibrium of the n -player game. We need the following additional assumptions.

Assumption 4.7. 1. $\xi^n \in \mathbb{R}_+^{n \times n}$ is a sequence of matrix with 0 diagonals such that $W_{\xi^n} \rightarrow W$ in strong operator topology, and

$$\lim_{n \rightarrow \infty} \frac{1}{n^3} \sum_{i, j=1}^n (\xi_{ij}^n)^2 = 0 \quad (10)$$

2. For each $t \in \mathbb{T}$, the transition dynamic $P_t : \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \times A \rightarrow \mathcal{P}(\mathbb{R}^d)$ is jointly continuous for all $t \in \mathbb{T}$.

The next main result demonstrates that the n -player game policy $\boldsymbol{\pi}^{n, \mathbf{u}^n}$ constructed from the graphon game equilibrium optimal policy $\hat{\pi}$ forms an approximate equilibrium, and it converges to the true equilibrium in an average sense as the number of players $n \rightarrow \infty$.

Theorem 4.8. Suppose assumption 4.3 and assumption 4.7 holds. For each $n \in \mathbb{N}_+$, let $\mathbf{U}^n := (U_1^n, \dots, U_n^n)$ where $U_i^n \sim \text{unif}(I_i^n)$ and U_i^n is independent of U_j^n for $i \neq j$. Then we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\epsilon_i^n(\mathbf{U}^n)] = 0 \quad (11)$$

The proof is in Appendix F. Equation (11) can be equivalently written as $\epsilon_{I^n}^n(\mathbf{U}^n) \rightarrow 0$ in probability, where $I^n \sim \text{unif}([n])$. Intuitively, for randomly assigned label \mathbf{U}^n , and a player I^n uniformly chosen on $[n]$, the error is small. As the number of player $n \rightarrow \infty$, the collection of \mathbf{U}^n and player label I^n such that the error cannot be controlled becomes a measure 0 set.

Remark 4.9. Equation (10) is a very mild graph denseness condition and is satisfied by many commonly-encountered

275 finite graphs. The assumption $W_{\xi^n} \rightarrow W$ also poses density
 276 restrictions on the underlying graphs of interaction
 277 matrix ξ^n , as the existence of a graphon limit implicitly
 278 implies that the sequence of finite graphs are dense enough.
 279 We give some examples and a detailed discussion on dense
 280 graph sequence in Appendix A.5.
 281

282 5. Learning Scheme and Sample Complexity

283 We now develop a scheme for learning the stationary equi-
 284 librium of infinite-horizon graphon games (Appendix A.1).
 285 Throughout the section we assume finite state space \mathcal{X} and
 286 action space A .
 287

288 5.1. Finite Classes of Label Space

289 To handle the continuous label space algorithmically, one
 290 generally needs function approximation techniques such as
 291 linear function approximation or neural networks, which is
 292 beyond the scope of this work. For the development and
 293 analysis of our algorithms, we discretize the label space
 294 $[0, 1]$ into D classes of types of players $\mathcal{U} \subset [0, 1]$ such that
 295 $|\mathcal{U}| = D < \infty$. We denote $\mathcal{U} := \{u_1, \dots, u_D\}$, and define
 296 projection mapping $\Pi_D : [0, 1] \rightarrow \mathcal{U}$. Denote the inverse
 297 image $I_{u_i} := \Pi_D^{-1}(u_i) \subset [0, 1]$. A simple example is the
 298 uniform quantization: $[0, 1]$ is divided into D bins $\{I_d^D\}_{d=1}^D$,
 299 and Π_D maps each bin to its midpoint:
 300

$$\Pi_D(u) = \sum_{i=1}^D \frac{2i-1}{2D} \mathbf{1}_{\{u \in I_i^D\}} \quad (12)$$

301 As we are only able to learn measures on the finite dis-
 302 cretization \mathcal{U} , we define $\Pi_D : \mathcal{P}(\mathcal{X})^\mathcal{U} \rightarrow \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$
 303 as follows: for any $M = \{M^{u_d}\}_{d=1}^D$, $\Pi_D M$ is the mea-
 304 sure $\text{Leb} \otimes \nu$, where ν is a probabilistic kernel given by
 $\nu(u) := \sum_{d=1}^D M^{u_d} \mathbf{1}_{\{u \in I_{u_d}\}}$.
 305

311 5.2. Approximate Fixed-Point Iteration

312 Our learning scheme follows fixed-point iteration (FPI)
 313 widely used for learning (G)MFGs (Guo et al., 2019; Cui
 314 & Koepll, 2022; Zhang et al., 2023). An FPI represents
 315 an update of the game: given the population measure, the
 316 representative player first finds the optimal policy in reac-
 317 tion to this population, i.e. $\Gamma_1 : \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \rightarrow \mathcal{A}_U$,
 $\Gamma_1(\mu) := \text{argmax}_{\pi \in \mathcal{A}_U} J_W(\mu, \pi)$. As everyone in the pop-
 318 ulation reacts similarly, the population is then updated to the induced state distribution of the acquired policy,
 319 i.e. $\Gamma_2 : \mathcal{A}_U \times \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \rightarrow \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$,
 $\Gamma_2(\pi, \mu) := \mathcal{L}(U, X^\pi)$. Then, the FPI is given by $\Gamma(\mu) := \Gamma_2(\Gamma_1(\mu), \mu)$, and the equilibrium population measure $\hat{\mu}$ satisfies $\hat{\mu} = \Gamma(\hat{\mu})$. However, the FPI operators can be hard to implement. As the environment (P and f) is unknown, Γ_1 and Γ_2 are not directly accessible and need to be approximated. The general approximate FPI scheme is presented in Algorithm 1.

Algorithm 1 provides a general framework that can incorporate various learning algorithms for the two evaluation steps as subroutines, with (i) and (ii) approximating Γ_1 and Γ_2 respectively. If the access to a state process generator (called an oracle) is assumed, we may generate the state variable under any control and population measure for arbitrary times, and Algorithm 1 recovers the algorithms used in prior work (Cui & Koepll, 2022; Zhang et al., 2023).

Algorithm 1 Approximate FPI for GMFGs

```

Initialize policy estimate  $\{\pi_d^0\}_{d=1}^D$  and population esti-  

mate  $\{M_d^0\}_{d=1}^D$  for all label classes  $d \in [D]$   

for  $k \leftarrow 0$  to  $K - 1$  do  

    for  $d \leftarrow 1$  to  $D$  do  

        (i) Evaluate approximate optimal policy  $\pi_d^{k+1}$  in re-  

        action to  $M^k$  (ii) Evaluate approximate population  

        measure  $M_d^{k+1}$  induced by  $\pi_d^{k+1}$   

    end for  

end for  

    Return  $\{\pi_d^K\}$  and  $\{M_d^K\}$ 

```

We next provide the first non-asymptotic analysis for D -class FPI scheme given the following assumptions.

Assumption 5.1. 1. The transition kernel and reward function are uniformly L_P, L_f Lipschitz w.r.t. the measure argument respectively:²

$$\begin{aligned} \sup_{x,a} |f(x, m_1, a) - f(x, m_2, a)| &\leq L_f \|m_1 - m_2\|_{\text{TV}}, \\ \sup_{x,a} \|P(x, m_1, a) - P(x, m_2, a)\|_{\text{TV}} &\leq L_P \|m_1 - m_2\|_{\text{TV}} \end{aligned}$$

2. There exists L_d such that

$$\sup_{u,v \in [0,1]} |W(u, v) - W(\Pi_D(u), v)| \leq L_d/D.$$

3. The FPI operator Γ is a contraction mapping: there exists $\kappa \in (0, 1)$ such that, $\|\Gamma(\mu_1) - \Gamma(\mu_2)\|_{\text{TV}} \leq (1 - \kappa) \|\mu_1 - \mu_2\|_{\text{TV}}, \forall \mu_1, \mu_2$.

5.1(2) ensures label classes \mathcal{U} are a good approximation of the label space $[0, 1]$. An example that satisfies this is the uniform quantization Π_D in equation (12) if the graphon is Lipschitz continuous in the first argument. An example policy operator satisfying 5.1(3) is the softmax operator, with its temperature parameter controlling the Lipschitz constant (Gao & Pavel, 2017). The contraction mapping assumption is limited but unfortunately necessary in complexity analysis proofs. We give a brief discussion on this assumption and different types of fixed-point theorems in Appendix A.6.

Suppose Assumptions 5.1 hold, Algorithm 1 with exact evaluation steps needs at most $D = O(\kappa^{-1} \epsilon^{-1})$ classes and $K = O(\kappa^{-1} \log \epsilon^{-1})$ iterations to achieve an ϵ -approximate equilibrium. This claim is formalized in Theorem 5.4.

²Finite signed measures on finite space can be equivalently expressed as a vector, and the total variation norm is equivalent to the ℓ_1 norm.

5.3. Online Oracle-Free Learning

We now present an online oracle-free subroutine for approximate evaluation steps in Algorithm 1 by specifying a concrete implementation of the two evaluation steps (i) and (ii). Specifically, we use SARSA (Sutton & Barto, 2018), a value-based reinforcement learning method, for policy estimation, and Markov chain Monte Carlo (MCMC) for population estimation.

For policy estimation, we maintain a Q-function: $\mathcal{U} \times \mathcal{X} \times A \rightarrow \mathbb{R}$, with entry $Q_d(x, a)$ estimating the expected return starting with the state-action pair (x, a) conditional on label being u_d . Let \mathcal{Q} be the collection of all Q-functions. To obtain the policy from a Q-function, we assume access to a Lipschitz continuous policy operator $\Gamma_\pi : \mathcal{Q} \rightarrow \mathcal{A}_U$, i.e., for any $Q_1, Q_2 \in \mathcal{Q}$, there exists a constant L_π such that

$$\sup_{u, x} \|(\Gamma_\pi(Q_1) - \Gamma_\pi(Q_2))(u, x)\|_{\text{TV}} \leq L_\pi \|Q_1 - Q_2\|_2 \quad (13)$$

An example policy operator satisfying (13) is the softmax function, with its temperature parameter controlling the constant L_π (Gao & Pavel, 2017). Given Γ_π , SARSA converges to the Q-function corresponding to the optimal policy in $\Gamma_\pi(\mathcal{Q}) \subset \mathcal{A}_U$ (Zou et al., 2019).

Remark 5.2. Utilizing a Lipschitz continuous policy operator, Γ_1 returns the optimal Q-function instead of a policy; and Assumption 5.1(3) can be relaxed to only requiring Γ_1 and Γ_2 to be Lipschitz continuous with constants L_1 and L_2 . Then, we can choose a sufficiently smooth policy operator such that $L_\pi L_1 L_2 < 1$, making the FPI operator $\Gamma(\mu) := \Gamma_2(\Gamma_\pi(\Gamma_1(\mu)), \mu)$ contractive.

For population estimation, we maintain a M-function: $\mathcal{U} \rightarrow \mathcal{P}(\mathcal{X})$, with entry M_d estimating the population measure of the representative player conditional on label u_d . Since $\mathcal{U} \times \mathcal{X} \times A$ is a finite space, both Q- and M-functions can be represented by tables. Being fully online, SARSA and MCMC can update the Q- and M-functions using the same online samples without the need of any oracle. Specifically, we execute H updates for the evaluation subroutine of Algorithm 1. At each step $\tau = 0, \dots, H-1$, the representative agent with label u_d at x_τ samples its action $a_\tau \sim \Gamma_\pi(Q_d^{k, \tau})$, reward $r_\tau = f(x_\tau, \mathbf{W}\Pi_D M^{k, 0}(u_d), a_\tau)$, next state $x_{\tau+1} \sim P(x_\tau, \mathbf{W}\Pi_D M^{k, 0}(u_d), a_\tau)$, and next action $a_{\tau+1} \sim \Gamma_\pi(Q_d^{k, \tau})$. Using these sample, the Q- and M-functions are updated simultaneously as follows:

$$\begin{aligned} Q_d^{k, \tau+1}(x_\tau, a_\tau) &\leftarrow (1 - \alpha_\tau) Q_d^{k, \tau}(x_\tau, a_\tau) \\ &\quad + \alpha_\tau \left(r_\tau + \gamma Q_d^{k, \tau}(x_{\tau+1}, a_{\tau+1}) \right), \\ M_d^{k, \tau+1} &\leftarrow (1 - \beta_\tau) M_d^{k, \tau} + \beta_\tau \delta_{x_{\tau+1}}, \end{aligned} \quad (14)$$

where the Q- and M-functions are indexed by the outer iteration k and the inner evaluation step t , and α_τ and β_τ are step sizes. Substituting the Q-function $Q_d^{k, \tau}$ with the optimal Q-function $Q^{\mu^{k, \tau}}$ associated with the population

measure $\mu^{k, \tau} = \Pi_D M^{k, \tau}$, we recover the FPI scheme in Algorithm 1. Substituting (i) and (ii) in Algorithm 1 with H updates using Equation (14), we obtain the first fully online algorithm for learning GMFGs. Notably, our method is oracle-free in the sense that we do not assume access to an optimal policy calculator or a state process generator. Additionally, in contrast to FPI-like methods in prior work where (i) and (ii) in Algorithm 1 are executed sequentially, Equation (14) updates both policy and population concurrently using the same samples, enhancing the sample efficiency. Algorithm 2 is an example of concrete realization of the aforementioned ideas.

Finally, we give the sample complexity of our method. As our method is fully online, we need the following ergodicity assumption (Zou et al., 2019).

Assumption 5.3. For any $\pi \in \Gamma_\pi(\mathcal{Q})$ and $M \in \mathcal{P}(\mathcal{X})^\mathcal{U}$, the Markovian state dynamic is ergodic: there exists $\mu \in \mathcal{P}_{\text{unif}}(\mathcal{U} \times \mathcal{X})$ and $c_1 > 0, c_2 \in (0, 1)$ such that

$$\sup_x \|\mathbb{P}(X_\tau \in \cdot | X_0 = x) - \mu\|_{\text{TV}} \leq c_1 c_2^\tau,$$

where the dynamic of X is determined by policy π and neighborhood measure $\mathbf{W}\Pi_D M$.

Theorem 5.4. Let $\hat{\mu}$ be the stationary equilibrium measure of the infinite horizon GMFG. Suppose Assumptions 5.1 and 5.3 hold. For any initial estimate $M^{0, 0} \in \mathcal{P}(\mathcal{X})^\mathcal{U}$, Algorithm 1, combined with Equation (14) and step sizes $\alpha_\tau, \beta_\tau \asymp 1/\tau$, finds an ϵ -approximate equilibrium distribution $M^{K, H}$ such that $\mathbb{E}\|\Pi_D M^{K, H} - \hat{\mu}\|_{\text{TV}} \leq \epsilon$, with the number of iteration being at most

$$K = O(\kappa^{-1} \log \epsilon^{-1}), \quad D = O(\kappa^{-1} \epsilon^{-1}), \\ H = O(\kappa^{-3} \epsilon^{-3} \log \epsilon^{-1}),$$

giving a total sample complexity of $O(\kappa^{-5} \epsilon^{-4} \log^2 \epsilon^{-1})$.

The proof of Theorem 5.4 and more details of our method are deferred to Appendix G.

6. Numerical Experiments

In this section, we apply our learning algorithm to three graphon game examples, namely, Flocking-, SIS- and Invest-Graphon. We first briefly introduce each game scenario, and present only the algorithm performance and GMFE for Flocking-Graphon due to space limit. The problem formulation of each game is in Appendix H. The detailed numerical results are in I, including algorithm performance (e.g., exploitability, convergence) and visualizations for GMFE. All numerical experiments are conducted on Mac Air M2.

Flocking-Graphon The flocking-graphon game (Lacker & Soret, 2022) studies the flocking behavior in a system where each agent makes decisions on its velocity which in turn determines its position. We consider the game with $\mathcal{X} = [0, 1]$, and a time horizon $\mathcal{T} = [0, 1]$, under proper discretization. The policy is in the form $\pi_t(u, x) \equiv \delta_{\alpha_t(u, x)}$,

where $\alpha_t(u, x)$ is the velocity of the agent conditional on label being u and position being x at time t . Each agent aims to minimize its own running cost determined by the velocity control and the agent's deviation from the population.

SIS-Graphon The SIS-Graphon game (Cui & Koepll, 2022) models an epidemic scenario where agents can choose take precautions to avoid being infected. The infected probability is determined by the agents' action (i.e., take precaution or not) and the number of infected neighbours.

Invest-Graphon In the Invest-Graphon game model (Cui & Koepll, 2022), each firm aims to maximize its own profit, which is determined by the firm's investment strategies and other firms' product quality.

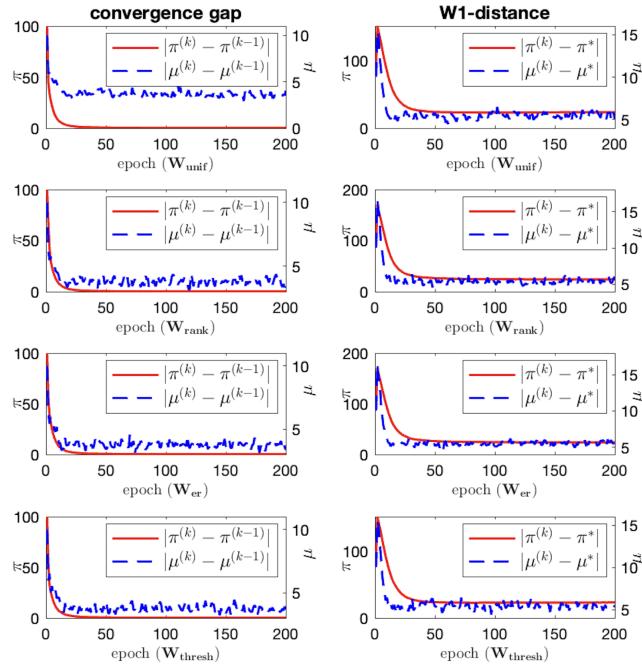


Figure 1. Algorithm performance (Flocking-Graphon)

We test four types of graphons: uniform attachment graphon ($W_{\text{unif}}(u, v) = 1 - \max(u, v)$), ranked attachment graphon ($W_{\text{rank}}(u, v) = 1 - uv$), Erdős-Rényi graphon ($W_{\text{er}}(u, v) = p$) and threshold graphon ($W_{\text{thresh}}(u, v) = \mathbf{1}_{u+v < 1}$). The details of each graphon is in Appendix A.5. Figure 1 demonstrates the algorithm performance to solve the Flock-Graphon. The x-axis denotes the epoch index k . We visualize the convergence gaps $|\mu^{(k)} - \mu^{(k-1)}|, |\pi^{(k)} - \pi^{(k-1)}|$, and the W1-distances $|\mu^{(k)} - \mu^*|, |\pi^{(k)} - \pi^*|$, which measures the closeness between the benchmark solution (π^*, μ^*) and results at each epoch. The benchmark solution is obtained by the equivalent class method (Cui & Koepll, 2022). The results show that it takes around 50 epochs for our algorithm to converge. The convergence performance remains consistent for all four graphons.

Figure 2 shows the obtained GMFE for Flocking-Graphon.

We visualize the policy and state density of agent with label $U = 1$ at equilibrium in a 3D plot. The x-axis denotes the space domain \mathcal{X} , and the y-axis is the time horizon \mathcal{T} . Agent with each label is initialized at $t = 0$ uniformly over \mathcal{X} . Note that the GMFE is time-dependent. We adapt our learning algorithm to solve GMFGs with finite horizons (See Algo 3 in Appendix G). The z-axis is the spatial-temporal velocity control $\alpha_t(1, x)$ and population density $\mu_t(1, x)$ of the agent with label 1. The numerical results show that GMFEs associated with W_{unif} and W_{thresh} are similar. The flock behavior occurs when agents gather together at position $x = 0.6$ and the population density μ reaches a red peak around 0.35 with velocity around 0.2. When the agent's velocity reaches the maximum velocity $\alpha_{\max} = 1$ (dark red), the population quickly dissipates (dark blue) and no flock behavior occurs.

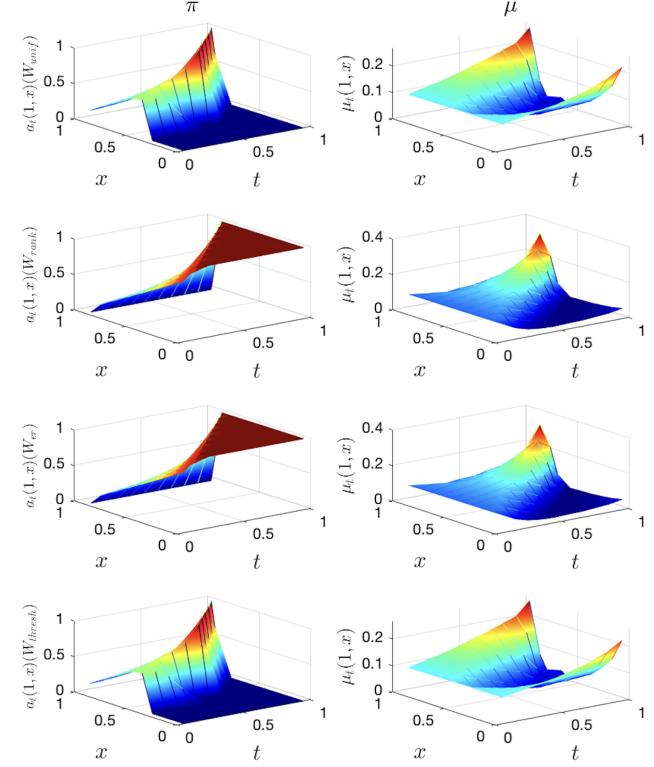


Figure 2. GMFE (Flocking-Graphon)

7. Conclusion

We offered a new general formulation of graphon games with one representative player in continuous state and action space. We gave a comprehensive analysis on the equilibrium properties with assumptions milder than previous works. We present a general approximate fixed-point iteration framework, and designed an oracle-free algorithm along with the sample complexity analysis.

440 Broader Impact

441 This work is motivated by the theoretical challenges in the
 442 analysis of graphon games. As a generalization to mean-
 443 field games, graphon games is capable of modeling hetero-
 444 geneous interactions among gaming participants, and this
 445 flexibility allows it to cover a broader range of applications
 446 in finance, economics, engineering and AI, including for
 447 example high-frequency trading, social opinion dynamics
 448 and autonomous vehicle driving. By addressing rigorously
 449 the technical issues faced by games on networks, this work
 450 proposes a conceptually and mathematically concise formu-
 451 lation. The analysis provides concrete theoretical founda-
 452 tion for the mathematical properties, on top of which the
 453 algorithms empower the solvability of the system. With
 454 the comprehensive and self-consistent analysis, this work is
 455 capable of modeling system of large amount of agents and
 456 remain computationally efficient.

458 References

460 Aliprantis, C. D. and Border, K. C. *Infinite Dimensional*
 461 *Analysis*. Springer, 3rd edition, 2006.

463 Anonymous, A. A single online agent can efficiently learn
 464 mean field games, 2024.

466 Aurell, A., Carmona, R., and Lauriere, M. Stochastic
 467 graphon games: II. the linear-quadratic case. *Applied*
 468 *Mathematics and Optimization*, 85, 06 2022.

470 Billingsley, P. *Probability and Measure*. John Wiley and
 471 Sons, 3rd edition, 1995.

473 Bris, P. L. and Poquet, C. A note on uniform in time mean-
 474 field limit in graphs, 2023.

476 Brunick, G. and Shreve, S. Mimicking an ito process by a
 477 solution of a stochastic differential equiation. *The Annals*
 478 *of Applied Probability*, pp. 1584–1628, 2013.

480 Caines, P. E. and Huang, M. Graphon mean field games and
 481 their equations. *SIAM Journal on Control and Optimiza-*
 482 *tion*, 59(6):4373–4399, 2021.

484 Cardaliaguet, P. and Hadikhanloo, S. Learning in mean field
 485 games: the fictitious play. *ESAIM: Control, Optimisation*
 486 *and Calculus of Variations*, 23, 07 2015. doi: 10.1051/
 487 cocv/2016004.

488 Carmona, R. and Delarue, F. *Probabilistic Theory of Mean*
 489 *Field Games with Applications I*. Springer, 2018.

491 Carmona, R., Cooney, D. B., Graves, C. V., and Laurière, M.
 492 Stochastic graphon games: I. the static case. *Mathematics*
 493 *of Operations Research*, 47(1):750–778, 2021.

Chen, X., Liu, S., and Di, X. A hybrid framework of rein-
 force learning and physics-informed deep learning
 for spatiotemporal mean field games. In *Proceedings*
 of the 2023 International Conference on Autonomous
 Agents and Multiagent Systems, pp. 1079–1087, 2023a.

Chen, X., Liu, S., and Di, X. Learning dual mean field
 games on graphs. In *Proceedings of the 26th European*
Conference on Artificial Intelligence, ECAI, 2023b.

Cui, K. and Koepll, H. Approximately solving mean field
 games via entropy-regularized deep reinforcement learn-
 ing. In *International Conference on Artificial Intelligence*
and Statistics, 2021.

Cui, K. and Koepll, H. Learning graphon mean field games
 and approximate nash equilibria. In *International Con-*
ference on Learning Representations, 2022.

Cui, K., KhudaBukhsh, W. R., and Koepll, H. Hyper-
 graphon mean field games. *Chaos: An Interdisciplinary*
Journal of Nonlinear Science, 32(11), 2022.

Delattre, S., Giacomin, G., and Luçon, E. A note on dy-
 namical models on random graphs and fokker-planck
 equations. *Publications of the Mathematical Institute of*
the Hungarian Academy of Sciences, 165:785–798, 2016.

Elie, R., P’erolat, J., Laurière, M., Geist, M., and Pietquin,
 O. On the convergence of model free learning in mean
 field games. In *AAAI*, 2020.

Erdős, P. and Rényi, A. On random graphs. i. *Publicationes*
Mathematicae, 6(3–4):290–297, 1959.

Erdős, P. and Rényi, A. On the evolution of random graphs.
Publicationes of the Mathematical Institute of the Hungar-
ian Academy of Sciences, 5:17–61, 1960.

Fabian, C., Cui, K., and Koepll, H. Learning sparse graphon
 mean field games. In *International Conference on Artifi-*
cial Intelligence and Statistics, pp. 4486–4514. PMLR,
 2023.

Gao, B. and Pavel, L. On the properties of the softmax func-
 tion with application in game theory and reinforcement
 learning. *arXiv preprint arXiv:1704.00805*, 2017.

Gao, S., Tchuendom, R. F., and Caines, P. E. Linear
 quadratic graphon field games. *Communications in Infor-*
mation and Systems, 21(3):341–369, 06 2021.

Guo, X., Hu, A., Xu, R., and Zhang, J. Learning mean-
 field games. *Advances in Neural Information Processing*
Systems, 32, 2019.

Huang, K., Di, X., Du, Q., and Chen, X. A game-theoretic
 framework for autonomous vehicles velocity control:
 Bridging microscopic differential games and macroscopic

- 495 mean field games. *Discrete and Continuous Dynamical*
 496 *Systems - Series B*, 25(12):4869–4903, 2020.
- 497
- 498 Huang, M., Malhamé, R. P., and Caines, P. E. Large pop-
 499 ulation stochastic dynamic games: closed-loop McKean-
 500 Vlasov systems and the Nash certainty equivalence prin-
 501 ciple. *Communications in Information & Systems*, 6(3):
 502 221–252, 2006.
- 503 Jabin, P.-E., Poyato, D., and Soler, J. Mean-field limit of
 504 non-exchangeable systems, 2022.
- 505
- 506 Lacker, D. Mean field games via controlled martingale
 507 problems: Existence of markovian equilibria. *Stochas-*
 508 *tic Processes and their Applications*, 23(4):2856–2894,
 509 2015.
- 510
- 511 Lacker, D. Mean field games and interacting particle sys-
 512 tems. *preprint*, 2018.
- 513
- 514 Lacker, D. and Soret, A. A label-state formulation of
 515 stochastic graphon games and approximate equilibria on
 516 large networks. *Mathematics of Operations Research*,
 517 2022.
- 518
- 519 Lacker, D., Ramanan, K., and Wu, R. Local weak conver-
 520 gence for sparse networks of interacting processes. *The*
 521 *Annals of Applied Probability*, 33(2):843 – 888, 2023.
- 522
- 523 Lasry, J.-M. and Lions, P.-L. Mean field games. *Japanese*
 524 *Journal of Mathematics*, 2(1):229–260, 2007.
- 525
- 526 Lauriere, M., Perrin, S., Girgin, S., Muller, P., Jain, A., Ca-
 527 bannes, T., Piliouras, G., Perolat, J., Elie, R., Pietquin, O.,
 528 and Geist, M. Scalable deep reinforcement learning algo-
 529 rithms for mean field games. In *Proceedings of the 39th*
 530 *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp.
 531 12078–12095. PMLR, 2022.
- 532
- 533 Littman, M. L. Markov Games As a Framework for Multi-
 534 agent Reinforcement Learning. In *Proceedings of the*
 535 *Eleventh International Conference on International Confer-*
 536 *ence on Machine Learning*, ICML’94, pp. 157–163,
 537 San Francisco, CA, USA, 1994. Morgan Kaufmann Pub-
 538 lishers Inc. ISBN 978-1-55860-335-6. event-place: New
 539 Brunswick, NJ, USA.
- 540
- 541 Lovász, L. *Large networks and graph limits*, volume 60.
 542 American Mathematical Soc., 2012.
- 543
- 544 Mitrophanov, A. Y. Sensitivity and convergence of uni-
 545 formly ergodic markov chains. *Journal of Applied Prob-*
 546 *ability*, 42(4):1003–1014, 2005.
- 547
- 548 Perrin, S., Perolat, J., Laurière, M., Geist, M., Elie, R., and
 549 Pietquin, O. Fictitious play for mean field games: Con-
 tinuous time analysis and applications. In *Proceedings of*
- the 34th International Conference on Neural Information Processing Systems, NIPS’20, 2020.
- Perrin, S., Laurière, M., Pérolat, J., Geist, M., Élie, R., and Pietquin, O. Mean field games flock! the reinforcement learning way. *arXiv preprint arXiv:2105.07933*, 2021.
- Perrin, S., Laurière, M., Pérolat, J., Élie, R., Geist, M., and Pietquin, O. Generalization in mean field games by learning master policies. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 9413–9421, 2022.
- Solan, E. and Vieille, N. Stochastic games. *Proceedings of the National Academy of Sciences*, 112(45):13743–13746, 2015.
- Sun, Y. The exact law of large numbers via fubini extension and characterization of insurable risks. *Journal of Economic Theory*, 126:31–69, 2006.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Tangpi, L. and Zhou, X. Optimal investment in a large pop-
 ulation of competitive and heterogeneous agents, 2023.
- Yang, J., Ye, X., Trivedi, R., Xu, H., and Zha, H. Deep mean
 field games for learning optimal behavior policy of large
 populations. In *International Conference on Learning Representations*, 2018a.
- Yang, Y., Luo, R., Li, M., Zhou, M., Zhang, W., and Wang,
 J. Mean Field Multi-Agent Reinforcement Learning. In
International Conference on Machine Learning, pp. 5571–
 5580, July 2018b.
- Zhang, F., Tan, V. Y., Wang, Z., and Yang, Z. Learning
 regularized monotone graphon mean-field games. *Neural Information Processing Systems*, 2023.
- Zou, S., Xu, T., and Liang, Y. Finite-sample analysis for
 sarsa with linear function approximation. *Advances in neural information processing systems*, 32, 2019.

Organization of Appendix

The appendix is outlined as follows.

Appendix A is a discussion section serving as a supplement to the concepts in the main paper. The topics include: the infinite horizon version of graphon game formulation (Appendix A.1), formulating the graphon game into a mean-field game with augmented state space (Appendix A.2), the degeneration of graphon game to mean-field games with trivial graphon (Appendix A.3), time-variant interaction intensities (Appendix A.4), dense graph sequence and examples (Appendix A.5), fixed point theorems and the contraction mapping assumption (Appendix A.6).

In Appendix B, we define the continuum-player formulation (Appendix B.1) and compare it with our representative-player formulation (Appendix B.2). In particular, we discuss in detail the aforementioned measurability issue residing in continuum-player formulation. We then give a toy example in Appendix C to demonstrate the difference on the two formulations.

The following three appendix are dedicated to the proof of analysis properties. The existence of equilibrium (Theorem 4.4) is proved in Appendix D. The uniqueness of equilibrium (Theorem 4.6) is proved in Appendix E. The approximate equilibrium (Theorem 4.8) is proved in Appendix F.

The sample complexity of learning algorithms (Theorem 5.4) are proved in Appendix G. Finally, we give the detailed problem setups for the numerical examples in Appendix H, and show the numerical results in Appendix I.

A. Additional Discussion

A.1. Infinite horizon formulation

In this section we define the infinite horizon version of the representative-player graphon game, as appose the the finite horizon version defined in section 4.1. All analysis results in section 4 regarding existence, uniqueness and approximate equilibrium holds by adjusting the assumptions accordingly.

Let the graphon $W \in L^1_+[0, 1]^2$ be given and fixed. Let $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ be a filtered probability space that support an \mathcal{F}_0 -measurable random variable U uniform on $[0, 1]$, and a Markov process X valued in \mathbb{R}^d . We understand U as the label for the representative player, and X as her state dynamic. Let the flow of label-state joint measures be $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$, where the path space $\mathcal{C} = \prod_{i=0}^{\infty} \mathbb{R}^d$ is now a countable product of \mathbb{R}^d . $\mu_t \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d)$ is the marginal under image $(u, x) \mapsto (u, x_t)$. Let the initial joint law $\lambda := \mathcal{L}(U, X_0) \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d)$ be given.

We still let \mathcal{A}_U denotes the collection of *time-variant* closed-loop (Markovian) policies $\mathbb{N}_+ \times [0, 1] \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$. For any $\pi \in \mathcal{A}_U$, (U, X) follows the transition dynamic

$$(U, X_0) \sim \lambda \\ a_t \sim \pi_t(U, X_t) \quad X_{t+1} \sim P(X_t, W\mu_t(U), a_t)$$

note that the transition law P is time-invariant. Let $f : \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \times A \rightarrow \mathbb{R}$ be the running reward and $\gamma \in (0, 1)$ be a known discount factor. The objective of the representative player is to choose $\pi \in \mathcal{A}_U$ to maximize

$$J_W(\mu, \pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t f(X_t^\pi, W\mu_t(U), a_t) \right]$$

Definition A.1. We say that $(\hat{\mu}, \hat{\pi}) \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C}) \times \mathcal{A}_U$ is a W -equilibrium if

$$J_W(\hat{\mu}, \hat{\pi}) = \sup_{\pi \in \mathcal{V}_U} J_W(\hat{\mu}, \pi) \\ \hat{\mu} = \mathcal{L}(U, X^{\hat{\pi}})$$

If we do not fix an initial distribution λ , we may define a stationary equilibrium which is time independent:

Definition A.2. We say that $(\hat{\mu}, \hat{\pi}) \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathbb{R}^d) \times \mathcal{A}_U$ is a stationary W -equilibrium if

$$J_W(\hat{\mu}, \hat{\pi}) = \sup_{\pi \in \hat{\mathcal{A}}_U} J_W(\hat{\mu}, \pi) \\ \hat{\mu} = \mathcal{L}(U, X_t^{\hat{\pi}}) \quad \forall t \geq 0$$

where \mathcal{A}_U now denotes the collection of time-invariant closed-loop policies $[0, 1] \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$.

Note that we need an additional ergodicity assumption of the Markov chain to show the existence of stationary W -equilibrium with the same proof argument in Appendix D, i.e. the Markov chain is ergodic (admits a limiting distribution) under any policy. This is made formal in Assumption 5.3.

A.2. Game with Augmented State Space

We give another view of graphon game by recasting it into a mean-field game with augmented state space. We view the label U as a coordinate of the state, and it remains at the same value a.s. Let $\bar{X} = \begin{pmatrix} U \\ X \end{pmatrix} \in \mathbb{R}^{d+1}$, where the state process space is augmented by one more dimension. Any fixed $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$ can now be regarded as an element in $\mathcal{P}(\mathcal{C}^{d+1})$ where $\mathcal{C}^{d+1} = (\mathbb{R}^{d+1})^{T+1}$ is the augmented path space. Given any graphon closed-loop policy $\pi \in \mathcal{A}_U$, for every $\bar{x} = \begin{pmatrix} u \\ x \end{pmatrix} \in \mathbb{R}^{d+1}$, define a mean-field closed-loop policy $\bar{\pi}$ and a mean-field Markovian transition law \bar{P} by

$$\begin{aligned}\bar{\pi}_t(\bar{x})(da) &:= \pi_t(u, x)(da) \\ \bar{P}_t(\bar{x}, m, a)(d\bar{y}) &:= \delta_u(dv) P_t(x, Wm(u), a)(dy) \quad \forall \bar{y} = \begin{pmatrix} v \\ y \end{pmatrix}\end{aligned}$$

respectively and let $\bar{\lambda}(d\bar{y}) = \lambda(dv, dy)$ for any $\bar{y} = \begin{pmatrix} v \\ y \end{pmatrix}$ be the mean-field initial condition. Then \bar{X} satisfies the dynamic

$$\begin{aligned}\bar{X}_0 &\sim \bar{\lambda} \\ a_t &\sim \bar{\pi}_t(\bar{X}_t) \quad \bar{X}_{t+1} \sim \bar{P}_t(\bar{x}, \mu_t, a_t)\end{aligned}$$

Define similarly for every $\bar{x} = \begin{pmatrix} u \\ x \end{pmatrix} \in \mathbb{R}^{d+1}$ the reward functions

$$\begin{aligned}\bar{f}_t(\bar{x}, m, a) &:= f_t(x, Wm(u), a) \\ \bar{g}(\bar{x}) &:= g(x, W\mu_t(u))\end{aligned}$$

for all $t \in \mathbb{T}$. The objective is recast into

$$J(\bar{\pi}) := \mathbb{E} \left[\sum_{t \in \mathbb{T}} \bar{f}_t(\bar{X}_t^{\bar{\pi}}, \mu_t, a_t) + \bar{g}(\bar{X}_T^{\bar{\pi}}, \mu_T) \right]$$

thus we have obtained a classic mean-field game problem associated with the new coefficients $\bar{\lambda}, \bar{P}_t, \bar{f}_t, \bar{g}$. Note that their implicit dependence on W .

However, it is worth noticing that in most of the proofs for graphon game, this translation into mean-field game with augmented state space does not simplify the mathematical analysis, and it is not appropriate to adapt the mean field game results directly. There are two main reasons (Lacker & Soret, 2022):

Firstly, it requires the graphon $W \in L_+^1[0, 1]^2$ to be continuous. To see this, recall that most of the results for classic mean-field games assume the joint continuity of reward function, see e.g. (Carmona & Delarue, 2018; Lacker, 2018). In particular, $\bar{f}_t(\bar{x}, m, a) := f_t(x, Wm(u), a)$ is assumed to be continuous in the augmented state variable (u, x) . This requires the graphon operator $W\mu$ viewed as a function

$$[0, 1] \ni u \mapsto W\mu(u) \in \mathcal{M}_+(E)$$

should be continuous, for any $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times E)$, which is achieved by a continuous graphon. However, the graphon is in general not a continuous function. Indeed many commonly encountered convergent graph sequence tends to a discontinuous graphon limit, see for instance examples in (Lovász, 2012) section 11.4.

Second, in the analysis of approximate equilibrium, the model setting for the finite-player game does not fit into this augmented state space framework. Consider an n -player game associated with interaction matrix $\xi \in \mathbb{R}_+^{n \times n}$, and assign player i the label $u_i \in I_i^n$. Recall the setting for finite player games in section 3, the running reward can be written as $f_t(X_t^i, W_\xi M(u_i), a_t^i)$, where M is the empirical label-state measure defined in (4). On the other hand, let $\bar{X}^i = \begin{pmatrix} u_i \\ X^i \end{pmatrix}$, and the running cost of player i in the aforementioned augmented state space framework is

$$\bar{f}_t(\bar{X}^i, M, a_t^i) := f_t(X^i, WM(u_i), a_t^i)$$

which is different from the original problem, as in the finite player game the graphon W needs to be replaced with the step graphon W_ξ . However, it is not possible to incorporate this change in the augmented state space framework. As a result the augmented state space transformation fails to provide an approximate equilibrium result, which is a strong justification of the reasonableness of graphon game formulation.

In continuous time setting (Lacker & Soret, 2022), the augmented state space formulation provides an equivalent forward-backward PDE system for the graphon game, and thus provides another perspective to the problem formulation.

Actually the continuum-player graphon games may be transformed to a mean-field game with augmented state space in a similar way, and many previous works on continuum-player formulation relied on this (Cui & Koeppl, 2022; Zhang et al., 2023) to show existence of equilibrium. However, they not only suffer from the two limitations mentioned above, but also encounter a critical measurability issue that representative-player formulation does not have, and this leads to difficulties in the proof. We will discuss this point in detail in Appendix B.

A.3. Degeneration to mean-field games under trivial graphon

When the graphon $W \equiv 1$, the interactions among players are symmetric, and we illustrate that our graphon game formulation degenerates to the classic mean-field game.

Let the population measure $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$ be given and take the product measure form: $\mu(du, dx) = du \times \nu(dx)$ for some $\nu \in \mathcal{P}(\mathcal{C})$, and let the initial distribution λ be a product measure with the path space marginal λ° . The graphon operator applied on μ degenerates to ν :

$$W\mu(u) = \int_{[0,1] \times \mathcal{C}} \delta_x m(dv, dx) = \int_{\mathcal{C}} \delta_x \nu(dx) \quad \forall u \in [0, 1]$$

Controlled by a closed-loop policy $\pi \in \mathcal{A}_U$ that depends only on the state variable, (U, X) follows the transition dynamic: $U \sim \text{unif}[0, 1]$, $X_0 \sim \lambda^\circ$ and

$$a_t \sim \pi_t \quad X_{t+1} \sim P_t(X_t, \nu_t, a_t)$$

note that U and X are now independent. The objective of the representative player becomes

$$J_W(\nu, \pi) = \mathbb{E} \left[\sum_{t \in \mathbb{T}} f_t(X_t^\pi, \nu_t, a_t) + g(X_T^\pi, \nu_T) \right]$$

where the expectation is w.r.t X only, thanks to the independence between U and X . In this way our label-state graphon game formulation degenerates to a classic mean-field game problem, and the equilibrium measure and controls indeed does not depend on the label.

A.4. Time-variant interaction intensity

It is possible to consider time-variant interaction intensity in our framework when the time horizon is finite. In definition of finite player games (section 3), we may replace ξ with a sequence of matrix $\{\xi^t\}_{t=0}^T$, where ξ^t is the interaction intensity of the n players at time t . The empirical weighted neighborhood measure of player i then becomes $M_t^i = \frac{1}{n} \sum_{j=1}^n \xi_{ij}^t \delta_{X_t^j}$, and it can be equivalently written as $W_{\xi^t} S_t(u_i)$, in the notation of section 3.

In the graphon game setting (section 4.1), we may work on a sequence of graphon $\{W_t\}_{t=0}^T$, where W_t is the interaction among a continuum types of players at time t . Note that the sequence $\{W_t\}_{t=0}^T$ should be non-random. By replacing every $W\mu_t$ with $W_t\mu_t$, it is ready to check that the existence (section 4.2) and uniqueness (section 4.3) results still holds. As for the approximate equilibrium result (section 4.4), we may change assumption 4.7(1) into the following: $W_{\xi^{n,t}} \rightarrow W_t$, and

$$\lim_{n \rightarrow \infty} \frac{1}{n^3} \sum_{i,j=1}^n (\xi_{ij}^{n,t})^2 = 0$$

for every $t = 0, \dots, T$. Then, the approximate equilibrium result still holds.

We present the main paper in terms of a time-invariant graphon W to avoid distraction from the main point we want to address.

715 A.5. Graph sequence and the convergence assumption 4.7

 716 Conceptually, a graph is dense if nearly every pair of vertices are connected by an edge. However rigorously, the denseness
 717 of graph is ill-defined, and different results require different denseness conditions.
 718

 719 We first demonstrate that equation (10) is indeed very mild. We may write $\text{Tr}((\xi^n)^2) = \sum_{i,j=1}^n (\xi_{ij}^n)^2$ where $\text{Tr}(\cdot)$ is
 720 the trace, and this is referred as second moment of square matrix. Here are several examples on commonly-encountered
 721 interaction matrix on networks.
 722

 723 **Complete graph.** Let $\xi_{ij}^n = 1$ for each $i \neq j$, and thus ξ^n is the adjacency matrix of a complete graph, and this recovers the
 724 mean-field case where the players interact symmetrically. We have $W_{\xi^n} \equiv 1$ for all n , and thus $W_{\xi^n} \rightarrow W$ for $W \equiv 1$. We
 725 have
 726

$$\frac{1}{n^3} \sum_{i,j=1}^n (\xi_{ij}^n)^2 \leq \frac{1}{n} \rightarrow 0$$
 727
 728

 729 **Threshold graph.** Consider a threshold graph on n vertices where vertex i and j are connected by an edge if $i + j < n$, and
 730 let $\xi_{ij}^n = 1_{i+j < n}$. It is easy to see that W_{ξ^n} converges in cut norm to a limit defined by $W(u, v) := 1_{u+v < 1}$. It is ready to
 731 check that assumption equation (10) is satisfied.
 732

 733 **Random walk on graph.** Consider a graph on n vertices where vertex i has degree d_i^n . Let $\xi_{ij}^n = \frac{1}{d_i^n} 1_{i \sim j}$, where $1_{i \sim j}$ is 1 if
 734 i and j are connected by an edge and 0 otherwise. Then ξ/n is a Markovian transition matrix of the random walk on the
 735 graph. We have
 736

$$\frac{1}{n^3} \sum_{i,j=1}^n (\xi_{ij}^n)^2 = \frac{1}{n} \sum_{i,j=1}^n \frac{1}{(d_i^n)^2} 1_{i \sim j} = \frac{1}{n} \sum_{i=1}^n \frac{1}{d_i^n}$$
 737
 738
 739

 740 and the assumption holds if $\sum_{i=1}^n \frac{1}{d_i^n} \rightarrow 0$, and intuitively this means the average of degrees diverge. In particular, if
 741 $d_1^n = \dots = d_n^n = d^n$, ξ^n becomes an interaction matrix on a d^n -regular graph, and we just need $\frac{1}{d^n} \rightarrow 0$, i.e. the degree d^n
 742 diverges to satisfies equation (10). However, not even every sequence of regular graphs has a graphon limit, and we will
 743 discuss this below.
 744

 745 **Erdős-Rényi graph.** Consider an Erdős-Rényi graph $G^n(p_n)$ (Erdős & Rényi, 1959) on n vertices, where every edge is
 746 connected with Bernoulli(p_n). Let $\xi_{ij}^n = \frac{1}{p_n} 1_{i \sim j}$, it is not hard to show that equation (10) holds in probability as long as
 747 $np_n \rightarrow \infty$. We understand np_n as the expected degree of any vertex, and this is an important quantity of Erdős-Rényi
 748 graphs that also implies connectivity (Erdős & Rényi, 1960). Moreover, when $p_n \rightarrow p$ for some $p \in (0, 1)$, $W_{\xi^n} \rightarrow W$ for
 749 $W \equiv 1$ in probability.
 750

 750 All examples mentioned above merely requires a diverging-average-degree type condition to be considered dense enough
 751 for our results to hold. These denseness conditions are attracting more and more awareness in the stochastic community and
 752 particularly in the studies of stochastic differential equation dynamics and heterogenous propagation of chaos on networks
 753 (Delattre et al., 2016; Bris & Poquet, 2023; Jabin et al., 2022).
 754

 755 The assumption $W_{\xi^n} \rightarrow W$ is also a denseness condition as the existence of a graphon limit implicitly implies that each
 756 graph in the converging sequence is dense. Actually in the sparse setting, vertices in a local neighborhood interact strongly
 757 with each other and do not become negligible as the number of vertices goes to infinity (Lacker et al., 2023). The propagation
 758 of chaos results also fail in this regime. Nevertheless, not every dense graph sequence necessarily admits a graphon limit,
 759 since the sequence is also required to preserve similar network structures. This can be formalized by graph homomorphism
 760 [(Lovász, 2012) chapter 5].
 761

 761 It is worth noticing that a sequence of sparse graphs may converge if they are sampled from the limiting graphon. There are
 762 works (Fabian et al., 2023) adopting this setting. However conceptually this means the finite player games are constructed
 763 from the graphon game, which is different from the view we take that graphon games are motivated by finite player games.
 764

 765 Finally, we demonstrate that the convergence of graphon in strong operator topology is weaker than converging in other
 766 norm. Recall the definition of integral operator in equation (2):
 767

$$W\phi(u) := \int_{[0,1]} W(u, v)\phi(v)dv \quad \forall \phi \in L_\infty[0, 1]$$
 768
 769

which maps $L_\infty[0, 1]$ to $L_1[0, 1]$. The integral operator norm is given by $\|W\|_{\infty \rightarrow 1} := \sup_{\|\phi\|_\infty \leq 1} \|W\phi\|_1$, where $\|\cdot\|_p$ is the L_p norm. It is known to be equivalent to the cut norm (Lovász, 2012) lemma 8.11 by $\|W\|_{\infty \rightarrow 1} \leq \|W\|_\square \leq 4\|W\|_{\infty \rightarrow 1}$, where the cut norm of a graphon is defined by

$$\|W\|_\square := \sup_{S, T \subset [0, 1]} \left| \int_{S \times T} W(u, v) dudv \right|$$

for measurable subsets S, T .

We may see the convergence in strong operator topology $W^n \rightarrow W$ as the strong convergence of the integral operator. As the operator norm convergence of bounded linear operators implies the strong convergence, we note that convergence in strong operator topology can be implied by converge in the cut norm. Indeed W^n converging to W in L_1 also implies $W^n \rightarrow W$, see Lemma D.3.

A.6. Fixed point theorems and contraction mapping assumption

There are two main stream fixed point theorems. The first type is based on contraction mapping, that if an operator is contraction in norm, then it admits a fixed point. An example of this catagory is the well-known Banach fixed point theorem. The second type, on the other hand, is usually based on the compact properties of the range space and operator, this includes Brouwer's fixed point theorem (compact, covex range space and continuous operators), Schauder's fixed point theorem (closed, bounded, convex range space and compact operators), and Kakutani-Fan-Glicksberg fixed point theorem for set-value functions, which is the one we will use in the proof of equilibrium existence (Appendix D).

Contraction based fixed point theorems usually have stronger assumptions, since the contraction in norm property is hard to verify. However, it provides clear approaches to find the fixed point when one exists: starting from an appropriate initial point, we may iteratively apply the operator and the result is guaranteed to converge to a fixed point. On the other hand, compact based fixed point theorems require weaker assumptions, but they fail to indicate how to find a fixed point rather than telling its theoretical existence.

As learning algorithms are designed to find the fixed point (equilibrium) of the game, they usually try to approximate the contraction mapping with estimations (since the environment is usually unknown) in order to demonstrates the convergences of algorithm. Thus the contraction mapping assumption is unfortunately necessary in complexity analysis proofs (see Assumption 5.1(3)), even though we do not make such assumptions in pure mathematical analysis in section 4.

B. Comparing Representative-player Games and Continuum-player Games

B.1. Continuum-Player Graphon Game

In this section we give a review on continuum-player graphon games in previous works. Consider a game with a continuum of players, labeled with $u \in [0, 1]$, and we assume the label space $[0, 1]$ is equipped with Borel- σ -algebra and Lebesgue measure. Each player u admits a state process X^u valued in \mathbb{R}^d . Let the population measure be a collection $\mu = \{\mu^u\}_{u \in [0, 1]}$, which is given and fixed, and it is usually an assumption that $u \mapsto \mu^u$ is a probabilistic kernel, i.e. $u \mapsto \mu^u(B)$ is a measurable function for any Borel subset $B \subset \mathcal{C}$.

Let $\mathcal{V} := \mathcal{P}(A)^T$ denote the space of all open-loop policy flows. Each element $\pi \in \mathcal{V}$ is a sequence of measures $\{\pi_t\}_{t \in \mathbb{T}}$ for $\pi_t \in \mathcal{P}(A)$. Alternatively, we may think \mathcal{V} as the collection of measures on $\mathbb{T} \times A$ with uniform first marginals. Let \mathcal{A} be the collection of all the feedback (closed-loop) policies $\mathbb{R}^d \rightarrow \mathcal{P}(A)$, and assume player u adopts a policy $\pi^u \in \mathcal{A}$, the state process follows

$$\begin{aligned} X_0^u &\sim \lambda^u \\ a_t &\sim \pi_t^u(X_t^u) \quad X_{t+1}^u \sim P_t(X_t^u, \mathbf{W}\mu_t(u), a_t) \end{aligned}$$

for some initial condition $\lambda^u \in \mathcal{P}(\mathbb{R}^d)$. Note that all the players' state dynamics are independent, in the following sense: for almost every $u \in [0, 1]$, X^u is independent of X^v for almost every $v \in [0, 1]$. Indeed this independence leads to a significant measurability issue that many papers ignore, and we will give a detailed discussion in Appendix B.2. Each player u aims to maximize an objective function

$$J^u(\mu, \pi^u) := \mathbb{E} \left[\sum_{t \in \mathbb{T}} f_t(X_t^{u, \pi^u}, \mathbf{W}\mu_t(u), a_t) + g(X_T^{u, \pi^u}, \mathbf{W}\mu_T(u)) \right]$$

where we denote X^{u,π^u} to emphasize the process X^u is controlled by policy π^u . The equilibrium is defined as a pair $(\hat{\mu}, \hat{\pi}) := (\{\hat{\mu}^u\}_{u \in [0,1]}, \{\hat{\pi}^u\}_{u \in [0,1]}) \in \mathcal{P}_{\text{unif}}([0,1] \times \mathcal{C}) \times \mathcal{A}^{[0,1]}$ such that

$$\begin{aligned} J^u(\hat{\mu}, \hat{\pi}^u) &= \sup_{\pi \in \mathcal{V}} J^u(\hat{\mu}, \pi) \\ \hat{\mu}^u &= \mathcal{L}(X^{u,\hat{\pi}^u}) \end{aligned}$$

for almost every $u \in [0, 1]$. This game "continuum-player formulation" since it involves a continuum of players.

B.2. Comparing representative-player games and continuum-player games

The representative-player graphon game we present in section 4.1 and the continuum-player graphon game in section B.1 are not mathematically equivalent. The representative-player formulation in section 4.1 provides some advantages both conceptually and technically.

Conceptually, our representative-player formulation inherits the spirit of mean-field game. We recall that there is only one representative player in the mean-field game, and all other players are abstracted into a population measure in $\mathcal{P}(\mathcal{C})$. Similarly, our game formulation is for one representative player, and the difference is that now the representative player is in addition assigned a random label, while all other players are abstracted into a label-state joint population measure on $\mathcal{P}_{\text{unif}}([0,1] \times \mathcal{C})$.

Mathematically, the representative-player formulation avoids significant measurability difficulties that the continuum-player formulation suffers from. For completeness, we first cite proposition 2.1 from (Sun, 2006) as follows:

Proposition B.1. Consider index space $(I, \mathcal{I}, \lambda)$ and probability space (Ω, \mathcal{F}, P) . Consider function $f : I \times \Omega \rightarrow E$ for some Polish space E . If f is measurable on the product space $(I \times \Omega, \mathcal{I} \otimes \mathcal{F}, \lambda \otimes P)$, equipped with the usual product σ -algebra, and for λ -almost every $j \in I$, f_j is independent of f_i for λ -almost every $i \in I$. Then, for λ -almost every $i \in I$, f_i is a constant random variable.

Intuitively, the product σ -algebra $\mathcal{I} \otimes \mathcal{E}$ fails to support the large amount of information when we require both the joint measurability of f , and the independence between f_i and f_j . This would lead to a problem when we consider a continuum of players, even if the state space is a finite space rather than \mathbb{R}^d , and even for a static game. Unfortunately many of the previous works on continuum-player setting ignored this measurability problem.

More precisely, let (Ω, \mathcal{F}, P) be a probability space that supports a collection of stochastic processes $\{X^u : u \in [0, 1]\}$, where X^u is a process on $\{0, 1, \dots, T\}$ valued in \mathcal{X} (which could be \mathbb{R}^d or a finite state space). X^u represents the state process of player with label u . From time $t - 1$ to t , X_t^u are generated independently for every $u \in [0, 1]$, and thus the mapping $u \rightarrow X^u(\omega)$ is not measurable for P -almost every $\omega \in \Omega$, similarly $u \mapsto \pi^u$ is not measurable. This measurability issue leads to significant difficulties in the proof, as the objective reward function may involve these mappings. For instance, as one attempts to transform the continuum-player graphon game into a mean-field game with augmented state space, the objective becomes

$$\mathbb{E} \int_{[0,1]} \left[\sum_{t \in \mathbb{T}} \bar{f}_t((X_t^{u,\pi^u}, u), \mu_t, a_t) + \bar{g}((X_T^{u,\pi^u}, u), \mu_T) \right] du$$

where the integral with respect to u over $[0, 1]$ is not well-defined since the integrand is not measurable. A similar argument demonstrates why we cannot aggregate the objective of all the players in a continuum-player graphon game, since the mapping $[0, 1] \ni u \mapsto J^u(\mu, \pi^u) \in \mathbb{R}$ is not measurable, the integral $\int_{[0,1]} J^u(\mu, \pi^u) du$ is not well-defined. Thus the continuum-player graphon game is not mathematically equivalent to our representative-player formulation.

This technical issue can be addressed by carefully enlarging the σ -algebra with rich Fubini extension [(Sun, 2006), section 2], allowing it to hold more information while ensuring the joint measurability and independence (Aurell et al., 2022; Tangpi & Zhou, 2023). However this approach is restricted to linear-quadratic problems.

On the contrary, our graphon game formulation considers only one representative player. Recall that for any $\mu \in \mathcal{P}_{\text{unif}}([0,1] \times \mathcal{C})$, the conditional law of X given U yielded by disintegration is uniquely defined Lebesgue almost surely, thus it encodes less information by only considering almost every label u , but this provides great technical convenience and allow us to consider the game for one player (Lacker & Soret, 2022).

C. A Toy Example on the Difference between the Two Formulations

In this section we compare two types of graphon games formulations on a toy example, inspired by the motivating example in (Cui & Koepll, 2021). The two types of formulations of graphon games lead to the same equilibrium in this particular one-shot game, while the representative-player graphon game is simpler in formulation. When a finite player game contains larger and continuous state and action spaces with more complex settings, our formulation would demonstrate more advantage in both analysis and computation. Note that this toy example focus on demonstrating the difference in formulation, and the measurability issue mentioned in Appendix B.2 is not the main point here as the example is simple enough to be solved explicitly, and no technical proofs are involved.

The interaction is defined by a threshold graph, where $\xi_{ij}^n = 1_{i+j < n}$. It is easy to see that W_{ξ^n} converges in cut norm to a limit defined by $W(u, v) := 1_{u+v < 1}$. Note that this graphon is discountinuous.

C.1. N-player Game

Consider a one-shot (single-stage) game for n players, and let the state and action space be $\mathcal{X} = A = \{-1, 1\}$, understood as left and right. Each player simultaneously chooses either left or right, and is punished by the weighted average of proportion of players that chose the same action. Precisely,

$$a^i = \begin{cases} 1 & \text{w.p. } p^i \\ -1 & \text{w.p. } 1 - p^i \end{cases} \quad X^i = a^i$$

where p^i is the probability player i choose right (state 1), and this characterizes the policy. Let $\mathbf{p} = (p^1, \dots, p^n)$ and let the terminal reward be $g(x, m) = -\langle m, 1_x \rangle$, where 1_x is the indicator function. Player i aims to maximize

$$\begin{aligned} J^i(\mathbf{p}) &= -\mathbb{E}\left(\sum_{j=1}^n \xi_{ij}^n 1_{X^i=X^j}\right) \\ &= -\sum_{j=1}^{n-i} (p^i p^j + (1-p^i)(1-p^j)) \end{aligned}$$

It can be verified that the equilibrium is given by $p^1 = \dots = p^n = \frac{1}{2}$.

C.2. Representative-player formulation

Consider a one-shot game for a single player, and let the state and action space be $\mathcal{X} = A = \{-1, 1\}$. Any population measure $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$ can be characterized by a function $q(u) := \mu(u, \{1\})$, $\forall u \in [0, 1]$. Let this population measure be fixed. The graphon operator is given by

$$W\mu(u) = \int_{[0,1]} W(u, v)(q(v)\delta_1 + (1-q(v))\delta_{-1})dv \in \mathcal{M}_+([-1, 1])$$

where δ is Dirac delta measure. The player is randomly assigned a label $U \sim \text{unif}[0, 1]$, and let $\pi(u)$ be her policy. Equivalently the policy can be characterized by $p(u) := \pi(u)(\{1\})$. Then she follows the dynamic

$$a = \begin{cases} 1 & \text{w.p. } p(U) \\ -1 & \text{w.p. } 1 - p(U) \end{cases} \quad X = a$$

The objective is

$$\begin{aligned} J(q, p) &= -\mathbb{E}\left(\langle W\mu(U), 1_X \rangle\right) \\ &= -\mathbb{E}\left(\int_{[0,1]} W(U, v)(q(v)1_{X=1} + (1-q(v))1_{X=-1})dv\right) \\ &= -\int_{u+v<1} (q(v)p(u) + (1-q(v))(1-p(u)))dvdu \end{aligned}$$

Solving this as a calculus of variation problem provides a necessary condition $\int_0^u q(v)dv = \frac{1}{2}$, $\forall u \in [0, 1]$, and thus the equilibrium is given by $p(u) = \frac{1}{2}$ for a.e. u , and $q(v) = \frac{1}{2}$ for a.e. v .

C.3. Continuum-player formulation

Consider a static game for a continuum of players with the same setting, and let the population measure be $\mathbf{q} := \{q^u\}_{u \in [0,1]}$ for $q^u = \mu^u(\{1\})$. Let this population measure be fixed. Each player $u \sim [0, 1]$ admits a policy π^u as the probability choosing 1, and denote $\mathbf{p} := \{p^u\}_{u \in [0,1]}$. Then the player u chooses the action

$$a^u = \begin{cases} 1 & \text{w.p. } p^u \\ -1 & \text{w.p. } 1 - p^u \end{cases} \quad X^u = a^u$$

and optimize the objective

$$\begin{aligned} J^u(\mathbf{q}, p^u) &= -\mathbb{E}\left(\langle \mathbf{W}\mu(u), 1_X \rangle\right) \\ &= -\mathbb{E}\left(\int_{[0,1]} W(u, v)(q^v 1_{X=1} + (1 - q^v) 1_{X=-1}) dv\right) \\ &= -\int_0^{1-u} (q^v p^u + (1 - q^v)(1 - p^u)) dv \end{aligned}$$

It is immediate that the equilibrium is given by $p^u = \frac{1}{2}$, and $q^v = \frac{1}{2}$ for almost every u, v . Note that the measurability issue is not a concern for this specific example, since it can be solved directly and thus doesn't involve technical analysis.

D. Proof for Existence

D.1. Preliminary Lemmas

Lemma D.1. [Lemma A.2 of (Lacker, 2015)] Let X_1 and X_2 be Polish spaces. Defined the coordinate projections $\Pi_i : \mathcal{X}_1 \times \mathcal{X}_2 \rightarrow \mathcal{X}_i$ for $i = 1, 2$. Then a set $S \subset \mathcal{P}(\mathcal{X}_1 \times \mathcal{X}_2)$ is tight if and only if the sets $S_1 = \{\mu \circ \Pi_1^{-1} : \mu \in S\}$ and $S_2 = \{\mu \circ \Pi_2^{-1} : \mu \in S\}$ are tight in $\mathcal{P}(X_1)$ and $\mathcal{P}(X_2)$ respectively.

Lemma D.2. [Corollary A.5 of (Lacker, 2015)] Let E, F, G be complete, separable metric spaces. $\phi : E \times F \times G \rightarrow \mathbb{R}$ is a bounded measurable function, with $\phi(x, \cdot, \cdot)$ being jointly continuous for any $x \in E$. Then the following mapping is continuous:

$$G \times \mathcal{P}(E \times F) \ni (z, P) \mapsto \int_{E \times F} \phi(x, y, z) P(dx, dy)$$

Lemma D.3. Let W^n, W be graphons. If $W^n \xrightarrow{L_1} W$, then, $W^n \rightarrow W$.

Proof. Given any $\psi \in L_\infty[0, 1]$,

$$\begin{aligned} \|\mathbf{W}^n \psi - \mathbf{W} \psi\|_1 &= \int_{[0,1]} |\mathbf{W}^n \psi(u) - \mathbf{W} \psi(u)| du \\ &= \int_{[0,1]} \left| \int_{[0,1]} W^n(u, v) \psi(v) - W(u, v) \psi(v) dv \right| du \\ &\leq \|\psi\|_\infty \int_{[0,1]^2} |W^n(u, v) - W(u, v)| dv du \\ &= \|\psi\|_\infty \|W^n - W\|_1 \rightarrow 0 \end{aligned}$$

□

Lemma D.4. [Lemma 4.2 of (Lacker & Soret, 2022)] Let E be any Polish space, and W be any graphon.

1. For a.e. $u \in [0, 1]$, the following map is continuous:

$$\mathcal{P}_{\text{unif}}([0, 1] \times E) \ni \mu \mapsto \mathbf{W}\mu(u) \in \mathcal{M}_+(E)$$

990 2. Suppose the map $[0, 1] \ni u \mapsto W(u, v)dv \in \mathcal{M}_+([0, 1])$ is continuous, then for any $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times E)$,

$$[0, 1] \ni u \mapsto W\mu(u) \in \mathcal{M}_+(E)$$

993 is continuous.

994 **Lemma D.5.** Let E, F be complete, separable metric space, and F is a regular measurable space. Consider a sequence of
995 probability measures on the product space $\{\nu_n\} \subset \mathcal{P}(E \times F)$. Suppose that ν_n admits disintegration
996

$$\nu_n(dx, dy) = \mu_n(dx)K(x, dy)$$

997 for some common kernel K , which is continuous as a mapping $E \rightarrow \mathcal{P}(F)$, i.e. any sequence $x_n \rightarrow x$ implies $K_{x_n} \Rightarrow K_x$.
998 Then if $\nu_n \Rightarrow \nu$, ν admits a disintegration $\nu(dx, dy) = \mu(dx)K(x, dy)$ for some $\mu \in \mathcal{P}(E)$.

1001 *Proof.* Let Π_1 be the projection to first coordinate, which is a continuous mapping. By continuous mapping theorem, the
1002 pushforward of a weak convergence measure seqnce under continuous mapping converge weakly:

$$\nu_n \circ \Pi_1^{-1} \Rightarrow \nu \circ \Pi_1^{-1} =: \mu$$

1006 Suppose ν admits disintegration $\nu(dx, dy) = \mu(dx)\bar{K}(x, dy)$ for some \bar{K} . Given $\forall \phi : E \times F \rightarrow \mathbb{R}$ bounded, jointly
1007 continuous, the mapping $E \ni x \mapsto \int_F \phi(x, y)K(x, dy) \in \mathbb{R}$ is bounded and continuous since for any $x_n \rightarrow x$,

$$\begin{aligned} & \left| \int_F \phi(x_n, y)K(x_n, dy) - \int_F \phi(x, y)K(x, dy) \right| \\ & \leq \left| \int_F \phi(x_n, y)K(x_n, dy) - \int_F \phi(x, y)K(x_n, dy) \right| + \left| \int_F \phi(x, y)K(x_n, dy) - \int_F \phi(x, y)K(x, dy) \right| \end{aligned}$$

1013 which converges to 0. Finally, $\langle \nu_n, \phi \rangle \rightarrow \langle \nu, \phi \rangle$, and on the other hand,

$$\langle \nu_n, \phi \rangle = \int_{E \times F} \phi(x, y)K(x, dy)\mu_n(dx) \longrightarrow \int_{E \times F} \phi(x, y)K(x, dy)\mu(dx)$$

1017 which holds for any ϕ bounded continuous, and we conclude that K is a version of \bar{K} . \square

D.2. Existence of Equilibrium

1021 Given any function $\phi : E \times A \rightarrow F$ for Polish space E, F and a measure $\pi \in \mathcal{P}(A)$, we may also abuse the notation by
1022 writing ϕ as a function $E \times \mathcal{P}(A) \rightarrow F$, defined by $\phi(x, \pi) = \langle \pi, \phi(x, \cdot) \rangle$ for each $x \in E$.

1023 Throughout the proof we fix a graphon W , and denote $\mathcal{V} = \mathcal{P}(A)^T$ the space of all policies. We fix any policy $\pi \in \mathcal{V}$, and
1024 construct the label-state joint measure of the representative player controlled by π as follows. Recall that at time t given
1025 $U = u, X_t = x, \alpha_t = a$, the law of next state X_{t+1} follows the probabilistic kernel $[0, 1] \times \mathbb{R}^d \times A \rightarrow \mathbb{R}^d$

$$\mathcal{L}(X_{t+1}|X_t = x, U_t = u, \alpha_t = a)(dy) = P_t(dy|x, W\mu_t(u), a) \quad \forall y \in \mathbb{R}^d$$

1029 and the control process α_t follows

$$\mathcal{L}(\alpha_t)(da) = \pi_t(da) \quad \forall a \in A$$

1032 We may thus consider

$$\widehat{P}_t^{\pi, \mu}(dy|u, x) := \mathcal{L}(X_{t+1}|X_t = x, U_t = u)(dy) = \int_A P_t(dy|x, W\mu_t(u), a)\pi_t(da) \quad \forall y \in \mathbb{R}^d$$

1036 and we use the superscript to emphasize that the law is controlled by the policy π . Note that $\mathcal{V}_U \ni \pi \mapsto \widehat{P}_t^{\pi, \mu}(u, x) \in \mathcal{P}(\mathbb{R}^d)$
1037 is measurable. The collection of kernels $\{\widehat{P}_t^\pi\}_{t \in \mathbb{T}}$ (recall $\mathbb{T} = \{0, 1, \dots, T-1\}$) along with the initial law λ implies a
1038 label-state joint law in $\mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$

$$\widehat{P}^{\pi, \mu}(du, dx) := \mathcal{L}(U, X)(du, dx) = \lambda(du, dx_0) \prod_{t \in \mathbb{T}} \widehat{P}_t^{\pi, \mu}(dx_{t+1}|u, x_t) \quad \forall (u, x) \in [0, 1] \times \mathcal{C}$$

1042 which is the label-state joint measure of the representative player, when her state dynamic is controlled by π . Since the
1043 space $[0, 1] \times \mathcal{C}$ is a standard measurable space, this is understood as a regular version of the kernel from \mathcal{V} to $[0, 1] \times \mathcal{C}$.

1045 **Lemma D.6.** Under assumption 4.3(5), for any $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$, $\pi \mapsto \widehat{P}^{\pi, \mu}$ is continuous. In particular, $\pi_t \mapsto$
 1046 $\widehat{P}_t^{\pi, \mu}(u, x)$ is continuous for every $(u, x) \in [0, 1] \times \mathbb{R}^d$.
 1047

1048 *Proof.* Let $\{\pi^n\} \subset \mathcal{V}$ be any sequence of policies such that $\pi^n \Rightarrow \pi$ for some $\pi \in \mathcal{V}$. For any $\phi : [0, 1] \times \mathcal{C} \rightarrow \mathbb{R}$ bounded
 1049 continuous,
 1050

$$\begin{aligned} & \int_{[0,1] \times \mathcal{C}} \phi(u, x) \widehat{P}^{\pi^n, \mu}(du, dx) \\ &= \int_{[0,1] \times \mathbb{R}^d} \left[\int_{A^T} \int_{(\mathbb{R}^d)^T} \phi(u, x_0, \dots, x_T) \prod_{t \in \mathbb{T}} P_t(dx_{t+1} | x_t, \mathbf{W}\mu_t(u), a_t) \pi^n(da_0, \dots, da_{T-1}) \right] \lambda(du, dx_0) \\ &=: \int_{[0,1] \times \mathbb{R}^d} \left[\int_{A^T} \psi(a_0, \dots, a_{T-1}) \pi^n(da_0, \dots, da_{T-1}) \right] \lambda(du, dx_0) \end{aligned}$$

1058 where

$$\psi(a_0, \dots, a_{T-1}) := \int_{(\mathbb{R}^d)^T} \phi(u, x_0, \dots, x_T) \prod_{t \in \mathbb{T}} P_t(dx_{t+1} | x_t, \mathbf{W}\mu_t(u), a_t)$$

1063 We know that $a_t \mapsto P_t(dx_{t+1} | x_t, \mathbf{W}\mu_t(u), a_t)$ is continuous for each $t \in \mathbb{T}$ by assumption 4.3(5), and since $(\mathbb{R}^d)^T$ is
 1064 separable, with standard measure theory argument for weak convergence on product space, for instance chapter 2 of
 1065 (Billingsley, 1995), the map ψ is continuous. Thus $\langle \pi^n, \psi \rangle \rightarrow \langle \pi, \psi \rangle$, and

$$\begin{aligned} & \int_{[0,1] \times \mathcal{C}} \phi(u, x) \widehat{P}^{\pi^n, \mu}(du, dx) \\ & \rightarrow \int_{[0,1] \times \mathbb{R}^d} \left[\int_{A^T} \psi(a_0, \dots, a_{T-1}) \pi(da_0, \dots, da_{T-1}) \right] \lambda(du, dx_0) \\ &= \int_{[0,1] \times \mathcal{C}} \phi(u, x) \widehat{P}^{\pi, \mu}(du, dx) \end{aligned}$$

1075 \square

1076 Define the probability space $\Omega := \mathcal{V} \times [0, 1] \times \mathcal{C}$, equipped with the product σ -algebra. A typical element of Ω is (π, u, x) ,
 1077 where we understood them as a policy, a label of the representative player and the player's path, respectively. Let the
 1078 coordinate maps be Λ, U, X respectively. The filtration is given by $\mathcal{F}_t = \sigma\{\Lambda|_{[t] \times A}, U, \{X_s\}_{0 \leq s \leq t}\}$.
 1079

1080 The collection of admissible laws $\mathcal{R}(\mu)$ is defined as the set

$$1082 \mathcal{R}(\mu) := \{R \in \mathcal{P}(\Omega) : R \text{ admits disintegration } R(d\pi, du, dx) = R_\Lambda(d\pi) \widehat{P}^{\pi, \mu}(du, dx) \text{ for some } R_\Lambda \in \mathcal{P}(\mathcal{V})\}$$

1084 Define a random variable $\Xi^\mu : \Omega \rightarrow \mathbb{R}$ by

$$1086 \Xi^\mu(\pi, u, x) := \sum_{t \in \mathbb{T}} \int_A f_t(\mathbf{W}\mu_t(u), x_t, a) \pi_t(da) + g(x_T, \mathbf{W}\mu_T(u)) \quad (15)$$

1089 where μ_t is the marginal obtained as the image by $(u, x) \mapsto (u, x_t)$. In particular, given a policy $\pi \in \mathcal{V}$, let
 1090 $R^{(\pi)}(d\tilde{\pi}, du, dx) := \delta_\pi(d\tilde{\pi}) \widehat{P}^{\tilde{\pi}, \mu}(du, dx)$ be an element of $\mathcal{R}(\mu)$, where δ is the Dirac measure. It holds that the ob-
 1091 jective can be rewritten as
 1092

$$1093 J_W(\mu, \pi) = \langle R^{(\pi)}, \Xi^\mu \rangle$$

1094 thus the expectation $\langle R, \Xi^\mu \rangle$ is a reformulation of the objective, and a single player's objective is to find the collection of
 1095 measures that maximize this expectation:

$$1098 \mathcal{R}^*(\mu) := \{R^* \in \mathcal{R}(\mu) : \langle R^*, \Xi^\mu \rangle \geq \langle R, \Xi^\mu \rangle, \forall R \in \mathcal{R}(\mu)\} \quad (16)$$

Define the correspondence (i.e. set valued function, see (Aliprantis & Border, 2006) for an overview) $\Phi : \mathcal{P}([0, 1] \times \mathcal{C}) \rightarrow 2^{\mathcal{P}([0, 1] \times \mathcal{C})}$, given by

$$\Phi(\mu) := \{R \circ (U, X)^{-1} : R \in \mathcal{R}^*(\mu)\}$$

The existence of W -equilibrium is divided into two steps: we first show the existence of optimizer to the optimization problem (16) over the probability measures, i.e. $\mathcal{R}^*(\mu)$ is non-empty for any μ ; Next, to obtain a W -equilibrium, we aim to find a fixed point for the correspondence Φ .

Proposition D.7. *For $\forall \mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$, the following optimization problem admits optimizer.*

$$\sup_{R \in \mathcal{R}(\mu)} \langle R, \Xi^\mu \rangle$$

Proof. We want to show that $R \mapsto \langle R, \Xi^\mu \rangle$ is a continuous mapping on compact space, and thus the maximum of this mapping is attained. With a direct application of lemma D.2, we immediately conclude that the following map is jointly continuous:

$$\text{Gr}(\mathcal{R}) \ni (\mu, R) \longmapsto \langle R, \Xi^\mu \rangle \in \mathbb{R} \quad (17)$$

where Gr denotes the graph of an operator.

It remains to prove that $\mathcal{R}(\mu)$ is compact. First we want to show $\mathcal{R}(\mu)$ is tight for any μ . By lemma D.1, it suffices to show that the following sets are tight: $\{R \circ X^{-1} : R \in \mathcal{R}(\mu)\}$, $\{R \circ U^{-1} : R \in \mathcal{R}(\mu)\}$, $\{R \circ \Lambda^{-1} : R \in \mathcal{R}(\mu)\}$. The last two follows immediately from the fact that $[0, 1]$ and A are compact spaces.

Fix any $\epsilon' > 0$, we could always find some ϵ such that $(1 - \epsilon)^{T+1} > 1 - \epsilon'$. By assumption 4.3(3) and 4.3(4), let $\{K_t\}_{t \in \mathbb{T}}$ be compact subsets of \mathbb{R}^d such that

$$\inf_{u \in [0, 1]} \lambda^u(K_0) > 1 - \epsilon \quad \inf_{\tilde{P}_t \in \zeta_t} \tilde{P}_t(K_{t+1}) > 1 - \epsilon \quad \forall t \in \mathbb{T}$$

Define $K = \prod_{t=0}^T K_t$, which is a compact subset of \mathcal{C} . Then for every $R \in \mathcal{R}(\mu)$, let $\hat{P}^{\pi, \mu}(du, dx)R_\Lambda(d\pi)$ be its disintegration,

$$\begin{aligned} (R \circ X^{-1})(K) &= R(\mathcal{V} \times [0, 1] \times K) \\ &= \int_{\mathcal{V} \times [0, 1] \times \mathcal{C}} 1_K(x) \hat{P}^{\pi, \mu}(du, dx) R_\Lambda(d\pi) \\ &= \int_{\mathcal{V}} \int_{[0, 1] \times \mathbb{R}^d} \left[\prod_{t=0}^{T-1} \int_A \int_{\mathbb{R}^d} 1_{K_{t+1}}(x_{t+1}) P_t(dx_{t+1} | x_t, \mathbf{W}\mu_t(u), a) \pi_t(da) \right] 1_{K_0}(x_0) \lambda(du, dx_0) R_\Lambda(d\pi) \\ &\geq \int_{\mathcal{V}} \int_{[0, 1]} \left[\prod_{t=0}^{T-1} \int_A (1 - \epsilon) \pi_t(da) \right] \int_{\mathbb{R}^d} 1_{K_0}(x_0) \lambda^u(dx_0) du R_\Lambda(d\pi) \\ &\geq \int_{\mathcal{V}} \int_{[0, 1]} (1 - \epsilon)^{T+1} du R_\Lambda(d\pi) \\ &= (1 - \epsilon)^{T+1} > 1 - \epsilon' \end{aligned}$$

and thus we have $\inf_{R \in \mathcal{R}(\mu)} (R \circ X^{-1})(K) > 1 - \epsilon'$, which implies the tightness of $\{R \circ X^{-1} : R \in \mathcal{R}(\mu)\}$. Note that if the state space \mathcal{X} of dynamic X is compact, then $\{R \circ X^{-1} : R \in \mathcal{R}(\mu)\}$ being tight is immediate. By Prokhorov's theorem, $\mathcal{R}(\mu)$ is precompact.

We conclude by showing that $\mathcal{R}(\mu)$ is closed. Let $\{R_n\} \subset \mathcal{R}(\mu)$, and $R_n \Rightarrow R$. Indeed each R_n admits disintegration $R_\Lambda^n(d\pi) \hat{P}^{\pi, \mu}(du, dx)$ for some $R_\Lambda^n \in \mathcal{P}(\mathcal{V})$, and the kernel $\pi \mapsto \hat{P}^{\pi, \mu}$ is continuous by lemma D.6. Then lemma D.5 implies that R admits disintegration $R_\Lambda(d\pi) \hat{P}^{\pi, \mu}(du, dx)$ and thus $R \in \mathcal{R}(\mu)$. \square

1155 Next we show that the correspondence Φ admits a fixed point, and thus the graphon game admits a W -equilibrium.

1156 **Proposition D.8.** *There exists a fixed point $\hat{\mu}$ for the correspondence Φ .*

1158 *Proof.* We aim to apply the Kakutani-Fan-Glicksberg fixed point theorem, which is a classic fixed point theorem for
 1159 correspondences, see for instance theorem 17.55 of (Aliprantis & Border, 2006). We need to show the following conditions:
 1160 $\exists K \subset \mathcal{P}([0, 1] \times \mathcal{C})$ nonempty, convex and compact, such that
 1161

- 1162 1. $\Phi(\mu) \subset K$ for each $\mu \in K$.
- 1163 2. $\Phi(\mu)$ is nonempty and convex for each $\mu \in K$.
- 1164 3. The graph $\text{Gr}(\Phi) = \{(\mu, \mu') : \mu \in K, \mu' \in \Phi(\mu)\}$ is closed.

1167 We start from defining K , note that λ is a fixed initial measure.

$$1169 K := \{\lambda \otimes \prod_{t=0}^{T-1} \hat{P}_t : \hat{P}_t \in \overline{\text{conv}}(\zeta_t)\}$$

1172 where $\overline{\text{conv}}(\cdot)$ denotes the closed convex hull of a set, and \otimes is the combinations of probabilistic kernels on product space.
 1173 K is obviously non-empty. By construction, K is the finite cartesian product of convex sets, thus K is convex. To show
 1174 K is compact, it suffices to show $\overline{\text{conv}}(\zeta_t)$ is compact for each $t \in \mathbb{T}$, since Tychonoff's theorem asserts that an arbitrary
 1175 product of compact spaces is again compact. Indeed, since ζ_t is tight, and thus precompact by Prokhorov's theorem, and the
 1176 closed convex hull of a precompact set is compact in a locally convex Hausdorff space. Again, if the value space \mathcal{X} of X is
 1177 compact, let $K = \mathcal{P}([0, 1] \times \mathcal{C})$ and K is compact automatically.

1178 For each $R \in \mathcal{R}(\mu)$, let it admit the disintegration $R = R_\Lambda \otimes \hat{P}$:

$$1180 R_\Lambda(d\pi) \hat{P}^{\pi, \mu}(du, dx) = \left[\lambda(du, dx_0) \prod_{t=0}^{T-1} \int_A P_t(dx_{t+1}|x_t, \mathbf{W}\mu_t(u), a) \pi_t(da) \right] R_\Lambda(d\pi)$$

1183 We claim that for any $t \in \mathbb{T}$,

$$1185 \hat{P}_t^{\pi, \mu}(dx_{t+1}|u, x_t) = \int_A P_t(dx_{t+1}|x_t, \mathbf{W}\mu_t(u), a) \pi_t(da) \in \overline{\text{conv}}(\zeta_t)$$

1187 since it is the limit of convex combinations of $P_t(\cdot|x_t, \mathbf{W}\mu_t(u), a) \in \zeta_t$. and thus for any $(\pi, \mu) \in \mathcal{V} \times \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$,
 1188 the measure $\hat{P}^{\pi, \mu} \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$ belongs to K . The pushforward of R onto (U, X) coordinate is
 1189

$$1190 (R \circ (U, X)^{-1})(du, dx) = \int_{\mathcal{V}} \hat{P}^{\pi, \mu}(du, dx) R_\Lambda(d\pi)$$

1193 which is also the limit of a sequence of convex combinations of $\hat{P}^{\pi, \mu}(du, dx)$, indexed by π . Thus by closedness and
 1194 compactness of K , $R \circ (U, X)^{-1} \in K$, and thus $\Phi(\mu) \subset K$.

1195 To show the convexity of $\Phi(\mu)$, we start with showing $\mathcal{R}(\mu)$ is convex since for any $R^1 = R_\Lambda^1 \otimes \hat{P}$ and $R^2 = R_\Lambda^2 \otimes \hat{P}$
 1196 and $\lambda \in [0, 1]$, $\lambda R^1 + (1 - \lambda) R^2 = (\lambda R_\Lambda^1 + (1 - \lambda) R_\Lambda^2) \otimes \hat{P} \in \mathcal{R}(\mu)$. Convexity of $\mathcal{R}^*(\mu)$ follows from the linearity of
 1197 $R \mapsto \langle R, \Xi^\mu \rangle$ and the convexity of $\mathcal{R}(\mu)$, and thus the convexity of $\Phi(\mu)$ follows from the linearity of map $R \mapsto R \circ (U, X)^{-1}$
 1198 and the convexity of $\mathcal{R}^*(\mu)$.

1199 It remains to show the closedness of the graph of Φ , and we first show the closedness of the following set:

$$1200 \{(\mu, R) : \mu \in K, R \in \mathcal{R}^*(\mu)\}$$

1203 Let $\mu_n \Rightarrow \mu$ and $R_n \Rightarrow R$ with $\mu_n, \mu \in K$, $R_n \in \mathcal{R}^*(\mu_n)$, and $R \in \mathcal{R}$. To show that $R \in \mathcal{R}^*(\mu)$, we use the continuity
 1204 (17), and for any $R' \in \mathcal{R}$,

$$1205 \langle R, \Xi^\mu \rangle = \lim_{n \rightarrow \infty} \langle R_n, \Gamma^{\mu_n} \rangle \geq \lim_{n \rightarrow \infty} \langle R', \Gamma^{\mu_n} \rangle = \langle R', \Xi^\mu \rangle$$

1207 and thus $\langle R, \Xi^\mu \rangle \geq \langle R', \Xi^\mu \rangle$ for any $R' \in \mathcal{R}$. The by the continuity of $R \mapsto R \circ (U, X)^{-1}$ and compactness of K , we have
 1208 the closedness of $\text{Gr}(\Phi)$. \square

D.3. Closed-loop equilibrium optimal policy

In this section we show the second part of theorem 4.4, that the equilibrium optimal open-loop policy can be made closed-loop.

Proposition D.9. *Let $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$, and $R \in \mathcal{R}(\mu)$. Then, there exists a closed-loop optimal policy, in the following sense: \exists a measurable function $\pi : \mathbb{T} \times [0, 1] \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$, and $R^0 \in \mathcal{R}(\mu)$, such that*

1. $R^0(\Lambda_t(da)) = \pi_t(U, X_t)(da) \quad \forall t \in \mathbb{T} = 1$
2. $\int_{\Omega} \Xi^{\mu} dR^0 \geq \int_{\Omega} \Xi^{\mu} dR$
3. $R^0 \circ (U, X_t)^{-1} = R \circ (U, X_t)^{-1} \quad \forall t \in \mathbb{T}$

Corollary D.10. *There exists a closed-loop equilibrium optimal policy to the graphon game.*

Proof. We first find a space $(\Omega^1, \mathcal{F}^1, R^1)$ supporting a random variable U^1 , an adapted process X^1 valued in \mathbb{R}^d , and a $\mathcal{P}(A)$ -valued adpated process Λ_t s.t.

$$\begin{aligned} (U^1, X_0^1) &\sim \lambda \quad X_{t+1}^1 \sim P_t(X_t^1, \mathbf{W}\mu_t(U^1), \Lambda_t) \\ R^1 \circ (U, X^1)^{-1} &= \mu \end{aligned}$$

The existence of such space is guarenteed by reasoning in Appendix D.2. We claim that there exists a measurable $\pi : \mathbb{T} \times [0, 1] \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$ such that

$$\pi(t, U^1, X_t^1) = \mathbb{E}^{R^1}(\Lambda_t | U^1, X_t^1) \quad R^1 - a.s. \quad \forall t \in \mathbb{T}$$

More precisely, for every bounded measurable $\phi : [0, 1] \times \mathbb{R}^d \times A \rightarrow \mathbb{R}$,

$$\int_A \phi(U^1, X_t^1, a) \pi(t, U^1, X_t^1)(da) = \mathbb{E}^{R^1} \left(\int_A \phi(U^1, X_t^1, a) \Lambda_t(da) \mid U^1, X_t^1 \right) \quad R^1 - a.s., \forall t \in \mathbb{T}$$

Define a collection of measures, $\{\eta_t\}_{t \in \mathbb{T}}$, $\eta_t \in \mathcal{P}([0, 1] \times \mathbb{R}^d \times A)$ by

$$\eta_t(C) := \mathbb{E}^{R^1} \left[\int_A 1_C(t, U_t^1, X_t^1, a) \Lambda_t(da) \right]$$

and let η_t adimits disintegration $\eta_t(du, dx, da) = \eta'_t(du, dx)\pi_t(u, x)(da)$, where η'_t is the marginal of η_t onto $[0, 1] \times \mathbb{R}^d$. Note that actually $\eta'_t(du, dx) = \mu_t$, since for any measurable $F \subset [0, 1] \times \mathbb{R}^d$,

$$\begin{aligned} \eta'_t(F) &= \eta_t(F \times A) = \mathbb{E}^{R^1} \left[\int_A 1_F(U_t^1, X_t^1) 1_A(a) \Lambda_t(da) \right] \\ &= \mathbb{E}^{R^1} [1_F(U_t^1, X_t^1)] = \langle R^1 \circ (U, X^1)^{-1}, 1_F \rangle \end{aligned}$$

Fix $\forall t$, for any bounded measurable $h : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}$,

$$\begin{aligned} &\mathbb{E}^{R^1} \left[h(U^1, X_t^1) \int_A \phi(U^1, X_t^1, a) \pi_t(U^1, X_t^1)(da) \right] \\ &= \int_{[0,1] \times \mathbb{R}^d} h(u, x) \int_A \phi(u, x, a) \pi_t(u, x)(da) \eta'_t(du, dx) \\ &= \int_{[0,1] \times \mathbb{R}^d \times A} h(u, x) \phi(u, x, a) \eta_t(du, dx, da) \\ &= \mathbb{E}^{R^1} \left[h(U^1, X_t^1) \int_A \phi(U^1, X_t^1, a) \Lambda_t(da) \right] \end{aligned}$$

1265 By definition of conditional expectation, the claim follows.

1266 Construct another probability space $(\Omega^2, \mathcal{F}^2, R^2)$ as follows: Let $\Omega^2 = [0, 1] \times \mathcal{C}$, U^2 and X^2 are the coordinate maps, and

$$1268 \quad 1269 \quad R^2 := R^1 \circ (U^1, X^1)^{-1} = \mu$$

1270 In the rest of the proof, we aim to show that U^2, X^2 follows the dynamic

$$1272 \quad 1273 \quad (U^2, X_0^2) \sim \lambda \quad X_{t+1}^2 \sim P_t(X_t^2, \mathbf{W}\mu_t(U^2), \pi(t, U^2, X_t^2))$$

1274 Fix any bounded continuous $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$. For any measurable $h : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}_+$,

$$\begin{aligned} 1277 \quad & \mathbb{E}^{R^2} [h(U^2, X_t^2) \psi(X_{t+1}^2)] \\ 1278 \quad &= \mathbb{E}^{R^1} [h(U^1, X_t^1) \psi(X_{t+1}^1)] \\ 1279 \quad &= \mathbb{E}^{R^1} \left[h(U^1, X_t^1) \mathbb{E}^{R^1} \left(\int_A \int_{\mathbb{R}^d} \psi(y) P_t(\mathbf{W}\mu_t(U^1), X_t^1, a)(dy) \Lambda_t(da) \mid U^1, X_t^1 \right) \right] \\ 1280 \quad &= \mathbb{E}^{R^1} \left[h(U^1, X_t^1) \int_A \int_{\mathbb{R}^d} \psi(y) P_t(\mathbf{W}\mu_t(U^1), X_t^1, a)(dy) \pi(t, U^1, X_t^1)(da) \right] \\ 1282 \quad &= \mathbb{E}^{R^2} \left[h(U^2, X_t^2) \int_A \int_{\mathbb{R}^d} \psi(y) P_t(\mathbf{W}\mu_t(U^2), X_t^2, a)(dy) \pi(t, U^2, X_t^2)(da) \right] \end{aligned}$$

1287 By definition of conditional expectation, we claim that

$$1289 \quad 1290 \quad \mathbb{E}^{R^2} [\psi(X_{t+1}^2) \mid U^2, X_t^2] = \int_A \int_{\mathbb{R}^d} \psi(y) P(t, U^2, \mathbf{W}\mu_t(U^2), X_t^2, a)(dy) \pi(t, U^2, X_t^2)(da)$$

1292 and note that this holds for any ψ bounded continuous. Finally, let $R^0 := R^2 \circ (\{\pi_t(U^2, X_t^2)\}_{t \in T}, U^2, X^2)^{-1}$, then
1293 $R^0 \in \mathcal{R}(\mu)$, and the objective value is preserved:

$$\begin{aligned} 1295 \quad & \int_{\Omega} \Xi^{\mu} dR^0 = \mathbb{E}^{R^2} \left[\sum_{t \in T} \int_A f(t, U^2, X_t^2, \mathbf{W}\mu_t(U^2), a) \pi(t, U^2, X_t^2)(da) + g(X_T^2, \mathbf{W}\mu_T(U^2)) \right] \\ 1296 \quad &= \mathbb{E}^{R^1} \left[\sum_{t \in T} \int_A f(t, U^1, X_t^1, \mathbf{W}\mu_t(U^1), a) \pi(t, U^1, X_t^1)(da) + g(X_T^1, \mathbf{W}\mu_T(U^1)) \right] \\ 1297 \quad &= \mathbb{E}^{R^1} \left[\sum_{t \in T} \int_A f(t, U^1, X_t^1, \mathbf{W}\mu_t(U^1), a) \Lambda_t(da) + g(X_T^1, \mathbf{W}\mu_T(U^1)) \right] \\ 1298 \quad &= \int_{\Omega} \Xi^{\mu} dR \end{aligned}$$

1306 \square

1307 *Remark D.11.* The proof is closely based on (Lacker, 2015), which utilized a remarkable result called Markovian projection
1309 theorem (or Mimicking theorem), originated from (Brunick & Shreve, 2013). However, the discrete time setting greatly
1310 simplifies the proof and just the definition of conditional expectation would work.

1312 E. Proof for Uniqueness

1314 Let (μ, π) and (ν, ρ) be two different W -equilibrium, and the Markovian state dynamic being X^{π} and X^{ρ} respectively. By
1315 construction the processes π and ρ must be different, since otherwise X^{π} and X^{ν} would be the same, and then μ and ν will
1316 be the same as well. Therefore, by uniqueness of optimal policy, we have

$$1318 \quad J_W(\mu, \pi) - J_W(\mu, \rho) > 0 \quad \text{and} \quad J_W(\nu, \rho) - J_W(\nu, \pi) > 0$$

note that the inequalities are strict. Adding them result in

$$J_W(\mu, \pi) - J_W(\nu, \pi) - (J_W(\mu, \rho) - J_W(\nu, \rho)) > 0 \quad (18)$$

Since the Markovian dynamic does not depend on the measure argument, when the population measure is μ , the dynamic controlled by policy ρ is the same pathwise as X^ρ . This is not true if the assumption is not satisfied, since

$$X_{t+1}^\pi \sim P_t(X_t^\pi, \mathbf{W}\mu_t(U), \pi_t) \quad X_{t+1}^\rho \sim P_t(X_t^\rho, \mathbf{W}\nu_t(U), \rho_t)$$

and under population measure μ , the process controlled by ρ follows the dynamic $X'_{t+1} \sim P_t(X'_t, \mathbf{W}\mu_t(U), \rho_t)$, which is not the same as X^ρ . Continue with the proof,

$$\begin{aligned} J_W(\mu, \pi) - J_W(\nu, \pi) &= \mathbb{E} \left[\sum_{t \in \mathbb{T}} (f_t^1(X_t^\pi, \mathbf{W}\mu_t(U)) - f_t^1(X_t^\pi, \mathbf{W}\nu_t(U))) \right. \\ &\quad \left. + \sum_{t \in \mathbb{T}} (f_t^2(X_t^\pi, \pi) - f_t^2(X_t^\pi, \pi)) + g(X_T^\pi, \mathbf{W}\mu_T(U)) - g(X_T^\pi, \mathbf{W}\nu_T(U)) \right] \\ &= \sum_{t \in \mathbb{T}} \int_{[0,1] \times \mathbb{R}^d} [f_t^1(x, \mathbf{W}\mu_t(u)) - f_t^1(x, \mathbf{W}\nu_t(u))] \mu_t(du, dx) \\ &\quad + \int_{[0,1] \times \mathbb{R}^d} [g(x, \mathbf{W}\mu_T(u)) - g(x, \mathbf{W}\nu_T(u))] \mu_T(du, dx) \end{aligned}$$

Similarly,

$$\begin{aligned} J_W(\mu, \rho) - J_W(\nu, \rho) &= \sum_{t \in \mathbb{T}} \int_{[0,1] \times \mathbb{R}^d} [f_t^1(x, \mathbf{W}\mu_t(u)) - f_t^1(x, \mathbf{W}\nu_t(u))] \nu_t(du, dx) \\ &\quad + \int_{[0,1] \times \mathbb{R}^d} [g(x, \mathbf{W}\mu_T(u)) - g(x, \mathbf{W}\nu_T(u))] \nu_T(du, dx) \end{aligned}$$

Taking difference,

$$\begin{aligned} J_W(\mu, \pi) - J_W(\nu, \pi) - (J_W(\mu, \rho) - J_W(\nu, \rho)) &= \sum_{t \in \mathbb{T}} \int_{[0,1] \times \mathbb{R}^d} [f_t^1(x, \mathbf{W}\mu_t(u)) - f_t^1(x, \mathbf{W}\nu_t(u))] (\mu_t - \nu_t)(du, dx) \\ &\quad + \int_{[0,1] \times \mathbb{R}^d} [g(x, \mathbf{W}\mu_T(u)) - g(x, \mathbf{W}\nu_T(u))] (\mu_T - \nu_T)(du, dx) \\ &\leq 0 \end{aligned}$$

by the assumed Larys-Lions monotonicity. However this contradicts (18), and we conclude that μ and ν should be the same.

F. Proof for Approximate Equilibrium

F.1. Comparable dynamics

Define $I_i^n := [(i-1)/n, i/n]$ for $i = 1, \dots, n-1$, $I_n^n := [(n-1)/n, 1]$, and $\mathbf{I}^n := I_1^n \times \dots \times I_n^n$. Let (μ, π) be a W -equilibrium, and X is the Markov chain controlled by policy π . Let X^u denote the state process conditional on $U = u$.

Fix $\forall n \in \mathbb{N}$, and any $\mathbf{u}^n = (u_1^n, \dots, u_n^n) \in [0, 1]^n$. Assign player i the policy

$$\hat{\pi}^{n, \mathbf{u}^n, i}(t, x_1, \dots, x_n) := \pi(t, u_i^n, x_i)$$

and let $\hat{\mathbf{X}}^{n, \mathbf{u}^n} = (\hat{X}^{n, \mathbf{u}^n, 1}, \dots, \hat{X}^{n, \mathbf{u}^n, n})$ be the state dynamic of all the players

$$\hat{X}_{t+1}^{n, \mathbf{u}^n, i} \sim P_t(\hat{X}_t^{n, \mathbf{u}^n, i}, \hat{M}_t^{n, \mathbf{u}^n, i}, \hat{\pi}_t^{n, \mathbf{u}^n, i}(\hat{\mathbf{X}}_t^{n, \mathbf{u}^n})) \quad \hat{X}_0^{n, \mathbf{u}^n, i} \sim \lambda_{u_i^n}$$

1375 where

$$1377 \quad \widehat{M}_t^{n,\mathbf{u}^n,i} := \frac{1}{n} \sum_{r=1}^n \xi_{ir}^n \delta_{\widehat{X}_t^{n,\mathbf{u}^n,r}}$$

1380 and $\widehat{M}_t^{n,\mathbf{u}^n,i}$ is the time t marginal. Let $\widehat{X}_t^{n,\mathbf{u}^n,\beta,j}$ denote the dynamic of player j when she change her policy from $\widehat{\pi}^{n,\mathbf{u}^n,j}$
 1381 to β . More specifically, player j follows

$$1383 \quad \widehat{X}_{t+1}^{n,\mathbf{u}^n,\beta,j} \sim P_t(\widehat{X}_t^{n,\mathbf{u}^n,j}, \widehat{M}_t^{n,\mathbf{u}^n,(\beta,j)}, \beta_t) \quad \widehat{X}_0^{n,\mathbf{u}^n,\beta,j} \sim \lambda_{u_j^n}$$

1385 and all other player $i \neq j$ follows

$$1386 \quad \widehat{X}_{t+1}^{n,\mathbf{u}^n,i} \sim P_t(\widehat{X}_t^{n,\mathbf{u}^n,i}, \widehat{M}_t^{n,\mathbf{u}^n,(\beta,j),i}, \widehat{\pi}_t^{n,\mathbf{u}^n,i}(\mathbf{X}_t^{n,\mathbf{u}^n,\beta,j})) \quad \widehat{X}_0^{n,\mathbf{u}^n,i} \sim \lambda_{u_i^n}$$

1388 where the empirical neighborhood measure is

$$1390 \quad \widehat{M}^{n,\mathbf{u}^n,(\beta,j),i} := \frac{1}{n} \left(\sum_{r \neq j} \xi_{ir}^n \delta_{\widehat{X}_t^{n,\mathbf{u}^n,r}} + \xi_{ij}^n \delta_{\widehat{X}_t^{n,\mathbf{u}^n,\beta,j}} \right)$$

1394 and $\mathbf{X}^{n,\mathbf{u}^n,\beta,j}$ denotes the vector $\mathbf{X}^{n,\mathbf{u}^n}$ with the j^{th} element replaced by $\widehat{X}^{n,\mathbf{u}^n,\beta,j}$.

1395 For $\forall u \in [0, 1]$, we define $X^{\pi,u}$ to be the process with marginal $U = u$, controloed by policy π , i.e.,

$$1397 \quad X_{t+1}^{\pi,u} \sim P_t(X_t^{\pi,u}, \mathbf{W}\mu_t(u), \pi(t, u, X_t^{\pi,u})) \quad X_0^{\pi,u} \sim \lambda_u$$

1400 **Proposition F.1.** Assume assumption 4.7 holds. Let $h : [0, 1] \times \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d) \rightarrow \mathbb{R}$ be a bounded measurable function
 1401 such that $h(u, \cdot, \cdot)$ is jointly continuous on $\mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d)$ for each fixed u . Then for each $t \in \mathbb{T}$,

$$1403 \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, \widehat{X}_t^{n,\mathbf{U}^n,i}, \widehat{M}_t^{n,\mathbf{U}^n,i})] \rightarrow \mathbb{E}[h(U, X_t, \mathbf{W}\mu_t(U))] \quad (19)$$

1406 *Proof.* Expand the underlying probability space such that it supports independent random elements $(U_i^n, Y^{n,i})$, $\forall i \in [n]$,
 1407 independent of $\widehat{\mathbf{X}}^{n,\mathbf{u}^n}$ and (U, X) , and the law satisfies

$$1409 \quad \mathcal{L}(Y^{n,i}|U_i^n = u) = \mathcal{L}(X|U = u) \quad \forall u \in I_i^n$$

1411 equivalently, this means for $\forall u \in I_i^n$, the conditional law satisfies

$$1413 \quad Y_{t+1}^{n,i}|(U_i^n = u) \sim P(t, Y_t^{n,i}, \mathbf{W}\mu_t(u), \pi_t(u, Y_t^{n,i}))$$

1414 In particular for every measurable $\phi : [0, 1] \times \mathcal{C} \rightarrow \mathbb{R}$,

$$1416 \quad \langle \mu, \phi \rangle = \mathbb{E}\phi(U, X) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}\phi(U_i^n, Y^{n,i}) \quad (20)$$

1419 Define empirical neighborhood measure

$$1421 \quad M^{n,i} := \frac{1}{n} \sum_{j=1}^n \xi_{ij}^n \delta_{Y^{n,j}} = \frac{1}{n} \sum_{j=1}^n W_{\xi^n}(U_i^n, U_j^n) \delta_{Y^{n,j}}$$

1424 and empirical label-state joint measure

$$1427 \quad \mu^n := \frac{1}{n} \sum_{j=1}^n \delta_{(U_i^n, Y^{n,i})}$$

1430 The theorem is then shown in the following two stages:

$$1432 \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \hat{M}_t^{n,\mathbf{U}^n,i})] \rightarrow \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, Y_t^{n,i}, M_t^{n,i})] \rightarrow \mathbb{E}[h(U, X_t, \mathbf{W}\mu_t(U))] \\ 1433 \\ 1434$$

1435
 1436 **Step i.** We first show that $\mathbf{W}_{\xi^n}\mu^n(U) \Rightarrow \mathbf{W}\mu(U)$ in probability. Fix a bounded continuous function $\phi : \mathbb{R}^d \rightarrow [-1, 1]$,
 1437 it suffices to show $\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle \rightarrow \langle \mathbf{W}\mu(U), \phi \rangle$ in probability. This is divided into two substeps. We first claim that
 1438 $\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle - \mathbb{E}[\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U] \rightarrow 0$ in probability. Note that
 1439

$$1440 \quad \langle \mathbf{W}_{\xi^n}\mu^n(u), \phi \rangle = \frac{1}{n} \sum_{j=1}^n W_{\xi^n}(u, U_j^n) \phi(Y^{n,i}) \\ 1441 \\ 1442$$

1443 For $u \in I_i^n$, by independence of $Y^{n,i}$,

$$1446 \quad \text{var}(\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U = u) = \text{var}\left(\frac{1}{n} \sum_{j=1}^n \xi_{ij}^n \phi(Y^{n,i})\right) \leq \frac{1}{n^2} \sum_{j=1}^n (\xi_{ij}^n)^2 \\ 1447 \\ 1448$$

1449 Then

$$1451 \quad \mathbb{E}[(\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle - \mathbb{E}[\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U])^2] \\ 1452 \\ 1453 \\ 1454 \\ 1455 \\ 1456 \\ 1457 \\ 1458 \\ 1459$$

$$\begin{aligned} &= \mathbb{E}[\text{var}(\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U)] \\ &= \sum_{i=1}^n \int_{I_i^n} \text{var}(\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U = u) du \\ &\leq \frac{1}{n^3} \sum_{i,j=1}^n (\xi_{ij}^n)^2 \rightarrow 0 \end{aligned}$$

1460 by assumption in equation 10, thus the convergence is in L^2 . In the second substep we show that $\mathbb{E}[\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U] \rightarrow$
 1461 $\langle \mathbf{W}\mu(U), \phi \rangle$ in probability. By independence of $(U_i^n, Y^{n,i})$,
 1462

$$1463 \\ 1464 \quad \mathbb{E}[\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U = u] = \mathbb{E}\left[\frac{1}{n} \sum_{j=1}^n W_{\xi^n}(u, U_j^n) \phi(Y^{n,i})\right] \\ 1465 \\ 1466 \\ 1467 \\ 1468 \\ 1469 \\ 1470$$

$$\begin{aligned} &= \mathbb{E}[W_{\xi^n}(u, U) \phi(X)] \\ &= \int_{[0,1]} W_{\xi^n}(u, v) \mathbb{E}[\phi(X) | U = v] dv \end{aligned}$$

1471 where we used identity (20). Similarly

$$1472 \\ 1473 \quad \langle \mathbf{W}\mu(u), \phi \rangle = \mathbb{E}[W(u, U) \phi(X)] = \int_{[0,1]} W(u, v) \mathbb{E}[\phi(X) | U = v] dv \\ 1474$$

1475 Thus

$$1477 \\ 1478 \quad \mathbb{E}[|\mathbb{E}[\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U] - \langle \mathbf{W}\mu(U), \phi \rangle|] = \int_{[0,1]} \left| \int_{[0,1]} (W_{\xi^n}(u, v) - W(u, v)) \mathbb{E}[\phi(X) | U = v] dv \right| du \\ 1479 \\ 1480 \\ 1481$$

$$= \|(\mathbf{W}_{\xi^n} - \mathbf{W})\phi\|_{L^1[0,1]}$$

1482 By the assumption that $W_{\xi^n} \rightarrow W$ in the strong operator topology, the right-hand side goes to 0 and thus
 1483 $\mathbb{E}[\langle \mathbf{W}_{\xi^n}\mu^n(U), \phi \rangle | U] \rightarrow \langle \mathbf{W}\mu(U), \phi \rangle$ in L^1 . This conclude the first step.
 1484

1485 Step ii. We next show by induction the following holds for each $t \in \mathbb{T}$:

$$1487 \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \hat{M}_t^{n,\mathbf{U}^n,i})] \rightarrow \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, Y_t^{n,i}, M_t^{n,i})] \quad (21)$$

1490 This is trivially true at time 0, since $\hat{X}_0^{n,\mathbf{U}^n,i}$ are initialized independently, we have $\mathcal{L}(U_i^n, \hat{X}_0^{n,\mathbf{U}^n,i}) = \mathcal{L}(U_i^n, Y_0^{n,i})$, and
1491 thus

$$1493 \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, \hat{X}_0^{n,\mathbf{U}^n,i}, \hat{M}_0^{n,\mathbf{U}^n,i})] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[h(U_i^n, Y_0^{n,i}, M_0^{n,i})]$$

1496 Now assume (21) holds for time $t - 1$. We have

$$\begin{aligned} 1498 & \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \hat{M}_t^{n,\mathbf{U}^n,i})] - \mathbb{E}[h(U_i^n, Y_t^{n,i}, M_t^{n,i})] \right) \\ 1499 & \leq \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \hat{M}_t^{n,\mathbf{U}^n,i})] - \mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, M_t^{n,i})]) \\ 1500 & \quad + \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, M_t^{n,i})] - \mathbb{E}[h(U_i^n, Y_t^{n,i}, M_t^{n,i})]) \\ 1503 & = \text{I} + \text{II} \end{aligned}$$

1508 Denote $\mathcal{F}_t^n := \sigma(\{U_i^n\}_{i=1}^n, \{\hat{X}_s^{n,\mathbf{U}^n,i}\}_{i=1}^n, \{Y_s^{n,i}\}_{i=1}^n, s \leq t)$. For term I, we note that $\hat{X}_t^{n,\mathbf{U}^n,i}$ and $\hat{X}_t^{n,\mathbf{U}^n,j}$ are independent conditional on \mathcal{F}_{t-1}^n , and

$$\begin{aligned} 1511 & \mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \hat{M}_t^{n,\mathbf{U}^n,i}) - h(u, \hat{X}_t^{n,\mathbf{U}^n,i}, M_t^{n,i}) \mid \mathcal{F}_{t-1}^n, \hat{X}_t^{n,\mathbf{U}^n,i}] \\ 1512 & = \int_{(\mathbb{R}^d)^{n-1}} h\left(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \frac{1}{n} \sum_{j=1}^n \xi_{ij}^n \delta_{x^j}\right) \prod_{j \neq i} \hat{P}_{t-1}^{n,\mathbf{U}^n,j}(dx^j) \\ 1513 & \quad - \int_{(\mathbb{R}^d)^{n-1}} h\left(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \frac{1}{n} \sum_{j=1}^n \xi_{ij}^n \delta_{y^j}\right) \prod_{j \neq i} P_{t-1}^{n,j}(dy^j) \end{aligned}$$

1518 where we use a shorthand notation

$$\begin{aligned} 1520 & \hat{P}_s^{n,\mathbf{U}^n,i} := P_s(\hat{X}_s^{n,\mathbf{U}^n,i}, \hat{M}_s^{n,\mathbf{U}^n,i}, \pi_s(U_i^n, \hat{X}_s^{n,\mathbf{U}^n,i})) \\ 1521 & P_s^{n,i} := P_s(Y_s^{n,i}, \mathbf{W}\mu_s(U_i^n), \pi_s(U_i^n, Y_s^{n,i})) \end{aligned}$$

1523 More specifically, define the function $h' : [0, 1] \times \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d)$,

$$1526 \quad h'(u, x, m) := \int_{(\mathbb{R}^d)^{n-1}} h\left(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \frac{1}{n} \sum_{j=1}^n \xi_{ij}^n \delta_{x^j}\right) \prod_{j \neq i} P_{t-1}(x, m, \pi_{t-1}(u, x))(dx^j)$$

1528 Then

$$1531 \quad \text{I} = \frac{1}{n} \sum_{i=1}^n \left(h'(U_i^n, \hat{X}_{t-1}^{n,\mathbf{U}^n,i}, \hat{M}_{t-1}^{n,\mathbf{U}^n,i}) - h'(U_i^n, Y_{t-1}^{n,i}, \mathbf{W}\mu_{t-1}(U_i^n)) \right)$$

1533 Similarly for II, note that $Y^{n,i}$ are independent,

$$\begin{aligned} 1536 & \mathbb{E}[h(u, \hat{X}_t^{n,\mathbf{U}^n,i}, M_t^{n,i}) - h(U_i^n, Y_t^{n,i}, M_t^{n,i}) \mid \mathcal{F}_{t-1}^n, Y_t^{n,j}, j \neq i] \\ 1537 & = \int_{(\mathbb{R}^d)^{n-1}} h(U_i^n, x^i, M_t^{n,i}) \hat{P}_{t-1}^{n,\mathbf{U}^n,i}(dx^i) - h(U_i^n, y^i, M_t^{n,i}) P_{t-1}^{n,i}(dy^i) \end{aligned}$$

1540 by defining the function $h'': [0, 1] \times \mathbb{R}^d \times \mathcal{M}_+(\mathbb{R}^d)$,

$$1542 h''(u, x, m) := \int_{(\mathbb{R}^d)^{n-1}} h(u, x^i, M_t^{n,i}) P_{t-1}(x, m, \pi_{t-1}(u, x))(dx^i)$$

1544 and

$$1546 \text{II} \leq \frac{1}{n} \sum_{i=1}^n \left(h''(U_i^n, \hat{X}_{t-1}^{n,\mathbf{U}^n,i}, \hat{M}_{t-1}^{n,\mathbf{U}^n,i}) - h''(U_i^n, Y_{t-1}^{n,i}, \mathbf{W}\mu_{t-1}(U_i^n)) \right)$$

1549 Note that by the assumption that h and P are continuous, $h'(t, \cdot, \cdot)$ and $h''(t, \cdot, \cdot)$ are jointly continuous for every $t \in \mathbb{T}$.
 1550 Combining I and II, by tower property,

$$\begin{aligned} 1552 & \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E}[h(U_i^n, \hat{X}_t^{n,\mathbf{U}^n,i}, \hat{M}_t^{n,\mathbf{U}^n,i})] - \mathbb{E}[h(U_i^n, Y_t^{n,i}, M_t^{n,i})] \right) \\ 1553 & \leq \text{I} + \text{II} \\ 1554 & \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(h' + h'')(U_i^n, \hat{X}_{t-1}^{n,\mathbf{U}^n,i}, \hat{M}_{t-1}^{n,\mathbf{U}^n,i})] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(h' + h'')(U_i^n, Y_{t-1}^{n,i}, M_{t-1}^{n,i})] \\ 1556 & + \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(h' + h'')(U_i^n, Y_{t-1}^{n,i}, M_{t-1}^{n,i})] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(h' + h'')(U_i^n, Y_{t-1}^{n,i}, \mathbf{W}\mu_{t-1}(U_i^n))] \\ 1557 & = \text{I}' + \text{II}' \end{aligned}$$

1563 by our assumption on time $t-1$, $\text{I}' \rightarrow 0$. In step (i) we proved that $\mathbf{W}\xi^n\mu^n(U) \rightarrow \mathbf{W}\mu(U)$. It is straightforward
 1564 to show $\mathbf{W}\xi^n\mu^n(U_i^n) \rightarrow \mathbf{W}\mu(U_i^n)$ with the same line of reasoning. Rewrite $\mathbf{W}\xi^n\mu^n(U_i^n) = M^{n,i}$, we actually have
 1565 $M^{n,i} \rightarrow \mathbf{W}\mu(U_i^n)$ in probability. Combined with the boundedness of integrand, the convergences is in L^1 and thus $\text{II}' \rightarrow 0$.

1566 **Step iii.** Finally, we aim to show,

$$1568 \frac{1}{n} \sum_{i=1}^n \mathbb{E}h(U_i^n, Y_t^{n,i}, M_t^{n,i}) \rightarrow \mathbb{E}[h(U, X, \mathbf{W}\mu(U))]$$

1571 This is justified with similar argument as in (Lacker & Soret, 2022) theorem 6.1, and this concludes the theorem. \square

F.2. Proof of theorem 4.8

1575 Recall the definition of $\epsilon^n(\mathbf{u}^n)$ in definition 3.1, we have

$$\begin{aligned} 1576 \epsilon_i^n(\mathbf{u}^n) &:= \sup_{\beta \in \mathcal{A}_n} J_i(\pi^{n,\mathbf{u}^n,1}, \dots, \pi^{n,\mathbf{u}^n,i-1}, \beta, \pi^{n,\mathbf{u}^n,i+1}, \dots, \pi^{n,\mathbf{u}^n,n}) - J_i(\pi^{n,\mathbf{u}^n}) \\ 1577 &\leq \sup_{\beta \in \mathcal{A}_n} \Delta_1^{n,i}(\beta, \mathbf{u}^n) + \sup_{\beta \in \mathcal{A}_n} \Delta_2^{n,i}(\beta, \mathbf{u}^n) + \sup_{\beta \in \mathcal{A}_n} \Delta_3^{n,i}(\beta, \mathbf{u}^n) + \Delta_4^{n,i}(\mathbf{u}^n) + \Delta_5^{n,i}(\mathbf{u}^n) \end{aligned}$$

1581 where

$$\begin{aligned} 1582 \Delta_1^{n,i}(\beta, \mathbf{u}^n) &:= \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, \hat{X}_t^{n,\mathbf{u}^n,\beta,i}, \hat{M}_t^{n,\mathbf{u}^n,(\beta,i),i}, \beta_t) + g(\hat{X}_T^{n,\mathbf{u}^n,\beta,i}, \hat{M}_T^{n,\mathbf{u}^n,(\beta,i),i}) \right] \\ 1583 &\quad - \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, \hat{X}_t^{n,\mathbf{u}^n,\beta,i}, \mathbf{W}\mu_t(u_i^n), \beta_t) + g(\hat{X}_T^{n,\mathbf{u}^n,\beta,i}, \mathbf{W}\mu_T(u_i^n)) \right] \\ 1584 \\ 1585 \Delta_2^{n,i}(\beta, \mathbf{u}^n) &:= \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, \hat{X}_t^{n,\mathbf{u}^n,\beta,i}, \mathbf{W}\mu_t(u_i^n), \beta_t) + g(\hat{X}_T^{n,\mathbf{u}^n,\beta,i}, \mathbf{W}\mu_T(u_i^n)) \right] \\ 1586 &\quad - \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, X_t^{\beta,u_i^n}, \mathbf{W}\mu_t(u_i^n), \beta_t) + g(X_T^{\beta,u_i^n}, \mathbf{W}\mu_T(u_i^n)) \right] \end{aligned}$$

$$\begin{aligned}
 \Delta_3^{n,i}(\beta, \mathbf{u}^n) &:= \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, X_t^{\beta, u_i^n}, \mathbf{W}_{\mu_t}(u_i^n), \beta_t) + g(X_T^{\beta, u_i^n}, \mathbf{W}_{\mu_T}(u_i^n)) \right] \\
 &\quad - \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, X_t^{u_i^n}, \mathbf{W}_{\mu_t}(u_i^n), \pi_t(u_i^n, \hat{X}_t^{n, \mathbf{u}^n, i})) + g(X_t^{u_i^n}, \mathbf{W}_{\mu_T}(u_i^n)) \right] \\
 \Delta_4^{n,i}(\mathbf{u}^n) &:= \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, X_t^{u_i^n}, \mathbf{W}_{\mu_t}(u_i^n), \pi_t(u_i^n, \hat{X}_t^{n, \mathbf{u}^n, i})) + g(X_t^{u_i^n}, \mathbf{W}_{\mu_T}(u_i^n)) \right] \\
 &\quad - \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, \hat{X}_t^{n, \mathbf{u}^n, i}, \mathbf{W}_{\mu_t}(u_i^n), \pi_t(u_i^n, \hat{X}_t^{n, \mathbf{u}^n, i})) + g(\hat{X}_T^{n, \mathbf{u}^n, i}, \mathbf{W}_{\mu_T}(u_i^n)) \right] \\
 \Delta_5^{n,i}(\mathbf{u}^n) &:= \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, \hat{X}_t^{n, \mathbf{u}^n, i}, \mathbf{W}_{\mu_t}(u_i^n), \pi_t(u_i^n, \hat{X}_t^{n, \mathbf{u}^n, i})) + g(\hat{X}_T^{n, \mathbf{u}^n, i}, \mathbf{W}_{\mu_T}(u_i^n)) \right] \\
 &\quad - \mathbb{E} \left[\sum_{t \in \mathbb{T}} f^i(t, \hat{X}_t^{n, \mathbf{u}^n, i}, \hat{M}_t^{n, \mathbf{u}^n, i}, \pi_t(u_i^n, \hat{X}_t^{n, \mathbf{u}^n, i})) + g(\hat{X}_T^{n, \mathbf{u}^n, i}, \hat{M}_T^{n, \mathbf{u}^n, i}) \right]
 \end{aligned}$$

and the β_t in these formula are short for $\beta_t(\hat{\mathbf{X}}_t^{n, \mathbf{u}^n, \beta, i})$, for closed-loop control $\beta : \mathbb{T} \times \mathbb{R}^d \rightarrow \mathcal{P}(A)$.

Lemma F.2. [Lemma 5.1 of (Lacker & Soret, 2022)] Fix $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{C})$, $u \in [0, 1]$. For any policy $\pi \in \mathcal{A}_1$ and $m \in \mathcal{P}(\mathbb{R}^d)$, define

$$X_{t+1}^{m, \pi} \sim P(t, X_t^{m, \pi}, \mathbf{W}_{\mu_t}(u), \pi(t, X_t^{m, \pi})) \quad X_0^{m, \pi} \sim m$$

and

$$J_W^{u, m}(\mu, \pi) := \mathbb{E} \left[\sum_{t \in \mathbb{T}} f(t, X_t^{m, \pi}, \mathbf{W}_{\mu_t}(u), \pi(t, X_t^{m, \pi})) + g(X_T^{m, \pi}, \mathbf{W}_{\mu_T}(u)) \right]$$

If $\pi \in \mathcal{A}_U$ is an optimal policy, in the sense that $J_W(\mu, \pi) \geq J_W(\mu, \beta)$ for all $\beta \in \mathcal{A}_U$, then

$$J_W^{u, \lambda_u}(\mu, \pi_u) = \sup_{\beta \in \mathcal{A}_1} J_W^{u, \lambda_u}(\mu, \beta) \tag{22}$$

where $\pi_u(t, x) := \pi(t, u, x)$.

Remark F.3. With similar notation as in the proof for equilibrium existence (Appendix D.2), we may denote

$$\mathcal{R}_{u, \lambda_u}(\mu) := \{R \in \mathcal{R}(\mu) \subset \mathcal{P}(\mathcal{V} \times [0, 1] \times \mathcal{C}) : R \circ U^{-1} = \delta_{u_i^n}, R \circ X_0^{-1} = \lambda_u\}$$

The joint law $\mathcal{L}(\pi, u, X^{\lambda_u, \pi}) \in \mathcal{R}_{u, \lambda_u}(\mu)$, and for any $\beta \in \mathcal{V}_U$, $\mathcal{L}(\beta, u, X^{\lambda_u, \beta}) \in \mathcal{R}_{u, \lambda_u}(\mu)$, note that β can be any open-loop policy. Thus equation (22) can be rewritten as³

$$\langle \mathcal{L}(\pi, u, X^{\lambda_u, \pi}), \Xi^\mu \rangle \geq \langle R, \Xi^\mu \rangle \quad \forall R \in \mathcal{R}_{u, \lambda_u}(\mu)$$

where Ξ^μ is defined in equation (15). This view simplifies the analysis of the following lemma.

Lemma F.4. $\sup_{\beta \in \mathcal{A}_n} \Delta_3^{n,i}(\beta, \mathbf{u}^n) \leq 0$ for a.e. $\mathbf{u}^n \in [0, 1]^n$ and all $i \in [n]$.

³Indeed it might not be directly obvious why $\langle \mathcal{L}(\pi, u, X^{\lambda_u, \pi}), \Xi^\mu \rangle \geq \langle R, \Xi^\mu \rangle$ holds for all $R \in \mathcal{R}_{u, \lambda_u}(\mu)$, since R might induce open-loop policies while the supremum in equation (22) is over \mathcal{A}_1 . This actually can be showed rigorously, however, and the reader may refer to Lemma 5.1 of (Lacker & Soret, 2022) for a proof.

1650 *Proof.* By construction, $\mathcal{L}(X^{u_i^n}) = \mathcal{L}(X^{\lambda_{u_i^n}, \pi_{u_i^n}})$, then the second term of $\Delta_3^{n,i}(\beta, \mathbf{u}^n)$ is actually $J_W^{u_i^n, \lambda_{u_i^n}}(\mu, \pi_u)$. On the
 1651 other hand, the joint law $\mathcal{L}(\beta, u_i^n, X^{\beta, u_i^n}) \in \mathcal{R}_{u_i^n, \lambda_{u_i^n}}(\mu)$. Thus by remark F.3, we deduce
 1652

$$1654 \sup_{\beta \in \mathcal{A}_n} \Delta_3^{n,i}(\beta, \mathbf{u}^n) \leq \sup_{\beta \in \mathcal{A}_1} J_W^{u_i^n, \lambda_{u_i^n}}(\mu, \pi_u^*) - J_W^{u_i^n, \lambda_{u_i^n}}(\mu, \pi_u)$$

1655 and following lemma F.2 equation (22), this is ≤ 0 for a.e. $\mathbf{u}^n \in [0, 1]^n$ and all $i \in [n]$. \square
 1656

1657 Take average, we have
 1658

$$1662 \frac{1}{n} \sum_{i=1}^n \epsilon_i^n(\mathbf{u}^n) \leq \frac{1}{n} \sum_{i=1}^n \sup_{\beta \in \mathcal{A}_n} \Delta_1^{n,i}(\beta, \mathbf{u}^n) + \frac{1}{n} \sum_{i=1}^n \sup_{\beta \in \mathcal{A}_n} \Delta_2^{n,i}(\beta, \mathbf{u}^n) + \frac{1}{n} \sum_{i=1}^n \Delta_4^{n,i}(\mathbf{u}^n) + \frac{1}{n} \sum_{i=1}^n \Delta_5^{n,i}(\mathbf{u}^n) \quad (23)$$

1666 By assumption 4.3, it's straightforward to see that $\{\mathcal{L}(\widehat{X}^{n, \mathbf{u}^n, \beta, i}) : (n, \mathbf{u}^n, \beta, i) \in \mathbb{N}_+ \times \mathbf{I}^n \times \mathcal{V} \times [n]\}$ is a tight collection
 1667 of measures in $\mathcal{P}(\mathcal{C})$. Let $K \subset \mathcal{C}$ be a compact subset s.t. $\sup_{n, \mathbf{u}^n, \beta, i} \mathbb{P}(\widehat{X}^{n, \mathbf{u}^n, \beta, i} \notin K) \leq \eta$ for some fixed $\eta > 0$. Define
 1668 function $h_1 : [0, 1] \times \mathcal{M}_+(\mathcal{C}) \rightarrow \mathbb{R}$ by
 1669

$$1671 h_1(u, m) := \sum_{t \in \mathbb{T}} \sup_{a \in A} \sup_{z \in K} (|f(t, z_t, \mathbf{W}\mu_t(u), a) - f(t, z_t, m_t, a)| + |g(z_T, \mathbf{W}\mu_T(u)) - g(z_T, m_T)|)$$

1674 Similarly, define
 1675

$$1677 h_2(u, x) := \sum_{t \in \mathbb{T}} \sup_{a \in A} \sup_{z \in K} (|\mathbb{E}f^i(t, z_t, \mathbf{W}\mu_t(u), a) - f^i(t, x_t, \mathbf{W}\mu_t(u), a)| - |\mathbb{E}g(z_T, \mathbf{W}\mu_T(u)) - g(x_T, \mathbf{W}\mu_T(u))|)$$

1680 Function h_1 and h_2 are bounded measurable since f and g are bounded continuous [(Aliprantis & Border, 2006), theorem
 1681 18.19]. Moreover, it follows from the compactness of A and K that $h_1(u, \cdot)$ is continuous on $\mathcal{M}_+(\mathcal{C})$ for a.e. u , and $h_2(u, \cdot)$
 1682 is continuous on E for a.e. u . Note that $(h_1, h_2)(U, X_t, \mathbf{W}\mu(U)) = 0$. We may thus use h_1 to bound Δ_1 and Δ_5 , use h_2 to
 1683 bound Δ_2 and Δ_4 . To address the region outside K , let C be a constant s.t. $\max(|f|, |g|) \leq C$, and equation (23) becomes
 1684

$$1685 \frac{1}{n} \sum_{i=1}^n \epsilon_i^n(\mathbf{u}^n) \leq \frac{2}{n} \sum_{i=1}^n \mathbb{E} [h_1(u_i^n, \widehat{M}_t^{n, \mathbf{u}^n, i})] + \frac{2}{n} \sum_{i=1}^n \mathbb{E} [h_2(u_i^n, \widehat{X}_t^{n, \mathbf{u}^n, i})] + 8\eta C$$

1688 by Proposition F.1,

$$1691 \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \epsilon_i^n(\mathbf{U}^n) \right] \leq \frac{2}{n} \sum_{i=1}^n \mathbb{E} [h_1(U_i^n, \widehat{M}_t^{n, \mathbf{U}^n, i})] + \frac{2}{n} \sum_{i=1}^n \mathbb{E} [h_2(U_i^n, \widehat{X}_t^{n, \mathbf{U}^n, i})] + 8\eta C$$

$$1694 \longrightarrow 2\mathbb{E} [h_1(U, \mathbf{W}\mu_t(U))] + 2\mathbb{E} [h_2(U, X_t)] + 8\eta C$$

$$1695 = 8\eta C$$

1697 the proof for theorem 4.8 is concluded by letting $\eta \rightarrow 0$.

1700 G. Proof for Online Learning Sample Complexity

1701 G.1. A concrete algorithm realization

1703 For clarity, we present Algorithm 1 combined with subroutine (14) in Algorithm 2.
 1704

1705 **Algorithm 2** Oracle-free Learning for GMFG
 1706 Initialize $Q^{0,0} = \{Q_d^{0,0}\}_{d=1}^D$ and $M^{0,0} = \{M_d^{0,0}\}_{d=1}^D$
 1707 **for** $k \leftarrow 0$ to $K - 1$ **do**
 1708 **for** $d \leftarrow 1$ to D **do**
 1709 Sample initial state $x_0 \sim M_d^{k,0}$, action $a_0 \sim \Gamma_\pi(Q_d^{k,0})$
 1710 **for** $\tau \leftarrow 0$ to H **do**
 1711 Sample reward $r_\tau = f(x_\tau, \mathbf{W}\Pi_D M^{k,0}(u_d), a_\tau)$, next state $x_{\tau+1} \sim P(x_\tau, \mathbf{W}\Pi_D M^{k,0}(u_d), a_\tau)$, and action
 1712 $a_{\tau+1} \sim \Gamma_\pi(Q_d^{k,\tau})$
 1713 Update Q-function:
 1714 $Q_d^{k,\tau+1}(x_\tau, a_\tau) \leftarrow (1 - \alpha_\tau)Q_d^{k,\tau}(x_\tau, a_\tau) + \alpha_\tau (r_\tau + \gamma Q_d^{k,\tau}(x_{\tau+1}, a_{\tau+1}))$
 1715 Update population measure:
 1716 $M_d^{k,\tau+1} \leftarrow (1 - \beta_\tau)M_d^{k,\tau} + \beta_\tau \delta_{x_{\tau+1}}$
 1717 **end for**
 1718 Let $Q_d^{k+1,0} = Q_d^{k,H}$ and $M_d^{k+1,0} = M_d^{k,H}$
 1719 **end for**
 1720 **end for**
 1721 Return policy $\pi^{(K)} := \Gamma_\pi(\Pi_D Q^{K,0})$ and population measure $\mu^{(K)} := \Pi_D M^{K,0}$, where $Q^{K,0} = \{Q_d^{K,0}\}_{d=1}^D$ and
 1722 $M^{K,0} = \{M_d^{K,0}\}_{d=1}^D$

1725
 1726 **Algorithm 3** Oracle-free Learning for GMFG - Finite Horizon
 1727
 1728 Initialize: time horizon T , $Q^{0,0:T} = \{Q_d^{0,0:T}\}_{d=1}^D$ and $M^{0,0:T} = \{M_d^{0,0:T}\}_{d=1}^D$
 1729 **for** $k \leftarrow 0$ to $K - 1$ **do**
 1730 **for** $d \leftarrow 1$ to D **do**
 1731 Sample initial state $x_0 \sim M_d^{k,0}$
 1732 **for** $t \leftarrow 0$ to $T - 1$ **do**
 1733 Choose action a_t from $Q_d^{k,t}(x, .)$
 1734 Sample reward $r_t = f(x_t, \mathbf{W}\Pi_D M^{k,t}(u_d), a_t)$, next state $x_{t+1} \sim P(x_t, \mathbf{W}\Pi_D M^{k,t}(u_d), a_t)$
 1735 Update population measure:
 1736 $M_d^{k,t+1} \leftarrow (1 - \beta_k)M_d^{k,t} + \beta_k \delta_{x_{t+1}}$
 1737 Update Q-function:
 1738 $Q_d^{k,t}(x_t, a_t) \leftarrow (1 - \alpha_k)Q_d^{k,t}(x_t, a_t) + \alpha_k (r_t + \gamma Q_d^{k,t+1}(x_{t+1}, a_{t+1}))$
 1739 **end for**
 1740 **end for**
 1741 **end for**
 1742 Return policy $\pi^{(K)} := \Gamma_\pi(\Pi_D Q^{K,0:T})$ and population measure $\mu^{(K)} := \Pi_D M^{K,0:T}$, where $Q^{K,0:T} = \{Q_d^{K,0:T}\}_{d=1}^D$
 1743 and $M^{K,0:T} = \{M_d^{K,0:T}\}_{d=1}^D$

1744
 1745
 1746 Algo. 3 is adapted from Algo. 2 to solve GMFGs with finite time horizons. The difference between two algorithms lies in
 1747 the learning rate. In Algo. 3, the learning rate has to capture each time step t in the time horizon T . Therefore, we have
 1748 $\beta_k = \frac{1}{1 + \#(t, k)}$ and $\alpha_k = \frac{1}{1 + \#(x, a, t, k)}$, where $\#(t, k)$ counts the number of visits to time step t up to epoch k . $\#(x, a, t, k)$
 1749 counts the number of visits to tuple (x, a, t) up to epoch k .
 1750

G.2. Discretization of label space

1751 Recall that $\mathcal{U} := \{u_1, \dots, u_D\}$ is the discretization of label space, and $\Pi_D : [0, 1] \rightarrow \mathcal{U}$ is the projection mapping. We
 1752 define the operator Π_D which maps operators defined on \mathcal{U} to operators defined on $[0, 1]$: for any operator $\tilde{\phi}$ defined on \mathcal{U} ,
 1753

$$\Pi_D \tilde{\phi}(u) := \sum_{d=1}^D \tilde{\phi}(u_d) 1_{\{u \in I_{u_d}\}},$$

generalizing it to a function on $[0, 1]$. In particular, for $\tilde{\mu} = \{\tilde{\mu}^{u_d}\}_{d=1}^D \in \mathcal{P}(\mathcal{X})^{\mathcal{U}}$, we regard $\Pi_D \tilde{\mu}$ as both the kernel $\nu : [0, 1] \rightarrow \mathcal{P}(\mathcal{X})$ given by

$$\nu(u) := \sum_{d=1}^D \tilde{\mu}^{u_d} 1_{\{u \in I_{u_d}\}}.$$

and also a measure in $\mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$, constructed by $\text{Leb} \otimes \nu$.

In addition to Π_D , we define a set value mapping $\Pi_D^\dagger : \mathcal{U} \rightarrow 2^{[0,1]}$ by $\Pi_D^\dagger(u_d) = I_d$ for any $u_d \in \mathcal{U}$. The operator Π_D^\dagger maps operators defined on $[0, 1]$ to operators defined on \mathcal{U} . For any operator ϕ defined on $[0, 1]$,

$$\Pi_D^\dagger \phi(u_d) := \phi(u_d) \quad u_d \in \mathcal{U}$$

Note that $\Pi_D^\dagger \Pi_D = \text{ID}_{\mathcal{U}}$, while the inverse is not necessarily true.

Lemma G.1. *The operator norm of $\Pi_D : \mathcal{P}(\mathcal{X})^{\mathcal{U}} \rightarrow \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$ is bounded by 1, where we equip the $\mathcal{P}(\mathcal{X})^{\mathcal{U}}$ with the norm $\|\tilde{\mu}\| = \sup_{u_d \in \mathcal{U}} \|\tilde{\mu}^{u_d}\|_{\text{TV}}, \forall \tilde{\mu}$.*

Proof. It holds for any $\tilde{\mu} \in \mathcal{P}(\mathcal{X})^{\mathcal{U}}$ that

$$\begin{aligned} \|\Pi_D \tilde{\mu}\|_{\text{TV}} &= \sup_{\|\phi\|_\infty \leq 1} \left| \int_{[0,1] \times \mathcal{X}} \phi(u, x) \Pi_D \tilde{\mu}(du, dx) \right| \\ &\leq \sum_{d=1}^D \sup_{\|\phi\|_\infty \leq 1} \int_{I_{u_d}} \left| \int_{\mathcal{X}} \phi(u, x) \tilde{\mu}^{u_d}(dx) \right| du \\ &\leq \sum_{d=1}^D \int_{I_{u_d}} \|\tilde{\mu}^{u_d}\|_{\text{TV}} du \leq \sup_{u_d \in \mathcal{U}} \|\tilde{\mu}^{u_d}\|_{\text{TV}} = \|\tilde{\mu}\| \end{aligned}$$

□

The following lemma ensures the discretization \mathcal{U} a good approximation of the label space.

Lemma G.2. *For any $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$, we have*

$$\sup_{u \in [0,1]} \|\mathbf{W}\mu(u) - \mathbf{W}\mu(\Pi_D(u))\|_{\text{TV}} \leq \frac{L_d}{D}$$

Proof. Recall the definition of total variation norm,

$$\begin{aligned} \sup_{u \in [0,1]} \|\mathbf{W}\mu(u) - \mathbf{W}\mu(\Pi_D(u))\|_{\text{TV}} &= \sup_u \sup_{\|\phi\|_\infty \leq 1} \left| \int_{[0,1] \times \mathcal{X}} (W(u, v) - W(\Pi_D(u), v)) \phi(x) \mu(dv, dx) \right| \\ &= \sup_u \int_{[0,1]} \left| (W(u, v) - W(\Pi_D(u), v)) \right| dv \\ &\leq \frac{L_d}{D}. \end{aligned}$$

□

G.3. Best response and induced population operator

Recall \mathcal{Q} is the collection of all $[0, 1] \times \mathcal{X} \times A \rightarrow \mathbb{R}$ functions, for any $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$, the Bellman (optimality) operator $\mathcal{T}_\mu : \mathcal{Q} \rightarrow \mathcal{Q}$ is defined by

$$\mathcal{T}_\mu q(u, x, a) = f(x, \mathbf{W}\mu(u), a) + \gamma \langle P(x, \mathbf{W}\mu(u), a), \sup_{a \in A} q(u, \cdot, a) \rangle$$

for any $q \in \mathcal{Q}$. It is known that \mathcal{T}_μ is a γ -contraction mapping, thus a unique fixed point exists and denote Q^μ the fixed point of \mathcal{T}_μ . Let the value function be $v^\mu(u, x) := \sup_{a \in A} Q^\mu(u, x, a)$.

1815 **BR and IP operator.** The FPI Γ , given by $\Gamma(\mu) = \Gamma_2(\Gamma_1(\mu), \mu)$, can be alternatively decompose into *best response*
 1816 (BR) w.r.t. the current population and the *induced population* (IP) w.r.t. the current policy. Define the BR operator
 1817 $\Gamma_{\text{BR}} : \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \rightarrow \mathcal{Q}$ by $\Gamma_{\text{BR}}(\mu) = Q^\mu$ where

$$1818 \quad Q^\mu(u, x, a) := f(x, \mathbf{W}\mu(u), a) + \langle P(x, \mathbf{W}\mu(u), a), v^\mu(u, \cdot) \rangle$$

1820 The IP operator $\Gamma_{\text{IP}} : \mathcal{Q} \times \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \rightarrow \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$ is defined by $\Gamma_{\text{IP}}(Q, \mu) = \mathcal{L}(U, X)$ where X follows
 1821 the Markov transition with population measure μ , and under policy $\Gamma_\pi(Q)$.

1822 Actually, $\Gamma_\pi \circ \Gamma_{\text{BR}}(\mu) = \Gamma_1(\mu)$, and $\Gamma_{\text{IP}}(Q, \mu) = \Gamma_2(\Gamma_\pi(Q), \mu)$, and we have $\Gamma(\mu) = \Gamma_2(\Gamma_1(\mu), \mu) = \Gamma_{\text{IP}}(\Gamma_{\text{BR}}(\mu), \mu)$.
 1823 However, both Γ_{BR} and Γ_{IP} are defined in terms of \mathcal{Q} , where the label space $[0, 1]$ is continuous. Thus we define the
 1824 following operators with Q-functions on \mathcal{U} .

1826 **Discretized BR and IP operator.** Let $\tilde{\mathcal{Q}}$ be the collection of all $\mathcal{U} \times \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ functions, we define the discretized BR
 1827 operator $\tilde{\Gamma}_{\text{BR}} : \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \rightarrow \tilde{\mathcal{Q}}$ by $\tilde{\Gamma}_{\text{BR}}(\mu) = \tilde{Q}^\mu$ where

$$1829 \quad \tilde{Q}^\mu(u_d, x, a) := f(x, \mathbf{W}\mu(u_d), a) + \langle P(x, \mathbf{W}\mu(u_d), a), v^\mu(u_d, \cdot) \rangle \quad \forall u_d \in \mathcal{U}$$

1830 $\tilde{\Gamma}_{\text{BR}}$ returns D best responses for labels in \mathcal{U} w.r.t. population distribution μ . In particular, \tilde{Q}^μ and Q^μ coincide at $\mathcal{U} \times \mathcal{X} \times \mathcal{A}$.

1832 The discretized IP operator $\tilde{\Gamma}_{\text{IP}} : \tilde{\mathcal{Q}} \times \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \rightarrow \mathcal{P}(\mathcal{X})^{\mathcal{U}}$ is defined by $\tilde{\Gamma}_{\text{IP}}(\tilde{Q}, \mu) = \mathcal{L}(U, X)$ where X follows
 1833 the Markov transition with population measure μ , and under policy $\Gamma_\pi(\tilde{Q})$, conditional on $U \in \mathcal{U}$. In other words, it is the
 1834 induced state distribution on $\mathcal{P}(\mathcal{X})^{\mathcal{U}}$ for the D classes.

1835 For notation simplicity, we denote $\Gamma_{\text{IP}}\Gamma_{\text{BR}}(\mu) = \Gamma_{\text{IP}}(\Gamma_{\text{BR}}(\mu), \mu)$, similarly for $\tilde{\Gamma}_{\text{IP}}\tilde{\Gamma}_{\text{BR}}$.

1837 **Algorithm operator.** Finally, the algorithm operator $\hat{\Gamma} : \mathcal{P}(\mathcal{X})^{\mathcal{U}} \rightarrow \mathcal{P}(\mathcal{X})^{\mathcal{U}}$ defined by

$$1839 \quad \hat{\Gamma} : \{M_d^{k,0}\}_{d=1}^D \mapsto \{M_d^{k,H}\}_{d=1}^D$$

1840 It returns the updated D -class population measure after an outer iteration of Algorithm 2, consisting of H online stochastic
 1841 updates to the D -class Q- and M-value functions.

1843 Given the initial D -class population estimate $M_0 := (M_d^{0,0})_{d=1}^D \in \mathcal{P}(\mathcal{X})^{\mathcal{U}}$, we can express Algorithm 2 as

$$1844 \quad \Pi_D \hat{\Gamma}^K M_0 = \Pi_D \left(\hat{\Gamma} \Pi_D^\dagger \Pi_D \right)^K M_0 = \left(\Pi_D \hat{\Gamma} \Pi_D^\dagger \right)^K \Pi_D M_0. \quad (24)$$

1847 G.4. Sample Complexity Analysis

1848 Our analysis follows the following illustration:

$$1849 \quad \underbrace{\Pi_D \hat{\Gamma}^K M_0}_{\text{Algorithm 2}} \xrightarrow{\text{approximates}} \underbrace{\left(\Pi_D \tilde{\Gamma}_{\text{IP}} \tilde{\Gamma}_{\text{BR}} \right)^K \Pi_D M_0}_{\text{Finite-label FPI}} \xrightarrow{\text{approximates}} \underbrace{\left(\Gamma_{\text{IP}} \Gamma_{\text{BR}} \right)^K \Pi_D M_0}_{\text{FPI}},$$

1853 We are ready to give the one-step approximation error of Algorithm 2.

1854 **Proposition G.3.** For any $\nu \in \Pi_D \mathcal{P}(\mathcal{X})^{\mathcal{U}}$, we have

$$1856 \quad \mathbb{E} \left\| \left(\Gamma_{\text{IP}} \Gamma_{\text{BR}} - \Pi_D \hat{\Gamma} \Pi_D^\dagger \right) \nu \right\|_{\text{TV}}^2 = O \left(\frac{D \log H}{H} + \frac{1}{D^2} \right).$$

1858 *Proof.* Consider the decomposition

$$1860 \quad \mathbb{E} \left\| \left(\Gamma_{\text{IP}} \Gamma_{\text{BR}} - \Pi_D \hat{\Gamma} \Pi_D^\dagger \right) \nu \right\|_{\text{TV}}^2 \leq 3 \mathbb{E} \underbrace{\left\| \Gamma_{\text{IP}} \left(\Gamma_{\text{BR}} - \Pi_D \tilde{\Gamma}_{\text{BR}} \right) \nu \right\|_{\text{TV}}^2}_{G_1} \\ 1861 \quad + 3 \mathbb{E} \underbrace{\left\| \left(\Gamma_{\text{IP}} \Pi_D - \Pi_D \tilde{\Gamma}_{\text{IP}} \right) \tilde{\Gamma}_{\text{BR}} \nu \right\|_{\text{TV}}^2}_{G_2} \\ 1862 \quad + 3 \mathbb{E} \underbrace{\left\| \Pi_D \left(\tilde{\Gamma}_{\text{IP}} \tilde{\Gamma}_{\text{BR}} - \hat{\Gamma} \Pi_D^\dagger \right) \nu \right\|_{\text{TV}}^2}_{G_3}.$$

1870 Note that though the kernel resulting from disintegration is only Lebesgue a.e. defined, we only consider those $\nu \in$
 1871 $\Pi_D \mathcal{P}(\mathcal{X})^{\mathcal{U}}$, i.e. there exists some $M \in \mathcal{P}(\mathcal{X})^{\mathcal{U}}$ such that $\nu = \Pi_D M$, and thus $\Pi_D^\dagger \nu = M$ is unique without ambiguity.

1872 Let $q := \Gamma_{\text{BR}} \nu$ and $\mu := \Gamma_{\text{IP}}(\nu, q) = \Gamma_{\text{IP}} \Gamma_{\text{BR}} \nu$. Similarly let $\tilde{q} := \tilde{\Gamma}_{\text{BR}} \nu$ and $\tilde{\mu} := \Gamma_{\text{IP}}(\nu, \Pi_D \tilde{q}) = \Gamma_{\text{IP}}(\nu, \Pi_D \tilde{\Gamma}_{\text{BR}} \nu)$. To
 1873 distinguish, $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$, $\tilde{\mu} \in \mathcal{P}(\mathcal{X})^{\mathcal{U}}$; $q \in \mathcal{Q}$, $\tilde{q} \in \tilde{\mathcal{Q}}$. Then, we have

$$\begin{aligned} \sqrt{G_1} &= \|\mu - \tilde{\mu}\|_{\text{TV}} \leq \sup_{\|\phi\|_\infty \leq 1} \int_{[0,1]} \left| \int_{\mathcal{X}} \phi(u, x) (\mu_u - \tilde{\mu}_u)(dx) \right| du \\ &\leq \int_{[0,1]} \sup_{u \in [0,1]} \|\mu_u - \tilde{\mu}_u\|_{\text{TV}} du \\ &= \int_{[0,1]} \sup_{u \in [0,1]} \|\mu_u - \tilde{\mu}_u\|_{\text{TV}} du \end{aligned}$$

1884 Since μ^u and $\tilde{\mu}^u$ are the law of process $X|U = u$ with the same transition kernel, by [Anonymous \(2024, Lemma 4\)](#), we
 1885 have for almost every u ,

$$\|\mu_u - \tilde{\mu}_u\|_{\text{TV}} \leq L_\pi \sigma \|q(u, \cdot) - q(\Pi_D(u), \cdot)\|_2 \leq L_\pi \sigma \sqrt{|\mathcal{X}| |A|} \|q(u, \cdot) - q(\Pi_D(u), \cdot)\|_\infty$$

1889 which gives

$$\sup_{u \in [0,1]} \|\mu_u - \tilde{\mu}_u\|_{\text{TV}} \leq L_\pi \sigma \sqrt{|\mathcal{X}| |A|} \|q - \Pi_D \tilde{q}\|_\infty.$$

1893 Therefore, by Lemma [G.6](#), we get

$$G_1 \leq \frac{|\mathcal{X}| |A| L_\pi^2 L_R^2 L_D^2 \sigma^2 (1 - \gamma + \|f\|_\infty)^2}{(1 - \gamma)^4 D^2}.$$

1898 For G_2 , by Lemma [G.5](#), we have

$$G_2 \leq \frac{L_P^2 L_D^2 \sigma^2}{D^2}.$$

1902 And Lemma [G.4](#) gives

$$G_3 = O\left(\frac{D |\mathcal{X}|^2 |A| \|f\|_\infty^2 L_\pi^2 \sigma^2 \log H}{\lambda_{\min}^2 (1 - \gamma)^4 H}\right).$$

1906 Plugging the above bounds on G_1 , G_2 , and G_3 into gives the desired result. \square

1908 Theorem [5.4](#) is immediate after combining Proposition [G.3](#) with the contraction assumption of FPI.

1910 *Proof of Theorem 5.4.* In this proof, we omit the subscript of the total variation norm for simplicity. We denote $M_k =$
 1911 $M^{k,0} = \{M_d^{k,0}\}_{d=1}^D$ and $\mu_k := \Pi_D M_k$ for $k = 0, \dots, K$. Note that $M_k = \widehat{\Gamma}^k M_0$. By Equation [\(24\)](#) and the definition of
 1912 the equilibrium population measure $\widehat{\mu}$, we have

$$\mathbb{E} \|\mu_K - \widehat{\mu}\|^2 = \mathbb{E} \left\| \Pi_D \widehat{\Gamma}^K M_0 - \widehat{\mu} \right\|^2 = \mathbb{E} \left\| \Pi_D \widehat{\Gamma} \Pi_D^\dagger \mu_{K-1} - \Gamma_{\text{IP}} \Gamma_{\text{BR}} \widehat{\mu} \right\|^2, \quad (25)$$

1917 Then, by Young's inequality, we have

$$\begin{aligned} &\mathbb{E} \left\| \Pi_D \widehat{\Gamma} \Pi_D^\dagger \mu_{K-1} - \Gamma_{\text{IP}} \Gamma_{\text{BR}} \widehat{\mu} \right\|^2 \\ &= \mathbb{E} \left\| \left(\Pi_D \widehat{\Gamma} \Pi_D^\dagger - \Gamma_{\text{IP}} \Gamma_{\text{BR}} \right) \mu_{K-1} + \Gamma_{\text{IP}} \Gamma_{\text{BR}} (\mu_{K-1} - \widehat{\mu}) \right\|^2 \\ &= (1 + 1/\kappa) \mathbb{E} \left\| \left(\Pi_D \widehat{\Gamma} \Pi_D^\dagger - \Gamma_{\text{IP}} \Gamma_{\text{BR}} \right) \mu_{K-1} \right\|^2 + (1 + \kappa) \mathbb{E} \|\Gamma_{\text{IP}} \Gamma_{\text{BR}} (\mu_{K-1} - \widehat{\mu})\|^2. \end{aligned}$$

1925 Using Proposition G.3 for the first term and the contracting FPI assumption for the second term, we get
 1926

$$\begin{aligned} 1927 \quad \mathbb{E} \left\| \boldsymbol{\Pi}_D \widehat{\Gamma} \boldsymbol{\Pi}_D^\dagger \mu_{K-1} - \Gamma_{\text{IP}} \Gamma_{\text{BR}} \widehat{\mu} \right\|^2 &\leq (1 + 1/\kappa) \cdot O \left(\frac{D}{H} + \frac{1}{D^2} \right) + (1 + \kappa)(1 - \kappa)^2 \mathbb{E} \|\mu_{K-1} - \widehat{\mu}\|^2 \\ 1928 \quad &\leq \frac{1}{\kappa} \cdot O \left(\frac{D}{H} + \frac{1}{D^2} \right) + (1 - \kappa) \mathbb{E} \|\mu_{K-1} - \widehat{\mu}\|^2. \\ 1929 \quad & \\ 1930 \quad & \\ 1931 \quad & \end{aligned}$$

1932 Recursively applying the above inequality to Equation (25) gives
 1933

$$\begin{aligned} 1934 \quad \mathbb{E} \|\mu_K - \widehat{\mu}\|^2 &\leq (1 - \kappa)^K \mathbb{E} \|\mu_0 - \widehat{\mu}\|^2 + \sum_{k=1}^K (1 - \kappa)^k \frac{1}{\kappa} \cdot O \left(\frac{D}{H} + \frac{1}{D^2} \right) \\ 1935 \quad & \\ 1936 \quad & \\ 1937 \quad = O \left(\exp(-\kappa K) + \frac{D}{\kappa^2 H} + \frac{1}{\kappa^2 D^2} \right). \\ 1938 \quad & \\ 1939 \quad & \\ 1940 \quad \text{Therefore, to find an } \epsilon\text{-approximation equilibrium population measure } \mu_K \text{ such that } \mathbb{E} \|\mu_K - \widehat{\mu}\| \leq \epsilon, \text{ we need at most} \\ 1941 \quad K = O(\kappa^{-1} \log \epsilon^{-1}), \quad D = O(\kappa^{-1} \epsilon^{-1}), \quad H = O(\kappa^{-3} \epsilon^{-3} \log \epsilon^{-1}). \\ 1942 \quad & \\ 1943 \quad & \square \\ 1944 \quad & \\ 1945 \quad \textbf{G.5. Auxiliary Lemmas} \\ 1946 \quad & \\ 1947 \quad \text{The following lemmas address } G_3, G_2, \text{ and } G_1 \text{ in Proposition G.3 respectively.} \\ 1948 \quad \textbf{Lemma G.4} \text{ (Online learning approximation error). Suppose Assumption 5.3 holds. With step sizes of } \alpha_\tau, \beta_\tau \asymp 1/\tau, \text{ for} \\ 1949 \quad \text{any } M \in \mathcal{P}(\mathcal{X})^{\mathcal{U}}, \text{ we have} \\ 1950 \quad \mathbb{E} \left\| \boldsymbol{\Pi}_D \left(\widetilde{\Gamma}_{\text{IP}} \widetilde{\Gamma}_{\text{BR}} \boldsymbol{\Pi}_D - \widehat{\Gamma} \right) M \right\|_{\text{TV}}^2 = O \left(\frac{D \log H}{H} \right). \\ 1951 \quad & \\ 1952 \quad & \\ 1953 \quad \textit{Proof.} \text{ We first denote } \widetilde{\mu} = \{\widetilde{\mu}^{u_d}\}_{d=1}^D := \widetilde{\Gamma}_{\text{IP}} \widetilde{\Gamma}_{\text{BR}} \boldsymbol{\Pi}_D M \text{ and } \widetilde{M} = \{\widetilde{M}^{u_d}\} := \widehat{\Gamma} M. \text{ Then, we know that } \widetilde{\mu}^{u_d} \text{ is the} \\ 1954 \quad \text{sationary distribution of the MDP dynamic with population measure } \boldsymbol{\Pi}_D M, \text{ conditional on } U = u_d \text{ at time 0 (i.e. the} \\ 1955 \quad \text{measure argument of reward function and transition kernel is } W \boldsymbol{\Pi}_D M(u_d), \text{ controlled by policy } \Gamma_\pi(\widetilde{\Gamma}_{\text{BR}} \boldsymbol{\Pi}_D M)(u_d, \cdot), \\ 1956 \quad \text{which is the optimal policy.} \\ 1957 \quad & \\ 1958 \quad \text{This is the same MDP in Algorithm 2 for label } u_d. \text{ Thus, by Anonymous (2024, Lemma 3), for any } u_d \in \mathcal{U}, \text{ we have} \\ 1959 \quad & \\ 1960 \quad \mathbb{E} \left\| \widetilde{\mu}^{u_d} - \widetilde{M}^{u_d} \right\|_2^2 = O \left(\frac{\|f\|_\infty^2 L^2 \sigma^2 |\mathcal{X}| |A| \log H}{\lambda_{\min}^2 (1 - \gamma)^4 H} \right), \\ 1961 \quad & \\ 1962 \quad & \\ 1963 \quad \text{where } \sigma := \hat{n} + c_1 c_2^{\hat{n}} / (1 - c_2), \hat{n} = \lceil \log_{c_2} c_1^{-1} \rceil, \text{ and } \lambda_{\min} \text{ is the lower bound of the probability of visiting any state-action} \\ 1964 \quad \text{pair under the steady distribution induced by any policy and kernel. Therefore, by Lemma G.1, we have} \\ 1965 \quad & \\ 1966 \quad \mathbb{E} \left\| \boldsymbol{\Pi}_D \left(\widetilde{\mu} - \widetilde{M} \right) \right\|_{\text{TV}}^2 \leq \mathbb{E} \left\| \widetilde{\mu} - \widetilde{M} \right\|^2 \leq D \sup_{u_d \in \mathcal{U}} \mathbb{E} \left\| \widetilde{\mu}^{u_d} - \widetilde{M}^{u_d} \right\|_{\text{TV}}^2 \\ 1967 \quad & \\ 1968 \quad \leq D |\mathcal{X}| \sup_{u_d \in \mathcal{U}} \mathbb{E} \left\| \widetilde{\mu}^{u_d} - \widetilde{M}^{u_d} \right\|_2^2 \\ 1969 \quad & \\ 1970 \quad = O \left(\frac{D \|f\|_\infty^2 L^2 \sigma^2 |\mathcal{X}|^2 |A| \log H}{\lambda_{\min}^2 (1 - \gamma)^4 H} \right). \\ 1971 \quad & \\ 1972 \quad & \\ 1973 \quad \text{where the total variation of measure on finite space is equivalent to } l_1 \text{ norm of the density vector.} \quad \square \\ 1974 \quad & \\ 1975 \quad \textbf{Lemma G.5} \text{ (Population discretization error). For any population distribution } \mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X}) \text{ and any } D\text{-class} \\ 1976 \quad Q\text{-value function } \widetilde{q}, \text{ we have} \\ 1977 \quad & \\ 1978 \quad \left\| \boldsymbol{\Pi}_D \widetilde{\Gamma}_{\text{IP}} (\mu, \widetilde{q}) - \Gamma_{\text{IP}} (\mu, \boldsymbol{\Pi}_D \widetilde{q}) \right\|_{\text{TV}} = O \left(\frac{1}{D} \right) \\ 1979 \quad & \\ 1980 \quad & \end{aligned}$$

1980 *Proof.* We first denote $\tilde{\nu} := \tilde{\Gamma}_{\text{IP}}(\mu, \tilde{q})$ and $\nu := \Gamma_{\text{IP}}(\mu, \Pi_D \tilde{q})$.

1981 Let ν admits disintegration $d\nu_u(dx)$. By construction, conditional on $U = u$, for a.e. $u \in I_d$, $\tilde{\nu}^{u_d}$ and ν_u are the invariant
1982 measures of two Markov processes that follow the same policy $\Gamma_\pi(\tilde{q}_{u_d})$, but w.r.t. different neighborhood measure. By
1983 Mitrophanov (2005, Corollary 3.1), we have for a.e. $u \in I_d$,

$$\begin{aligned} \|\tilde{\nu}^{u_d} - \nu_u\|_{\text{TV}} &\leq \sigma \sup_{x,a} \|P(x, \mathbf{W}\mu(u_d), a) - P(x, \mathbf{W}\mu(u), a)\|_{\text{TV}} \\ &\leq \sigma L_P \|\mathbf{W}\mu(u_d) - \mathbf{W}\mu(u)\|_{\text{TV}} \leq \frac{\sigma L_P L_D}{D} \end{aligned}$$

1989 and thus

$$\begin{aligned} \|\Pi_D \tilde{\nu} - \nu\|_{\text{TV}} &= \sup_{\|\phi\|_\infty \leq 1} \left| \int_{[0,1] \times \mathcal{X}} \phi(u, x) (\Pi_D \tilde{\nu} - \nu)(du, dx) \right| \\ &\leq \sum_{d=1}^D \sup_{\|\phi\|_\infty \leq 1} \int_{I_{u_d}} \left| \int_{\mathcal{X}} \phi(u, x) (\tilde{\nu}^{u_d}(dx) - \nu_u(dx)) \right| du \\ &\leq \sum_{d=1}^D \int_{I_{u_d} \times \mathcal{X}} \|\tilde{\nu}^{u_d} - \nu_u\|_{\text{TV}} du \leq \frac{\sigma L_P L_D}{D} \end{aligned}$$

2000 \square

2002 **Lemma G.6** (Population discretization error). *For any population distribution $\mu \in \mathcal{P}_{\text{unif}}([0, 1] \times \mathcal{X})$, let $q_* := \Gamma_{\text{BR}}\mu$ and
2003 $\tilde{q}_* := \tilde{\Gamma}_{\text{BR}}\mu$. We have*

$$\sup_{\mu} \|q_* - \Pi_D \tilde{q}_*\|_\infty = O\left(\frac{1}{D}\right)$$

2008 *Proof.* Denote value function associated with a policy $\rho : [0, 1] \times \mathcal{X} \rightarrow \mathcal{P}(A)$ by

$$v^{\rho_{u_0}}(u_1, u_2, x) = \mathbb{E} \left[\sum_{\tau \geq 0} \gamma^\tau f(X_\tau^{\rho_{u_0}}, \mathbf{W}\mu(u_2), \alpha_\tau^{\rho_{u_0}}) \mid X_0^{\rho_{u_0}} = x, U = u_1 \right]$$

2013 and with a slight abuse of notation, we denote

$$q^{\rho_{u_0}}(u_1, u_2, x, a) := f(x, \mathbf{W}\mu(u_2), a) + \gamma \langle P(x, \mathbf{W}\mu(u_1), a), v^{\rho_{u_0}}(u_1, u_2, \cdot) \rangle$$

2017 where ρ_{u_0} is to fix u_0 as the first argument of ρ . In words, $v^{\rho_{u_0}}(u_1, u_2, x)$ and $q^{\rho_{u_0}}(u_1, u_2, x, a)$ are generalization of
2018 typical value function and Q functions, where the the policy follows label u_0 , state transition follows u_1 , and the reward
2019 follows u_2 .

2020 Note that $q_* \in \mathcal{Q}$, and $\tilde{q}_* \in \tilde{\mathcal{Q}}$. Let $\pi = \Gamma_\pi(q_*)$, and by definition of Γ_{BR} and $\Gamma_{\text{BR}}^{(D)}$ we know that

$$\begin{aligned} q_*(u, x, a) &= q^{\pi_u}(u, u, x, a) \iff v^{\pi_u}(u, u, x) = \sup_{a \in A} q_*(u, x, a) \\ \tilde{q}_*(u_d, x, a) &= q^{\pi_{u_d}}(u_d, u_d, x, a) \iff v^{\pi_{u_d}}(u_d, u_d, x) = \sup_{a \in A} \tilde{q}_*(u_d, x, a) \quad 1 \leq d \leq D \end{aligned}$$

2027 Note that q_* and \tilde{q}_* coincide on the space $\mathcal{U} \times \mathcal{X} \times A$, by definition of Γ_{BR} and $\tilde{\Gamma}_{\text{BR}}$. On $([0, 1] \setminus \mathcal{U}) \times \mathcal{X} \times A$, the Q-function
2028 q_* is strictly larger than $\Pi_D \tilde{q}_*$ by its optimality. With this in mind, we have

$$\begin{aligned} \|q_* - \Pi_D \tilde{q}_*\|_\infty &= \sup_{x,a} \sup_{1 \leq d \leq D} \sup_{u \in I_{u_d}} (q_*(u, x, a) - \Pi_D \tilde{q}_*(u, x, a)) \\ &= \sup_{x,a} \sup_{1 \leq d \leq D} \sup_{u \in I_{u_d}} (q^{\pi_u}(u, u, x, a) - q^{\pi_{u_d}}(u_d, u_d, x, a)). \end{aligned}$$

2035 where

$$\begin{aligned}
 q^{\pi_u}(u, u, x, a) - q^{\pi_{u_d}}(u_d, u_d, x, a) &\leq |q^{\pi_u}(u, u, x, a) - q^{\pi_u}(u, u_d, x, a)| \\
 &\quad + |q^{\pi_u}(u, u_d, x, a) - q^{\pi_u}(u_d, u_d, x, a)| \\
 &\quad + (q^{\pi_u}(u_d, u_d, x, a) - q^{\pi_{u_d}}(u_d, u_d, x, a)) \\
 &= \text{I} + \text{II} + \text{III}.
 \end{aligned}$$

2043 **Term I.** we use the Lipschitzness of the reward function, and obtain

$$\begin{aligned}
 \text{I} &\leq |f(x, \mathbf{W}\mu(u), a) - f(x, \mathbf{W}\mu(u_d), a)| + \gamma |\langle P(x, \mathbf{W}\mu(u), a), v^{\pi_u}(u, u, \cdot) - v^{\pi_u}(u, u_d, \cdot) \rangle| \\
 &\leq L_f \|\mathbf{W}\mu(u) - \mathbf{W}\mu(u_d)\|_{\text{TV}} \\
 &\quad + \gamma \left\langle P(x, \mathbf{W}\mu(u), a), \mathbb{E} \left[\sum_{\tau \geq 0} \gamma^\tau |f(X_\tau^{\pi_u}, \mathbf{W}\mu(u), \alpha_\tau^{\pi_u}) - f(X_\tau^{\pi_u}, \mathbf{W}\mu(u_d), \alpha_\tau^{\pi_u})| \middle| X_0^{\pi_u} = \cdot, U = u \right] \right\rangle \\
 &\leq L_f \|\mathbf{W}\mu(u) - \mathbf{W}\mu(u_d)\|_{\text{TV}} + \gamma \left\langle P(x, \mathbf{W}\mu(u), a), \mathbb{E} \left[\sum_{\tau \geq 0} \gamma^\tau L_f \|\mathbf{W}\mu(u) - \mathbf{W}\mu(u_d)\|_{\text{TV}} \middle| X_0^{\pi_u} = \cdot, U = u \right] \right\rangle \\
 &\leq \frac{L_f}{1 - \gamma} \|\mathbf{W}\mu(u) - \mathbf{W}\mu(u_d)\|_{\text{TV}} \\
 &\leq \frac{L_f L_d}{(1 - \gamma) D}
 \end{aligned}$$

2059 **Term II.** we first define iteratively for $t \geq 1$ the measure of state-action pair under a policy $\pi_u : \mathcal{X} \rightarrow \mathcal{P}(A)$ as

$$\begin{aligned}
 \underline{P}_\tau^{\pi_u}(x_0, m, a_0) &:= \mathcal{L}(X_\tau^{\pi_u}, \alpha_\tau^{\pi_u} | X_0^{\pi_u} = x_0, a_0^{\pi_u} = a_0) \\
 &= \int_{\mathcal{X}^2 \times A^2} [\delta_{x_\tau} \delta_{a_\tau} \pi_{u,x_\tau}(da_\tau) P(x_{\tau-1}, m, a_{\tau-1})(dx_\tau)] \underline{P}_{\tau-1}^{\pi_u}(x_0, m, a_0)(dx_{\tau-1}, da_{\tau-1}) \\
 &\in \mathcal{P}(\mathcal{X} \times A)
 \end{aligned}$$

2066 we claim that for any $\pi_u : \mathcal{X} \rightarrow \mathcal{P}(A)$, any $(x_0, a_0) \in \mathcal{X} \times A$, any $m_1, m_2 \in \mathcal{P}(\mathcal{X})$ and any time $t \geq 1$,

$$\|\underline{P}_\tau^{\pi_u}(x_0, m_1, a_0) - \underline{P}_\tau^{\pi_u}(x_0, m_2, a_0)\|_{\text{TV}} \leq \tau L_P \|m_1 - m_2\|_{\text{TV}} \tag{26}$$

2069 It is trivial that

$$\underline{P}_1^{\pi_u}(x_0, m, a_0) = P(x_0, m, a_0)$$

2073 is uniformly Lipschitz in measure argument. Assuming equation (26) holds for $t - 1$, we now show it holds for t with the
2074 add and subtract trick again.

$$\begin{aligned}
 &\|\underline{P}_\tau^{\pi_u}(x_0, m_1, a_0) - \underline{P}_\tau^{\pi_u}(x_0, m_2, a_0)\|_{\text{TV}} \\
 &\leq \sup_{\|\phi\|_\infty \leq 1} \int_{A \times \mathcal{X}^2} \phi(a_\tau, x_\tau) \pi_{u,x_\tau}(da_\tau) \left[P_{x_{\tau-1}, m_1, a_{\tau-1}}(dx_\tau) \underline{P}_{\tau-1}^{\pi_u}(x_0, m_1, a_0)(dx_{\tau-1}, da_{\tau-1}) \right. \\
 &\quad \left. - P_{x_{\tau-1}, m_2, a_{\tau-1}}(dx_\tau) \underline{P}_{\tau-1}^{\pi_u}(x_0, m_2, a_0)(dx_{\tau-1}, da_{\tau-1}) \right] \\
 &\leq \sup_{\|\phi\|_\infty \leq 1} \int_{A \times \mathcal{X}^2} \phi(a_\tau, x_\tau) \pi_{u,x_\tau}(da_\tau) \left[P_{x_{\tau-1}, m_1, a_{\tau-1}} - P_{x_{\tau-1}, m_2, a_{\tau-1}} \right](dx_\tau) \underline{P}_{\tau-1}^{\pi_u}(x_0, m_1, a_0)(dx_{\tau-1}, da_{\tau-1}) \\
 &\quad + \sup_{\|\phi\|_\infty \leq 1} \int_{A \times \mathcal{X}^2} \phi(a_\tau, x_\tau) \pi_{u,x_\tau}(da_\tau) P_{x_{\tau-1}, m_2, a_{\tau-1}}(dx_\tau) \left[\underline{P}_{\tau-1}^{\pi_u}(x_0, m_1, a_0) - \underline{P}_{\tau-1}^{\pi_u}(x_0, m_2, a_0) \right](dx_{\tau-1}, da_{\tau-1}) \\
 &\leq (\tau - 1) L_P \|m_1 - m_2\|_{\text{TV}} + L_P \|m_1 - m_2\|_{\text{TV}} \\
 &= \tau L_P \|m_1 - m_2\|_{\text{TV}}
 \end{aligned}$$

2090 With this claim, we have

$$\begin{aligned}
 \text{II} &\leq \left| q^{\pi_u}(u, u_d, x, a) - q^{\pi_u}(u_d, u_d, x, a) \right| \\
 &\leq \sum_{\tau \geq 0} \gamma^\tau \left| \left\langle P_\tau^{\pi_u}(x, \mathbf{W}\mu(u), a) - P_\tau^{\pi_u}(x, \mathbf{W}\mu(u_d), a), f(\cdot, \mathbf{W}\mu(u_d), \cdot) \right\rangle \right| \\
 &\leq \sum_{\tau \geq 0} \gamma^\tau \|f\|_\infty \|P_\tau^{\pi_u}(x, \mathbf{W}\mu(u), a) - P_\tau^{\pi_u}(x, \mathbf{W}\mu(u_d), a)\|_{\text{TV}} \\
 &\leq L_P \|f\|_\infty \|\mathbf{W}\mu(u) - \mathbf{W}\mu(u_d)\|_{\text{TV}} \sum_{\tau \geq 0} \tau \gamma^\tau \\
 &\leq \frac{\|f\|_\infty L_P L_D}{(1 - \gamma)^2 D}
 \end{aligned}$$

2104 **Term III.**

$$2108 \quad \text{III} = q^{\pi_u}(u_d, u_d, x, a) - q^{\pi_{u_d}}(u_d, u_d, x, a) \leq 0$$

2110 is immediate as π_{u_d} is the optimizer of $v^{\pi_{u_d}}(u_d, u_d, \cdot)$.

2112 Finally we conclude

$$2114 \quad \|q_* - \mathbf{\Pi}_D \tilde{q}_*\|_\infty \leq \frac{L_f L_D (1 - \gamma + \|f\|_\infty)}{(1 - \gamma)^2 D}$$

2117 \square

2120 **H. Experiment Setup**

2121 **H.1. Experiment 1: Flocking-Graphon**

2123 The flocking graphon game (Lacker & Soret, 2022) studies the flocking behavior in a system where each player makes
2124 decisions regarding velocity control. The transition dynamic is a continuous time state process, which is given by:

$$2126 \quad dx_t = \alpha_t dt + \sigma dB_t$$

2128 where, $x_t \in \mathcal{X}$. \mathcal{X} is the one-dimensional space. α_t is the velocity control at position x at time t . B_t is a one-dimensional
2129 Brownian motion. The player aims to optimize the running cost. Mathematically,

$$2131 \quad -\mathbb{E} \left[\int_{t=0}^T \alpha_t^2 dt + c |x_T - G^\mu(U)|^2 \right]$$

2135 where, $c > 0$ is a constant, and

$$2137 \quad G^\mu(u) := \langle \mathbf{W}\mu_T(u), \text{Id} \rangle = \int_{[0,1] \times \mathbb{R}} W(u, v) x \mu_T(dv, dx)$$

2140 $G^\mu(u)$ is interpreted as the centroid of the population over the space domain \mathcal{X} , with Id being the identity mapping. More
2141 specifically, $G^\mu(u)$ is the average of the state distribution of the population μ , weighted from the perspective of player with
2142 label u . Intuitively, the running cost arises from change in the velocity, and the terminal cost is associated with deviation
2143 from the centroid at terminal time.

H.2. Experiment 2: SIS-Graphon

(Cui & Koepll, 2022) Consider a game that models pandemic evolution. It admits state space $\mathcal{X} = \{x_S, x_I\}$ where x_S represents a safe state, and x_I represents an infection state. The action space is taken to be $A = \{a_U, a_D\}$, where a_U represents keeping interaction with others and a_D represents taking a quarantine. The terminal time is set to $T = 50$. The transition probability is

$$\begin{aligned}\mathbb{P}(x_S|x_I, m, a) &= \frac{1}{2} & \forall(m, a) \in \mathcal{M}_+(\mathcal{X}) \times A \\ \mathbb{P}(x_I|x_S, m, a_U) &= \frac{4}{5}m(x_I) & \forall m \in \mathcal{M}_+(\mathcal{X}) \\ \mathbb{P}(x_I|x_S, m, a_D) &= 0 & \forall m \in \mathcal{M}_+(\mathcal{X})\end{aligned}$$

An infected agent may turn safe with half probability each time step, regardless of the action. The probability a safe agent is infected is proportion to the infected individuals in her neighborhood when she keep interaction with others, and is 0 when she takes a quarantine. The reward function is given by

$$f(x, m, a) = -2 \cdot \mathbf{1}_{x_I}(x) - 0.5 \cdot \mathbf{1}_{a_D}(a)$$

An agent takes cost from both being infected and taking quarantine action.

H.3. Experiment 3: Investment-Graphon

(Cui & Koepll, 2022) In Investment-Graphon game, The terminal time is set to $T = 50$. Each agent is viewed as a firm, and let $\mathcal{X} = \{0, 1, \dots, 9\}$ be the quality of products this firm provides. With action space given by $A = \{a_I, a_O\}$, the transition kernel is defined by

$$\begin{aligned}\mathbb{P}(x+1|x, m, a_I) &= \frac{9-x}{10} & \forall m \in \mathcal{M}_+(\mathcal{X}) \\ \mathbb{P}(x|x, m, a_I) &= \frac{1+x}{10} & \forall m \in \mathcal{M}_+(\mathcal{X}) \\ \mathbb{P}(x|x, m, a_O) &= 1 & \forall m \in \mathcal{M}_+(\mathcal{X})\end{aligned}$$

We interpret a_I as investment, and a_O as not investing. A firm may improve the product quality by investing, and the probability of a successful investment decrease as the current quality is already high. Initially, every firm starts from quality 0. The reward function is given by

$$f(x, m, a) = \frac{0.3x}{1 + \sum_{x' \in \mathcal{X}} x' m(x')} - 2 \cdot \mathbf{1}_{a_I}(a)$$

A firm's profit is proportion to the quality of product, and decrease with the average product quality within its neighborhood.

I. Experiment Results

In this section, we present detailed numerical results for three graphon games utilized in the main body. The experiment results include algorithm performance (convergece gap, W1-distance, exploitability) and GMFE.

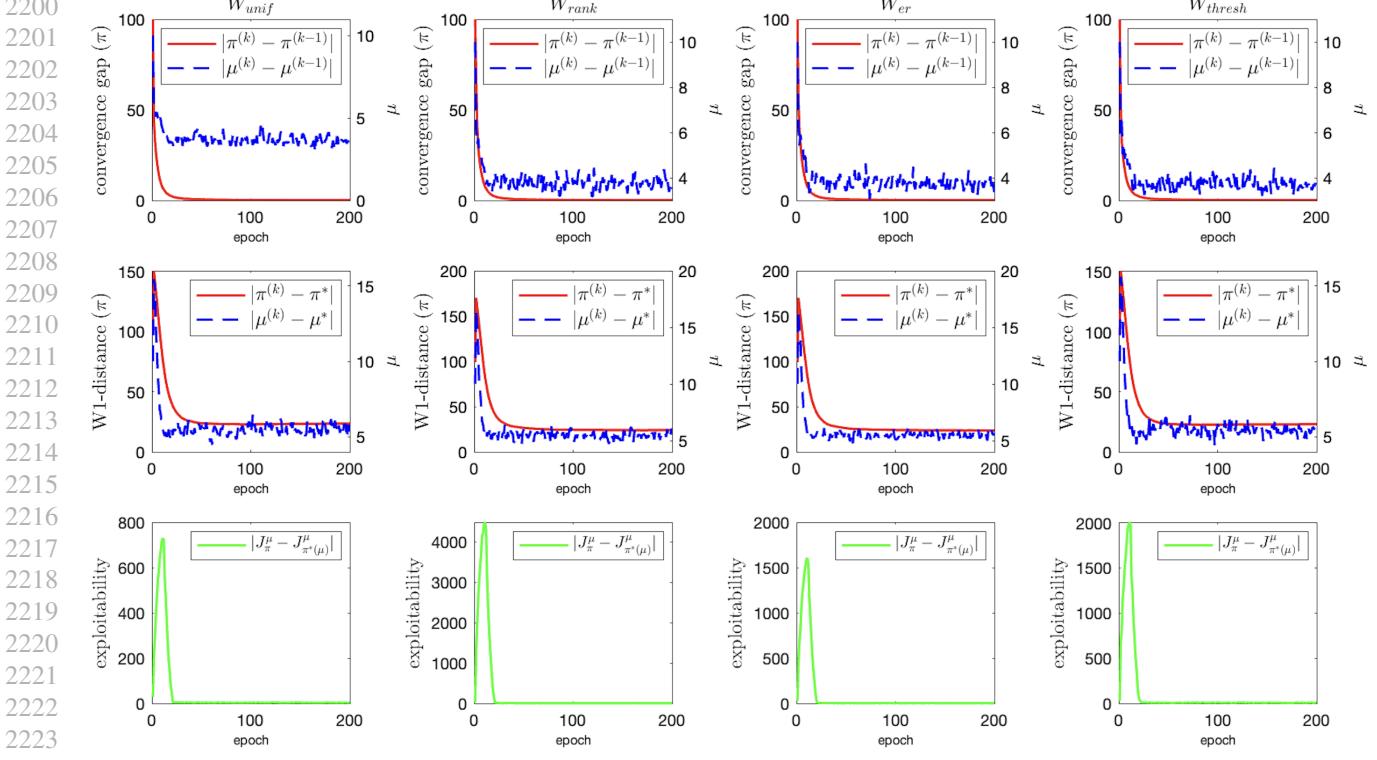


Figure 3. Flocking-Graphon: Algorithm performance. We demonstrate the convergence gap (top), W1-distance (middle) and exploitability (bottom) corresponding to four types of graphs. The exploitability indicates how an agent can improve be deviating from the policy used by the rest of the population. Mathematically, the exploitability is calculated as $|J_\pi^\mu - J_{\pi^*(\mu)}^\mu|$. It measures the gap between the policy adopted by the population and the best policy that an agent can achieve in response to the population state.

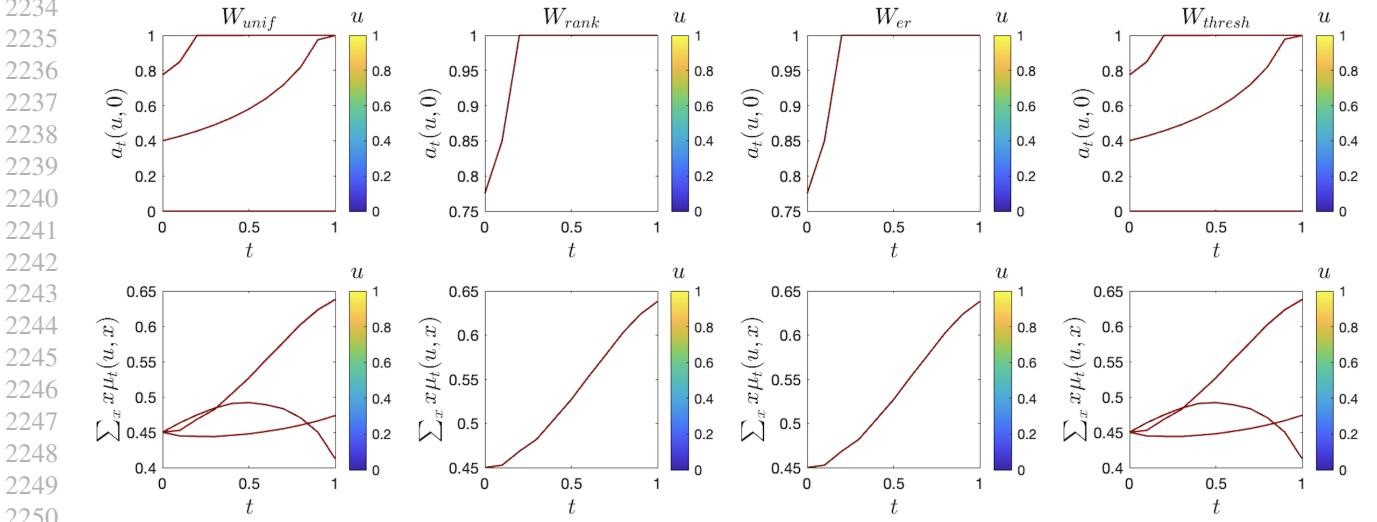


Figure 4. Flocking-Graphon: GMFE. Top: The velocity control at position $x = 0$. The x-axis denotes the time horizon and the y-axis denotes the velocity at equilibrium. The color bar denotes the label state. Bottom: The expected position x across the time. It can be regarded as the centroid of the population.

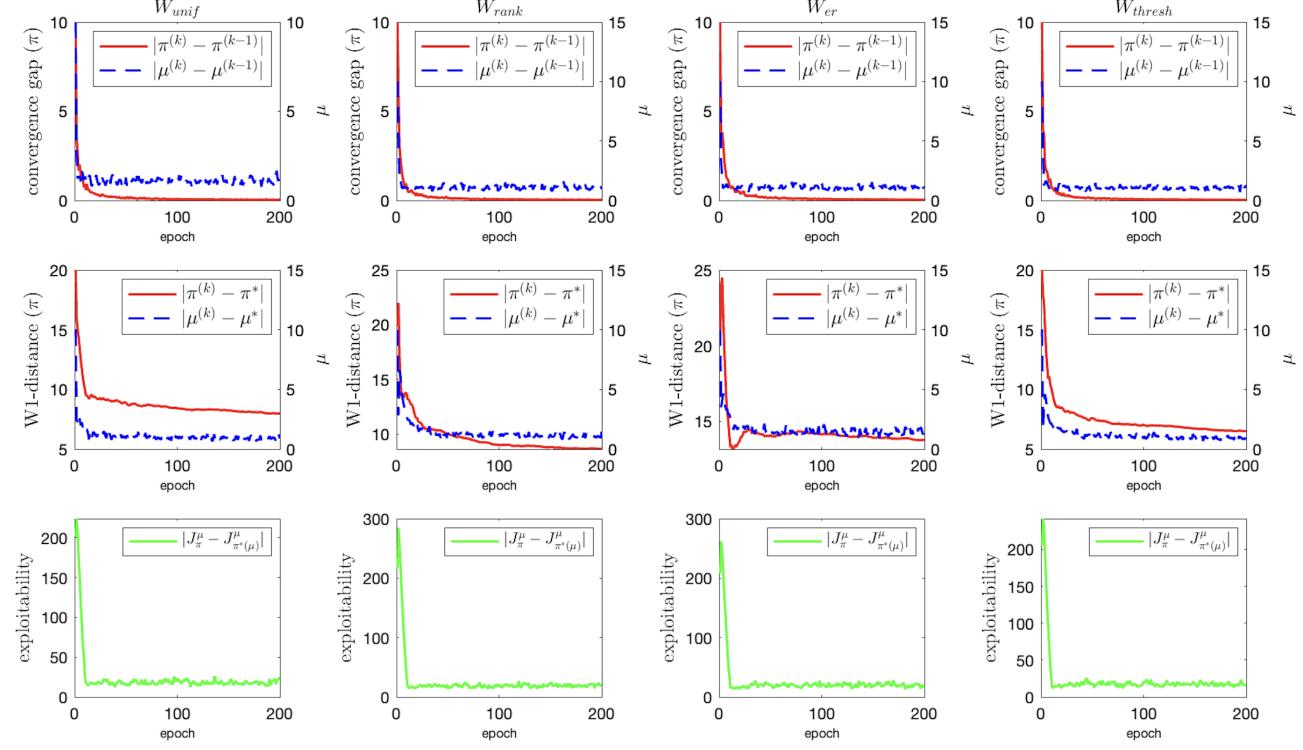


Figure 5. **SIS-Graphon**: Algorithm performance. We demonstrate the convergence gap (top), W1-distance (middle) and exploitability (bottom) corresponding to four types of graphs.

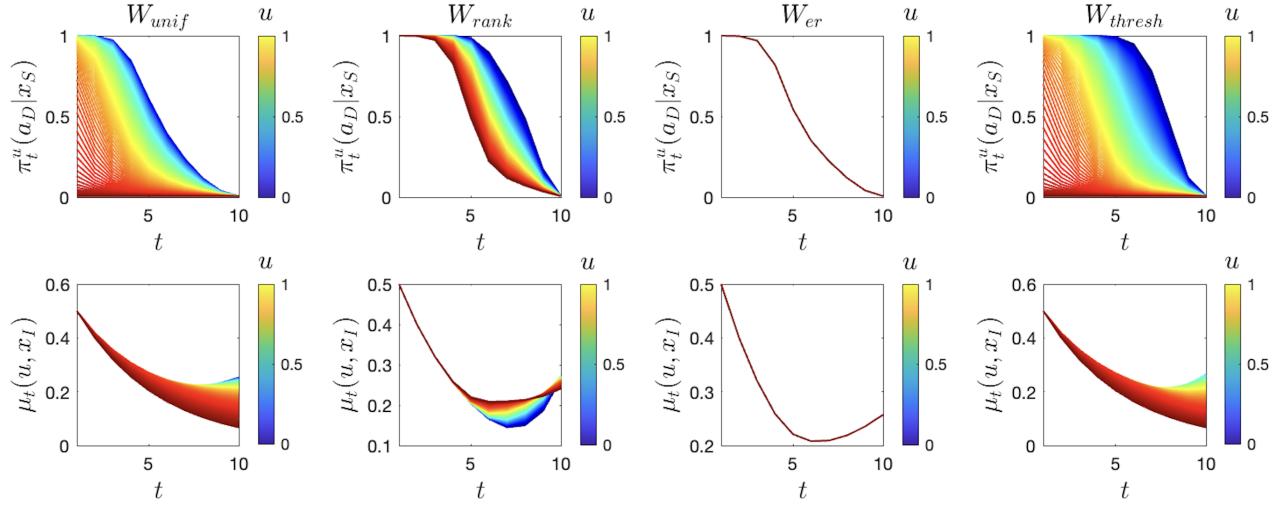


Figure 6. **SIS-Graphon**: GMFE. Top: The probability of taking precautions when healthy. The results for graphs W_{unif} , W_{rank} and W_{er} is consistent with (Cui & Koeppl, 2022). We add the results for graph W_{thresh} . It is shown that the GMFE with W_{thresh} is similar to W_{unif} . Bottom: The population being infected. Agents with a higher u have less connections with others. It means they are less likely infected by the population in a comparison to others. Thus, they take less precautions.

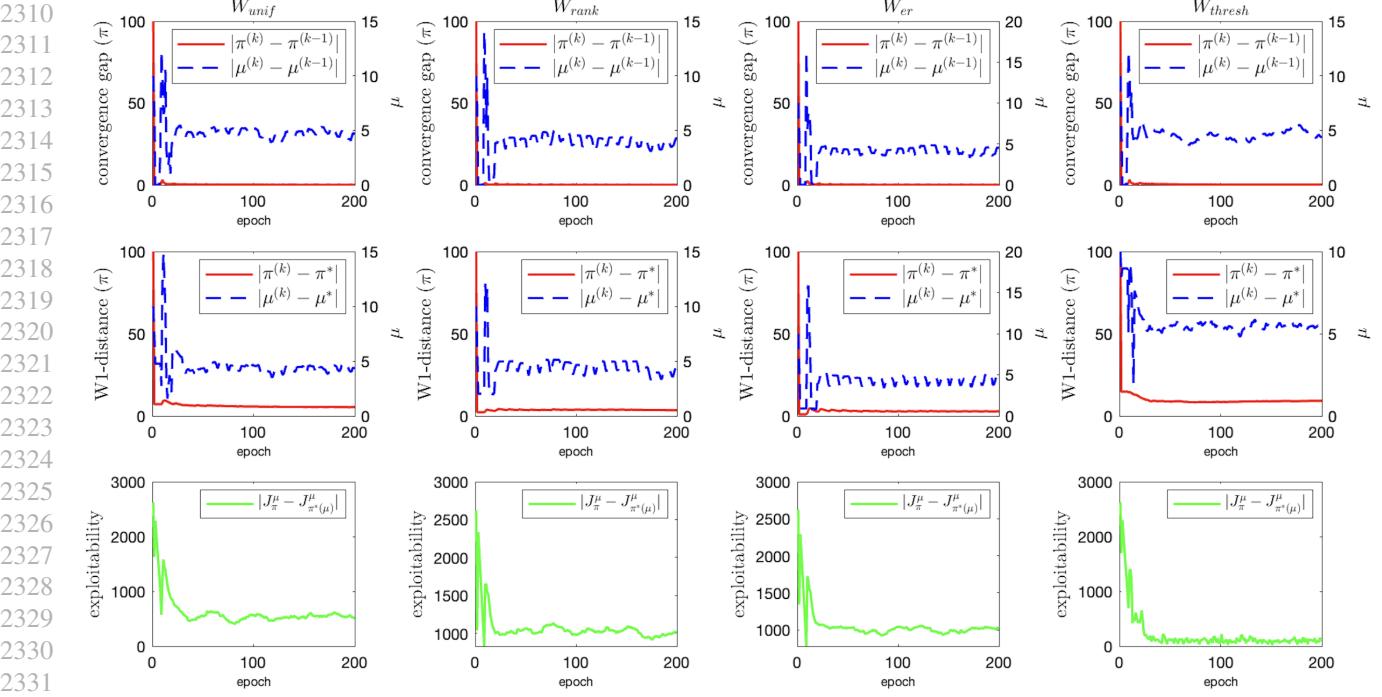


Figure 7. **Invest-Graphon**: Algorithm performance. We demonstrate the convergence gap (top), W_1 -distance (middle) and exploitability (bottom) corresponding to four types of graphs.

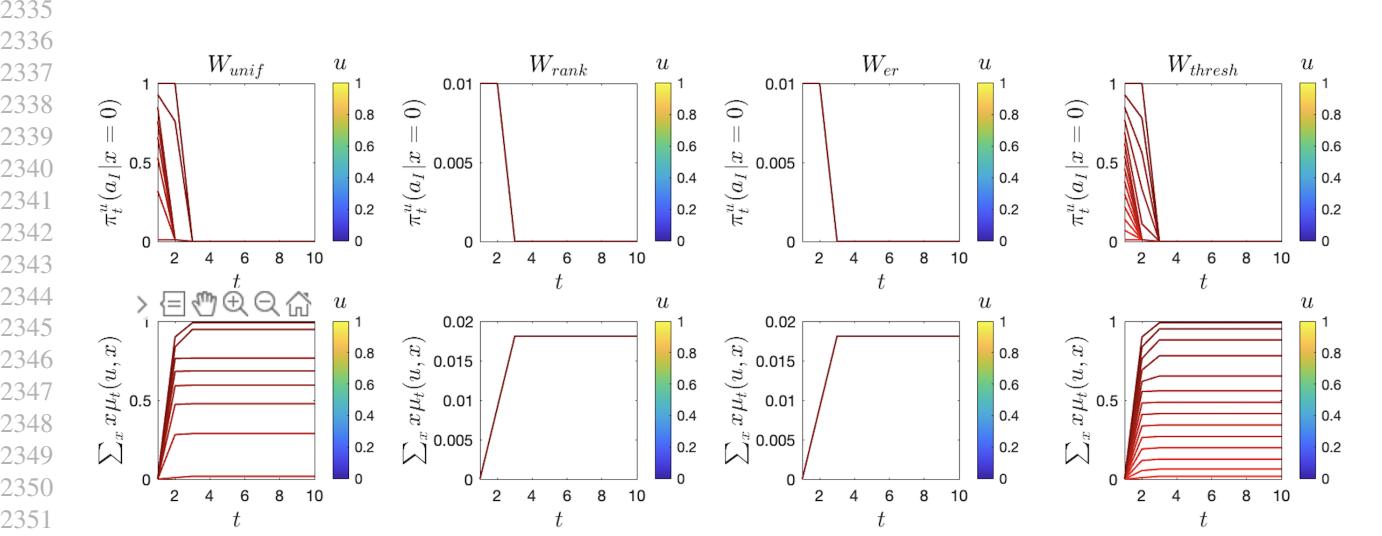


Figure 8. **Invest-Graphon**: GMFE. Top: the probability of investing on product quality when $x = 0$. Bottom: The expected product quality across the time.