# Model-free Reinforcement Learning for Non-stationary Mean Field Games

Rajesh K Mishra, Deepanshu Vasal, and Sriram Vishwanath

*Abstract*— In this paper, we consider a finite horizon, non-stationary, mean field game (MFG) with a large population of homogeneous players, sequentially making strategic decisions, where each player is affected by other players through an aggregate population state termed as *mean field state*. Each player has a private type that only it can observe, and a mean field population state representing the empirical distribution of other players' types, which is shared among all of them. Recently, authors in [1] provided a sequential decomposition algorithm to compute mean field equillibrium (MFE) for such games which allows for the computation of equilibrium policies for them in linear time than exponential, as before. In this paper, we extend it for the case when state transitions are not known, to propose a reinforcement learning algorithm based on Expected Sarsa with a policy gradient approach that learns the MFE policy by learning the dynamics of the game simultaneously. We illustrate our results using cyber-physical security example.

## I. Introduction

There is an increasing number of applications today that involve *large scale* interactions among strategic agents, such as smart grid, autonomous vehicles, cyber-physical systems, Internet of Things (IoT), renewable energy markets, electric vehicle charging, ride sharing apps, financial markets, crypto-currencies and many more. Such applications can be modeled through MFGs introduced by Huang et al [2] and Lasry and Lions [3], where each player is affected by other players not individually, but through a 'mean field' that is an aggregate of the population states.

Multi-agent systems have become ubiquitous owing to their variety of applications. In the recent decade, several multi-agent reinforcement learning (MARL) approaches have been proposed to learn optimal strategies for the agents in the system. However, these approaches scale poorly with size and the dynamic nature of the environment makes the learning of these strategies rather difficult. In contrast, a planning framework based on the mean field approximation has shown potential due to the remarkable property wherein the agents decouple from one another with the increase in their number. As such the agents interact only through the mean field which makes dynamics of the system tractable. MARL systems work well with probabilistic models but fail with large number of agents and this is where the mean field concept has proven useful.

mean field approximation has been in use in game theory applications to study large number of non-cooperative players. In such games, given the mean field states, players' MFE strategies are computed backward recursively through dynamic programming (or HJB equation in continuous time), whereas given players' strategies, mean field states are computed forward recursively using Mckean Vlasov equation (or Fokker-Plank equation). Overall, MFE strategies and mean field states are coupled across time through a fixed-point equation. Recently, in [1], the author presented a sequential decomposition algorithm that computes such equilibrium strategies by decomposing this fixed-point equation across time, and reducing the complexity of this problem from exponential to linear in time. It involves solving a smaller fixed-point equation for each time instant $t$.

The concept of a generalized MFG is discussed by authors in [4], where they prove the uniqueness and existence of Nash equillibrium (NE). They also propose a Q-learning algorithm with Boltzmann policy and analyze its convergence properties and complexity. In [5], the authors propose a posterior sampling based approach where each agent samples a transition probability from a previous transition and converges to the optimal oblivious strategy for MFGs. Authors in [6] consider a multi-agent setting with agents coupled by the average action of the agents. They show that the mean field Q algorithm that converges to the NE. In [7], the authors propose reinforcement learning (RL) algorithms for stationary MFGs that compute MFE and social-welfare optimal solution strategies. In [8], the authors consider a non stationary MFG and propose a fictitious play iterative learning algorithm to devise optimal strategies for mean field states. They argue that if the agent plays the best response to the observed population state flow, they eventually converge to the NE.

Most prior research assume a full knowledge of the dynamics of the system. They also assume a stationary mean field approximation while computing the MFE. In our paper, however, we consider a non stationary MFG with no prior knowledge of the system dynamics, to derive the optimal policy as a function of the mean field state. The main contribution of our paper is an RL algorithm that solves the fixed point equation that was proposed by the authors in [1] while learning the dynamics of the Markov decision process (MDP) simultaneously. We prove the convergence of our algorithm to the MFE analytically. Finally, we show convergence to the mean field equilibrium strategy for a stylized problem of *Malware Spread* where the strategy derived from our algorithm coincides with the optimal strategy obtained assuming the full knowledge of the system.

The paper is structured as follows. In Section II, we present our model. In Section III, we present the sequential decomposition algorithm with backward recursion presented in [1] to compute MFE of the game. In Section IV, we present our reinforcement learning algorithm and prove its convergence to MFE polices. In Section VI, we present

a cyber-physical system security example. We conclude in Section VII.

### A. Notation

We use uppercase letters for random variables and lowercase for their realizations. For any variable, subscripts represent time indices and superscripts represent player identities. We use notation $-i$ to denote all players other than the player $i$ i.e. $-i = \{1, 2, \ldots, i-1, i+1, \ldots, N\}$. We use notation $a_{t:t'}$ to represent the vector $(a_t, a_t + 1, \ldots, a_{t'})$ when $t' \geq t$ or an empty vector if $t' < t$. We use $a_t^{-i}$ to mean $(a_t^1, a_t^2, \ldots, a_t^{i-1}, a_t^{i+1}, \ldots, a_t^N)$. We remove superscripts or subscripts if we want to represent the whole vector, for example $a_t$ represents $(a_t^1, \ldots, a_t^N)$. We denote the indicator function of any set $A$ by $\mathbb{1}_{\{A\}}$. For any finite set $\mathcal{S}$, $\mathcal{P}(\mathcal{S})$ represents space of probability measures on $\mathcal{S}$ and $|\mathcal{S}|$ represents its cardinality. We denote by $P^\sigma$ (or $E^\sigma$) the probability measure generated by (or expectation with respect to) strategy profile $\sigma$ and the space for all such strategies as $\mathcal{K}^\sigma$. We denote the set of real numbers by $\mathbb{R}$. All equalities and inequalities involving random variables are to be interpreted in *a.s.* sense.

### II. MODEL

We consider a finite horizon discrete-time large population sequential game with $N$ homogeneous players, where $N$ tends to $\infty$. We denote the set of homogeneous players by $\mathcal{N}$ and the set of time periods by $\mathcal{T}$. In each time instant $t \in \mathcal{T}$, a player $i \in \mathcal{N}$ observes a private type $x_t^i \in \mathcal{X} = \{1, 2, \cdots, N_x\}$ and a common observation of the mean field population state $z_t \in \mathcal{Z}$, where $z_t = (z_t(1), z_t(2), \ldots, z_t(N_x))$ is the fraction of population having type $x \in \mathcal{X}$ at time $t$ given as.

$$z_t(x) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}_{\{x_t^i = x\}} \tag{1}$$

with $\sum_{i=1}^{N_x} z_t(i) = 1$. Consequently, the player $i$ takes an action $a_t^i \in \mathcal{A} = \{1, 2, \cdots, N_a\}$ based on policy $\sigma_t^i(.|z_{1:t}, x_t^i)$, and receives a reward $R(x_t^i, a_t^i, z_t)$ which is a function of its current type $x_t^i$, action $a_t^i$ and the common observation $z_t$. Player $i$'s type evolves as a controlled Markov process,

$$x_{t+1}^i = f_x(x_t^i, a_t^i, z_t, w_t^i). \tag{2}$$

The random variables $(w_t^i)_{i \in \mathcal{N}, t \in \mathcal{T}}$ are assumed to be mutually independent across players and across time. We can also express the above update of $x_t^i$ through a kernel, $x_{t+1}^i \sim \tau(\cdot|x_t^i, a_t^i, z_t)$, where $\tau(\cdot|\cdot)$ represents the transition probabilities of the MDP. In this paper, we assume that $\tau$ is unknown. We develop a model-free RL algorithm to derive the equilibrium policy where $\tau$ is used to simulate the model and generate samples for learning. The idea of MFGs is to approximate the finite population game with an infinite population game such the the mean field state $z_t$ converges to the statistical MFE of the game [7].

In any time period $t$, player $i$ observes $(z_{1:t}, x_{1:t}^i)$ and takes action $a_t^i$ according to a behavioral strategy $\sigma^i = (\sigma_t^i)_t$,

where $\sigma_t^i : \mathcal{Z}^t \times \mathcal{X}^t \to \mathcal{P}(\mathcal{A})$ defined over the space $\mathcal{K}^\sigma$ which implies $A_t^i \sim \sigma_t^i(\cdot|z_{1:t}, x_{1:t}^i)$. In the case of a finite time-horizon game, $\mathbb{G}_T$, each player wants to choose the strategy $\tilde{\sigma}$ that maximizes its total expected discounted reward over a time-horizon $T$, discounted by discount factor $0 < \delta \leq 1$, given as

$$J_t^{\tilde{\sigma}} = \mathbb{E}^{\tilde{\sigma}_{t:T}} \left[ \sum_{k=t}^{T} \delta^{k-t} R(X_k^i, A_k^i, z_k) | z_t, x_t^i \right]. \tag{3}$$

### III. PRELIMINARIES

### A. Solution concept: Markov perfect equillibrium (MPE)

The Nash equilibrium of $\mathbb{G}_T$ is defined as strategies $\tilde{\sigma} = (\tilde{\sigma}_t^i)_{i \in \mathcal{N}, t \in \mathcal{T}}$ that satisfy, for all $i \in \mathcal{N}$,

$$\mathbb{E}^{\tilde{\sigma}^i, \tilde{\sigma}^{-i}} \left[ \sum_{t=1}^{T} \delta^{t-1} R(X_t^i, A_t^i, Z_t) \right]$$

$$\geq \mathbb{E}^{\sigma^i, \tilde{\sigma}^{-i}} \left[ \sum_{t=1}^{T} \delta^{t-1} R(X_t^i, A_t^i, Z_t) \right]. \tag{4}$$

For sequential games, however, a more appropriate equilibrium concept is MPE [9], which is used in this paper. We note that although a Markov perfect equilibrium is also a Nash equilibrium of the game, the opposite might not be true always. An MPE $(\tilde{\sigma})$ satisfies sequential rationality such that for $\mathbb{G}_T$, $\forall i \in \mathcal{N}$, $\forall t \in \mathcal{T}$, $\forall h_t^i \in \mathcal{H}_t^i$, where $\mathcal{H}_t^i$ is the space of all possible mean field trajectories till time $t$, and $\forall \sigma^i \in \mathcal{K}^\sigma$,

$$\mathbb{E}^{\tilde{\sigma}^i \tilde{\sigma}^{-i}} \left[ \sum_{n=t}^{T} \delta^{n-t} R(X_n^i, A_n^i, Z_n) | z_{1:t}, x_{1:t}^i \right]$$

$$\geq \mathbb{E}^{\sigma^i \tilde{\sigma}^{-i}} \left[ \sum_{n=t}^{T} \delta^{n-t} R(X_n^i, A_n^i, Z_n) | z_{1:t}, x_{1:t}^i \right]. \tag{5}$$

### B. A solution concept: MFE

MFE can be defined as the combination of the optimal policy $\tilde{\sigma} := \{\tilde{\sigma}_t\}_{t \in \mathcal{T}}$ and the mean field states $z := \{z_t\}_{t \in \mathcal{T}}$ that satisfy the following:

1) A policy $\tilde{\sigma}$ for some $z_{1:T}$ such that

$$\mathbb{E}^{\tilde{\sigma}^i} \left[ \sum_{k=t}^{T} \delta^{k-t} R(X_k^i, A_k^i, Z_k) | z_{1:T}, x_{1:t}^i \right]$$

$$\geq \mathbb{E}^{\sigma^i} \left[ \sum_{k=t}^{T} \delta^{k-t} R(X_k^i, A_k^i, Z_k) | z_{1:T}, x_{1:t}^i \right]. \tag{6}$$

2) A function $\Phi[z] : \mathcal{S}_z \to 2^{\mathcal{S}_\sigma}$ such that

$$\Phi(z) = \{\tilde{\sigma} \in S_\sigma : \tilde{\sigma} \text{ is optimal for } z\} \tag{7}$$

3) A mapping $\Lambda : \mathcal{S}_\sigma \to \mathcal{S}_z$ with $\tilde{\sigma} \in S_\sigma$, $z = \Lambda(\tilde{\sigma})$ is constructed recursively as,

$$z_{t+1}(y) = \sum_{x_t, a_t} \tau(y|x, a, z_t) \tilde{\sigma}_t(a|z_t, x_t) z_t(dx_t) \tag{8}$$

Then, the pair $(\tilde{\sigma}, z)$ can be called a MFE which is a good approximation for the MPE when the population grows large.

## C. A methodology to compute MFE

In this section, we summarize the backward recursive methodology based on sequential decomposition that would be used to compute the MPE. It is worth noting that, [1] provides a sequential decomposition algorithm used for the case when model of the MDP us known. Here, we modify the algorithm so that it could be used for the model-free case as well. We consider Markovian equilibrium strategies of player $i$ which depend on the common information $z_t$ at time $t$, and on its current private type $x_t^i$. Equivalently, player $i$ takes action of the form $A_t^i \sim \sigma_t^i(\cdot|z_t, x_t^i)$. Similar to the common agent approach [10], an alternate and equivalent way of defining the strategies of the players is as follows. We consider a common fictitious agent that views the common information $z_t$ and generates a prescription function $\gamma_t^i : \mathcal{X} \to \mathcal{P}(\mathcal{A})$ as a function of $z_t$ through an equilibrium generating function $\theta_t^i : \mathcal{Z} \to (\mathcal{X} \to \mathcal{P}(\mathcal{A}))$ such that $\gamma_t^i = \theta_t^i[z_t]$. Then action $A_t^i$ is generated by applying this prescription function $\gamma_t^i$ on player $i$'s current private information $x_t^i$, i.e. $A_t^i \sim \gamma_t^i(\cdot|x_t^i)$. Thus $A_t^i \sim \sigma_t^i(\cdot|z_t, x_t^i) = \theta_t^i[z_t](\cdot|x_t^i)$.

We are only interested in symmetric equilibria of such games such that $A_t^i \sim \gamma_t(\cdot|x_t^i) = \theta_t[z_t](\cdot|x_t^i)$ i.e. there is no dependence of $i$ on the strategies of the players.

For a given symmetric prescription function $\gamma_t = \theta[z_t]$, the statistical mean field $z_t$ evolves according to the discrete-time McKean Vlasov equation, $\forall y \in \mathcal{X}$:

$$z_{t+1}(y) = \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} z_t(x)\gamma_t(a|x)\tau(y|x, a, z_t), \quad (9)$$

which implies

$$z_{t+1} = \phi(z_t, \gamma_t). \quad (10)$$

## D. Backward recursive algorithm for $\mathbb{G}_T$

This section summarizes the proposed a novel model-free algorithm to compute the optimum policy function $\tilde{\theta}_t$ as a function of mean field $z_t$ where equilibrium generating function $(\tilde{\theta}_t)_{t \in [T]}$ is defined as $\tilde{\theta}_t : \mathcal{Z} \to \{\mathcal{X} \to \mathcal{P}(\mathcal{A})\}$, and for each $z_t$, we generate $\tilde{\gamma}_t = \tilde{\theta}_t[z_t]$. In addition, we generate an action value function $Q_t$ defined as $Q_t : \Pi \times \mathcal{Z} \times \mathcal{X} \times \mathcal{A} \to \mathbb{R}$ that captures the expected sum of returns at time $t$ following certain action from a state and then continuing with the optimal policy $\tilde{\sigma}_{t+1}$ from time $t+1$ onwards. As per our knowledge, this is the first attempt to solve a fixed point equation using the action value function in a model-free algorithm and determining the corresponding equilibrium policy when the mean field is non-stationary. The algorithm can be summarized as follows:

1) Initialize $\forall z_{T+1}, x_{T+1}^i \in \mathcal{X}, A_{T+1}^i \in \mathcal{A}, \tilde{\gamma}_{T+1} \in \Pi$,

$$V_{T+1}[z_{T+1}, x_{T+1}^i] \triangleq 0, \quad (11)$$

$$\tilde{\theta}_{T+1}[z_{T+1}] \triangleq 0. \quad (12)$$

2) For $t = T, T-1 \ldots 1$ and $\forall z_t$, $\tilde{\theta}_t[z_t]$ is generated through the following steps.

a) Compute $Q_t$, $\forall x_t^i \in \mathcal{X}$, $\forall a_t^i \in \mathcal{A}$, and $\forall \tilde{\gamma}_t \in \Pi$ as,

$$Q_t(z_t, x_t^i, a_t, \tilde{\gamma}_t) \triangleq R(z_t, x_t^i, a_t^i) + \delta \mathbb{E}[V_{t+1}(z_{t+1}, X_{t+1}^i)|z_t, x_t^i], \quad (13)$$

where the expectation in (13) is with respect to the random variable $X_{t+1}^i$ through the measure $\tau(x_{t+1}^i|x_t^i, a_t^i, z_t)$. The mean state $z_{t+1} = \phi(z_t, \tilde{\gamma}_t)$.

b) Set $\tilde{\theta}_t[z_t] = \tilde{\gamma}_t$, where $\tilde{\gamma}_t$ is the solution to the following fixed-point equation at all $x_t^i \in \mathcal{X}$ and $\forall i \in \mathcal{N}$

$$\tilde{\gamma}_t(\cdot|x_t^i) \in$$
$$\arg\max_{\gamma_t(\cdot|x_t^i)} \mathbb{E}^{\gamma_t(\cdot|x_t^i)}[Q_t(z_t, x_t^i, A_t^i, \tilde{\gamma}_t)|z_t, x_t^i]$$
$$(14)$$

where expectation in (14) is with respect to random variable $A_t^i$ through the measure $\gamma_t(a_t^i|x_t^i)$.

c) The value function $V_t$ is computed $\forall x_t^i \in \mathcal{X}$ as,

$$V_t(z_t, x_t^i) = \mathbb{E}^{\tilde{\gamma}_t(\cdot|x_t^i)}[Q_t(z_t, x_t^i, A_t^i, \tilde{\gamma}_t)|z_t, x_t^i]$$
$$(15)$$

Then, an equilibrium strategy is defined as

$$\tilde{\sigma}_t^i(a_t^i|z_{1:t}, x_{1:t}^i) = \tilde{\theta}[z_t](a_t^i|x_t^i), \quad (16)$$

where $\tilde{\theta}[z_t] = \tilde{\gamma}_t$. The proof for the existence of the solution to the fixed point equation in (14) which is the MPE of the game has already been provided in [1] and will be revisited when we prove the convergence of our algorithm.

## IV. REINFORCEMENT ALGORITHM

In this section, we describe our proposed RL algorithm that computes the optimal policy which maximizes the expected sum of returns as specified in (3). The optimal policy $\tilde{\sigma}$ is defined for a discretized set of mean states $z \in \mathcal{Z}$, given as $\tilde{\theta}[z]$. At any time $t$, and for any current mean state $z_t \in \mathcal{Z}$, $\tilde{\theta}_t[z_t]$ maps to a function $\tilde{\gamma}_t$ that prescribes the probabilistic action $a_t^i$, an agent $i$ should take, given the state $x_t^i$.

We implement the RL algorithm based on Expected Sarsa [11] and without the explicit knowledge of the the transition probabilities $\tau(x_{t+1}^i|x_t^i, a_t^i)$. The algorithm basically computes the $Q_t$-values at each instant and then learns the optimal policy $\tilde{\gamma}_t$ by solving the fixed point equation in (14). We use Expected Sarsa to update the $Q_t$ values, from the current reward $r_t$ and $V_{t+1}$, the value at the future state. This update can be expressed as

$$Q_t(z_t, x_t, a_t, \tilde{\gamma}_t) = (1-\alpha)Q_t(z_t, x_t^i, a_t^i) + \alpha(r_t + \delta V_{t+1}(z_{t+1}, x_{t+1}^i)) \quad (17)$$

The $Q$-values are a function not only of the states and the actions but also of the current mean field state $z_t$ and the current optimal policy $\tilde{\theta}[z_t] = \tilde{\gamma}_t$. The current optimal policy and current mean state determine the next mean state $z_{t+1}$ which determines the value function $V_{t+1}(z_{t+1}, x_{t+1}^i)$ at the future state. Therefore, the functions $Q_t$ is defined

over all possible equilibrium policies $\tilde{\gamma}_t\left(\cdot|x_t\right) \in \Pi$, where $\Pi$ is the space of all possible strategies from a given state $x_t$. The value function $V_{t+1}\left(z_{t+1}, x_{t+1}^i\right)$ is determined using functional approximation. In addition, due to the non-stationarity of the mean states, the equation in (14), which is used to solve for the optimal policy, is not just a single step optimization, but a fixed point equation which needs to be solved with repeated iterations. In our paper, we use a policy gradient approach to solve for the optimal policy at each mean state and for every time iteration repeatedly. The entire RL algorithm described here is summarized in Algorithm 1.

At each time instant $t$, the policy iteration algorithm computes the equilibrium policy based on the action value function $Q_t$ through a policy gradient approach. In other words, the solution to the fixed point equation in (14) is the policy where $Q_t$ has the highest gradient. Given that the the function $Q_t$ is a function of the optimal policy itself, the new found policy changes the $Q_t$. Therefore, this process is repeated over several iterations in order to arrive at the required prescription function $\tilde{\gamma}_t$. This is repeated at all the mean states $z_t \in \mathcal{Z}$ so that we get the final equilibrium function $\tilde{\theta}\left[z\right]$.

## V. CONVERGENCE

In this section, we prove the convergence of the proposed RL algorithm to the equilibrium strategy of the statistical MFG. Using backward recursion and sequential decomposition, the RL algorithm is able to arrive at the equilibrium strategy for player $i$ at each time $t$. In other words, we show that the a player $i$ has no incentive to deviate from the equilibrium strategy $\tilde{\sigma}$ given that the other players are playing the equilibrium strategy. Before proving convergence, we establish two lemmas related to the main theorems. In the first lemma we show that the value function $V_t$ captures the expected sum of rewards accumulated by playing the $\tilde{\sigma}$ at time $t$ by the $i^{th}$ player. We follow it with the proof of the next lemma establishing the optimality of the $V$ value over the expected sum of rewards accumulated by playing any other strategy other than $\tilde{\sigma}$.

*Lemma 1:* $\forall t \in [T], \forall z_t, x_t^i \in \mathcal{X}$,

$$V_t\left(z_t, x_t^i\right) = \mathbb{E}^{\tilde{\sigma}_{t:T}}\left[\sum_{k=t}^{T} \delta^{k-t} R\left(x_k^i, A_k^i, z_k\right)|z_t, x_t^i\right] \quad (18)$$

$$= J_t^{\tilde{\sigma}} \quad (19)$$

where $\tilde{\sigma}_t$ is the equilibrium policy at time $t$ and $J_t^{\tilde{\sigma}}$ is the accumulated optimal returns from $t$ till $T$ by following the equilibrium policy.

*Proof:* We prove the lemma using the theory of mathematical induction.

At $t = T$, from (15),

$$V_T\left(z_T, x_T^i\right) = \mathbb{E}^{\tilde{\sigma}_T}\left[Q_T\left(z_T, x_T^i, A_T^i, \tilde{\sigma}_T\left(\cdot, z_T\right)\right)\right] \quad (20)$$

$$= \mathbb{E}^{\tilde{\sigma}_T}\left[R\left(x_T^i, A_T^i, z_T\right)\right], \quad (21)$$

---

**Algorithm 1:** Equilibrium Policy

**Input:**
L: Batch Size for Sarsa
I: Policy Iterations
T: Time Length

**Output:**
$\tilde{\sigma}$: Optimal Policy

1 Initialize: $V_{T+1}$
2 Initialize: $\tilde{\theta}_{T+1}$
3 **for** $t = T \ldots 1$ **do**
4     **for** $z_t \in \mathcal{Z}, \tilde{\gamma}_t \in \Pi$ **do**
5         Next mean state: $z_{t+1} \sim \phi\left(z_t, \tilde{\gamma}_t\right)$
6         Optimum policy: $\tilde{\gamma}_{t+1} = \tilde{\theta}_{t+1}\left[z_{t+1}\right]$
7         **for** $l = 1, 2, \ldots L$ **do**
8             **for** $\left(x_l^i, a_l^i\right) \in \mathcal{X} \times \mathcal{A}$ **do**
9                 Sample: $x_{l+1}^i \sim \tau\left(.|x_l^i, a_l^i, z_t\right)$
10                 Sarsa Target:
                  $G = R(x_l^i, a_l^i, z_t) + \delta V_{t+1}\left(z_{t+1}, x_{l+1}^i\right)$
11                 $Q_t\left(z_t, x_l^i, a_l^i, \tilde{\gamma}_t\right) =$
                $\left(1 - \alpha\right) Q_t\left(z_t, x_l^i, a_l^i, \tilde{\gamma}_t\right) + \alpha G$
12             **end**
13         **end**
14     **end**
15     **for** $z \in \mathcal{Z}$ **do**
16         **for** $n = 1 \ldots I$ **do**
17             $\tilde{\gamma}_n = \text{PG}\left(Q_t\left(z_t, \cdot, \cdot, \tilde{\gamma}_{n-1}\right)\right)$
18             Increment: $n = n + 1$
19         **end**
20         $\tilde{\theta}\left[z_t\right] = \tilde{\gamma}_n$
21     **end**
22     **for** $z \in \mathcal{Z}$ **do**
23         **for** $x^i \in \mathcal{X}$ **do**
24             $\tilde{\gamma} = \tilde{\theta}\left[z\right]$
25             $V_t\left(z_t, x^i\right) = \mathbb{E}^{\tilde{\gamma}}\left[Q_t\left(z_t, x^i, A^i, \tilde{\gamma}\right)\right]$
26         **end**
27     **end**
28 **end**
29 $\tilde{\sigma} = \tilde{\theta}[z_t] \,\forall z_t$

**Result:** $\tilde{\sigma}$

---

which is true from (13). For any mean state $z_t$, $\tilde{\gamma}_t = \tilde{\sigma}_t\left(\cdot|\cdot, z_t\right)$, therefore, we have,

$$V_T\left(z_T, x_T^i\right) = \mathbb{E}^{\tilde{\gamma}_T}\left[R\left(x_T^i, A_T^i\right)\right]. \quad (22)$$

which is the maximum returns the agents can receive at $t = T$ because $\tilde{\gamma}_t$ is the solution to the fixed point equation in (14) that maximizes (22).

Now assuming that the proposition is true for $t = t + 1$, we get,

$$V_{t+1}\left(z_{t+1}, x_{t+1}^i\right) = J_{t+1}^{\tilde{\sigma}} \quad (23)$$

$$= \mathbb{E}^{\tilde{\sigma}_{t+1:T}}\left[\sum_{k=t+1}^{T} \delta^{k-(t+1)} R\left(x_k^i, A_k^i, z_k\right)|z_{t+1}, x_{t+1}^i\right]$$
$$\quad (24)$$

$$\text{(25)}$$

At time $t = t$, we have,

$$V_t\left(z_t, x_t^i\right) = \mathbb{E}^{\tilde{\sigma}_t}\left[Q_t\left(z_t, x_t^i, A_t^i, \tilde{\sigma}_t\left(\cdot|\cdot, z_t\right)\right)|z_t, x_t^i\right] \quad \text{(26a)}$$

$$= \mathbb{E}^{\tilde{\sigma}_t}[R\left(x_t^i, z_t, A_t^i\right) +$$
$$\delta V_{t+1}\left(z_{t+1}, x_{t+1}^i\right)|z_t, x_t^i] \quad \text{(26b)}$$

(26a) is from the definition in (15) while (26b) is from the definition in (13). Using the assumption in (24), we get,

$$V_t\left(z_t, x_t^i\right) = \mathbb{E}^{\tilde{\sigma}_t}\Bigg[R\left(x_t^i, A_t^i, z_t\right) +$$

$$\delta\mathbb{E}^{\tilde{\sigma}_{t+1:T}}\left[\sum_{k=t+1}^{T}\delta^{k-(t+1)}R\left(x_k^i, z_k\right)|z_{t+1}, x_{t+1}^i\right]$$

$$|z_t, x_t^i\Bigg] \quad \text{(27)}$$

Using the the expression for expected sum of returns in (3), we get,

$$V_t\left(z_t, x_t^i\right) = \mathbb{E}^{\tilde{\sigma}_t}[R\left(x_t^i, z_t, A_t^i\right) + \delta J_{t+1}] \quad \text{(28)}$$

■

*Lemma 2:* $\forall i \in \mathcal{N}, t \in [T], \forall z_t, x_t^i \in \mathcal{X}, a_t^i \in \mathcal{A}$,

$$V_t\left(z_t, x_t^i\right) \geq \mathbb{E}^{\sigma_t^i, \tilde{\sigma}_t^{-i}}\left[Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right)|z_t, x_t^i\right] \quad \text{(29)}$$

where $z_{t+1} = \phi\left(z_t, \tilde{\sigma}_t\left(\cdot|z_t, \cdot\right)\right)$.

*Proof:* Given that at time $t+1$ the equilibrium policy is $\tilde{\sigma}_{t+1}$ is the solution to the fixed point equation in (14), let us assume that the player $i$ plays a different policy $\widehat{\sigma}_{t+1}$ such that $\forall x_{t+1}^i \in \mathcal{X}$,

$$\widehat{\gamma}_t\left(\cdot|x_{t+1}^i\right) \notin \arg\max_{\gamma_t}\mathbb{E}\left[Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right)|z_t, x_t^i\right] \quad \text{(30)}$$

with $\widehat{\sigma}_t\left(\cdot|z_t, x_t^i\right) = \widehat{\gamma}_t\left(\cdot|x_t^i\right)$.

Now, the equation on the other side of the inequality in(29) can be changed assuming $\widehat{\sigma}$ as the suboptimal policy as follows:

$$\mathbb{E}^{\sigma_t^i, \tilde{\sigma}_t^{-i}}\left[Q_t\left(z_{t+1}, X_t^i, A_t^i, \tilde{\gamma}_t\right)|z_t, x_t^i, a_t^i\right]$$
$$= \mathbb{E}^{\widehat{\gamma}_t^i, \tilde{\sigma}_t^{-i}}\left[Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right)|z_t, x_t^i, a_t^i\right] \quad \text{(31a)}$$

Considering (30), we can proceed as,

$$\mathbb{E}^{\widehat{\gamma}_t^i, \tilde{\sigma}_t^{-i}}\left[Q_t\left(z_t, X_t^i, A_t^i, \tilde{\gamma}_t\right)|z_t, x_t^i, a_t^i\right]$$
$$\leq \mathbb{E}^{\tilde{\gamma}_t^i, \tilde{\sigma}_t^{-i}}\left[Q_t\left(z_{t+1}, X_t^i, A_{t+1}^i, \tilde{\gamma}_t\right)|z_t, x_t^i, a_t^i\right] \quad \text{(32a)}$$
$$= V_t\left(z_t, x_t^i\right) \quad \text{(32b)}$$

where the last statement is from the definition in (13).

■

*Theorem 1:* A strategy $(\tilde{\sigma})$ constructed from the above algorithm is an MFE of the game. i.e

$$\mathbb{E}^{\tilde{\sigma}}\left[\sum_{k=t}^{T}\delta^{k-t}R\left(X_k^i, A_k^i, Z_k\right)|z_{1:t}, x_{1:t}^i\right]$$

$$\geq \mathbb{E}^{\sigma^i, \tilde{\sigma}^{-i}}\left[\sum_{k=t}^{T}\delta^{k-t}R\left(X_k^i, A_k^i, Z_k\right)|z_{1:t}, x_{1:t}^i\right] \quad \text{(33)}$$

*Proof:* We prove it through the technique of mathematical induction and will use the results that were proved before in Lemma 1 and Lemma 2.

For the base case, we consider $t = T$. The expected sum of returns, when the player $i$ follows the equilibrium policy $\tilde{\sigma}$ is given as

$$J_T = \mathbb{E}^{\tilde{\sigma}_T}\left[R\left(x_T^i, A_T^i, z_T\right)|z_{1:T}, x_{1:T}^i\right] \quad \text{(34)}$$
$$= V_T\left(z_T, x_T^i\right), \quad \text{(35)}$$

which is true from Lemma 1.

Now, from Lemma 2, we have,

$$V_T \geq \mathbb{E}^{\sigma_T^i, \tilde{\sigma}_T^{-i}}\left[Q_T\left(z_T, x_T^i, A_T^i, \tilde{\gamma}_T\right)|z_t, x_t^i\right] \quad \text{(36)}$$
$$= \mathbb{E}^{\sigma_T^i, \tilde{\sigma}_T^{-i}}\left[R\left(x_T^i, A_T^i, z_T\right)|z_{1:T}, x_{1:T}^i\right] \quad \text{(37)}$$

Assuming that the condition in (33) holds at $t = t + 1$, we get,

$$\mathbb{E}^{\tilde{\sigma}}\left[\sum_{k=t+1}^{T}\delta^{k-t-1}R\left(x_k^i, A_k^i, z_k\right)|z_{1:t+1}, x_{1:t+1}^i\right]$$

$$\geq \mathbb{E}^{\sigma^i\tilde{\sigma}^{-i}}\left[\sum_{k=t+1}^{T}\delta^{k-t-1}R\left(x_k^i, A_k^i, z_k\right)|z_{1:t+1}, x_{1:t+1}^i\right]$$
$$\text{(38)}$$

We need to prove that the expression in (33) holds for $t = t$ as well. Let us represent the left hand side of (33) as $L$, i.e.

$$L = \mathbb{E}^{\tilde{\sigma}}\left[\sum_{n=t}^{T}\delta^{k-t}R\left(X_k^i, A_k^i, Z_k\right)|z_{1:t}, x_{1:t}^i\right] \quad \text{(39)}$$

The expectation at $t + 1$ is independent of the rewards at time $t$, therefore, we can rewrite (39) as

$$L = \mathbb{E}^{\tilde{\sigma}_t}\Bigg[R\left(x_t^i, a_t^i, z_t\right) +$$

$$\delta\mathbb{E}^{\tilde{\sigma}_{t+1:T}}\left[\sum_{k=t+1}^{T}\delta^{k-t-1}R\left(X_k^i, A_k^i, Z_k\right)|z_{1:t+1}, x_{1:t+1}^i\right]\Bigg]$$
$$\text{(40)}$$

Using Lemma 1, and the definition of $V_t$ in (15) we get,

$$L = \mathbb{E}^{\tilde{\sigma}_t}\left[R\left(x_t^i, a_t^i, z_t\right) + \delta V_{t+1}\left(z_{1:t+1}, x_{1:t+1}^i\right)\right] \quad \text{(41a)}$$
$$= \mathbb{E}^{\tilde{\sigma}_t}\left[Q_t\left(z_t, x_t^i, a_t^i, \tilde{\sigma}_t\left(\cdot|z_t, \cdot\right)\right)\right] \quad \text{(41b)}$$
$$= V_t\left(z_t, x_t^i\right) \quad \text{(41c)}$$

Now, from Lemma 2,

$$L \geq \mathbb{E}^{\sigma_t^i, \tilde{\sigma}_t^{-i}}\left[Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right)|z_t, x_t^i\right] \quad \text{(42a)}$$
$$= \mathbb{E}^{\sigma_t^i, \tilde{\sigma}_t^{-i}}\left[R(x_t^i, A_t^i, z_t) + \delta V_{t+1}\left(z_{t+1}, X_{t+1}^i\right)|z_t, x_t^i\right]$$
$$\text{(42b)}$$

which is true as per the definition in (13). The expression can now be expanded using Lemma 1 and use the assumption we made in (38).

$$L \geq \mathbb{E}^{\sigma_t^i, \tilde{\sigma}_t^{-i}}\Bigg[R(z_t, x_t^i, a_t^i) +$$

$$\delta \mathbb{E}^{\tilde{\sigma}} \left[ \sum_{k=t}^{T} \delta^{k-t} R\left(x_k^i, A_k^i, z_k\right) | z_t, x_t^i \right] \right] \tag{43a}$$

$$\geq \mathbb{E}^{\sigma_t^i, \tilde{\sigma}_t^{-i}} \left[ R(z_t, x_t^i, a_t^i) + \right.$$

$$\left. \delta \mathbb{E}^{\sigma^i \tilde{\sigma}^{-i}} \left[ \sum_{k=t+1}^{T} \delta^{k-t-1} R\left(x_k^i, A_k^i, z_k\right) | z_{1:t+1}, x_{1:t+1}^i \right] \right] \tag{43b}$$

$$= \mathbb{E}^{\sigma^i, \tilde{\sigma}^{-i}} \left[ \sum_{k=t}^{T} \delta^{k-t} R\left(X_k^i, A_k^i, z_k\right) | z_{1:t}, x_{1:t}^i, a_{1:t}^i \right] \tag{43c}$$

■

*Theorem 2:* Let $\tilde{\sigma}$ be an MFE of the mean field game. Then there exists an equilibrium generating function $\tilde{\theta}$ that satisfies (14) in the backward recursion algorithm such that $\tilde{\sigma}$ is defined using $\tilde{\theta}$.

*Proof:* Given that $\tilde{\sigma}$ is the MFE of the game, the equilibrium generating function such that

$$\tilde{\theta}[z] = \tilde{\sigma}\left(\cdot | z, \cdot\right) \quad \forall z \tag{44}$$

At time $t$, let us assume that $\tilde{\gamma}_t \left(= \tilde{\theta}\left[z_t\right]\right)$ be the prescription function that gives the equilibrium action but does not satisfy the fixed point equation in (14), i.e.

$$\tilde{\gamma}_t \notin \arg\max_{\gamma_t} \mathbb{E}^{\gamma_t} \left[ Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right) \right] \tag{45}$$

Let us assume that $\widehat{\gamma}_t$ be a solution to the fixed point equation in the RL algorithm, then

$$\mathbb{E}^{\widehat{\gamma}_t} \left[ Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right) \right] \geq \mathbb{E}^{\tilde{\gamma}_t} \left[ Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right) \right] \tag{46}$$

But, we know

$$\mathbb{E}^{\widehat{\gamma}_t} \left[ Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right) \right] \tag{47a}$$

$$= \mathbb{E}^{\widehat{\sigma}_t} \left[ R_t\left(x_t^i, A_t^i, z_t\right) + V_{t+1}\left(z_{t+1}, x_{t+1}^i\right) \right] \tag{47b}$$

$$= \mathbb{E}^{\widehat{\sigma}_t, \tilde{\sigma}_{t+1:T}} \left[ \sum_{n=t}^{T} \delta^{n-t} R\left(X_n^i, A_n^i, Z_n\right) | z_{1:t}, x_{1:t}^i, a_{1:t}^i \right] \tag{47c}$$

$$\geq \mathbb{E}^{\tilde{\gamma}_t} \left[ Q_t\left(z_t, x_t^i, A_t^i, \tilde{\gamma}_t\right) \right] \tag{47d}$$

$$= V_t\left(z_t, x_t^i\right) \tag{47e}$$

which is contradictory on the consideration that $\tilde{\sigma}$ is a MPE of the game. ■

*Assumption 1:* Let the reward function $R\left(x_t^i, a_t^i, z_t\right)$ and the state transition matrix $\tau\left(x_{t+1}^i | x_t^i, a_t^i, z_t\right)$ be continuous functions in $z_t$.

*Theorem 3:* Under the assumption 1, there exists a solution to the fixed point equation in (14).

*Proof:* It has already been shown that when the reward function is bounded, there exists a MFE for finite horizon games. It can also been shown that the solutions to the fixed point equations are MFE of the game. Thus there exists an equilibrium policy for the game achieved at each time $t$. ■

## VI. NUMERICAL EXAMPLE

### A. Security of cyber-physical system: Malware Spread

We consider a security problem in a cyber physical network with positive externalities. It is a discrete version of the malware problem presented in [12], [13], [14], [15]. Some other applications of this model include flu vaccination, entry and exit of firms, investment, network effects. In this model, we suppose there are large number of cyber-physical nodes where each node has a private state $x_t^i \in \{0, 1\}$ where $x_t^i = 0$ represent 'healthy' state and $x_t^i = 1$ is the infected state. Each node can take an action $a_t^i \in \{0, 1\}$, where $a_t^i = 0$ implies "do nothing" whereas $a_t^i = 1$ implies "repair". The dynamics($= Q_x\left(\cdot | \cdot\right)$) are given by

$$x_{t+1}^i = \begin{cases} x_t^i + (1 - x_t^i) w_t^i & \text{for} \quad a_t^i = 0 \\ 0 & \text{for} \quad a_t^i = 1 \end{cases}$$

where $w_t^i \in \{0, 1\}$ is a binary valued random variable with $P(w_t^i = 1) = q$ representing the probability of a node getting infected. Thus if a node doesn't do anything, it could get infected with certain probability, however, if it takes repair action, it comes back to the healthy state. Each node gets a reward

$$r(x_t^i, a_t^i, z_t) = -(k + z_t(1)) x_t^i - \lambda a_t^i. \tag{48}$$

where $z_t(1)$ is the mean field population state being 1 at time $t$, $\lambda$ is the cost of repair and $(k + z_t(1))$ represents the risk of being infected. We pose it as an infinite horizon discounted dynamic game. We consider parameters $k = 0.2, \lambda = 0.5, \delta = 0.9, q = 0.9$ for numerical results presented in Figures 1-3.
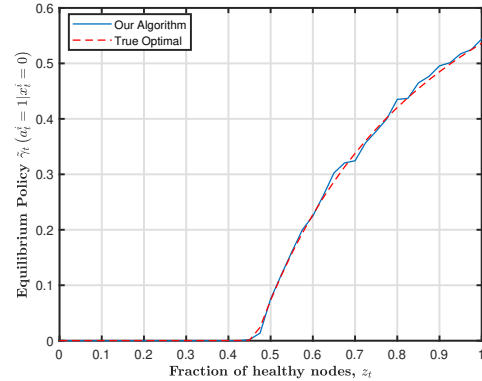


Fig. 1. $\gamma(1|0)$: Probability of choosing action 1, given $x^i = 0$

The learning parameter $\alpha$ for the sarsa update was set at 0.1. The overall length of the time-horizon $T$ was chosen to be 60 iterations long.

Figure 1 and Figure 2 show the equilibrium policies $\tilde{\gamma}$ at different values of the mean state $z_t$ for states $x_t^i = 0$ and $x_t^i = 1$. The plotted graphs are the probabilities with which we choose action $a_t^i = 1$. The plots of our algorithm are compared across the true strategy that was obtained by assuming the knowledge of the dynamics of MDP and then solving the fixed point equation. The strategies estimated
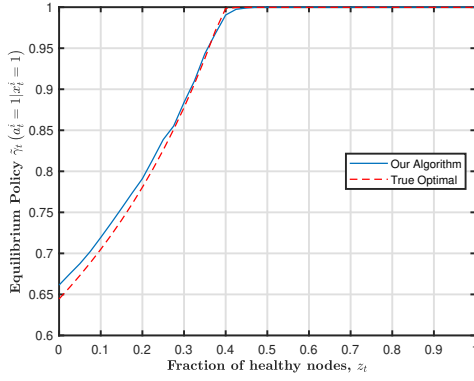
Fig. 2. $\gamma(1|1)$: Probability of choosing action 1, given $x^i = 1$

using the proposed RL algorithm coincides with the true strategies establishing the accuracy of our algorithm.

## VII. CONCLUSION

We considered a finite horizon discrete-time sequential MFG with infinite homogeneous players. The players had access to their private type and the common information of mean population state. A fixed point decomposition method was suggested in an earlier paper that computes the equilibrium strategy at different mean states but with the knowledge of the dynamics of the MDP. Here, we proposed a RL algorithm that employs Expected Sarsa to learn the dynamics of the game and solve the fixed point equation, iteratively, to arrive at the equilibrium strategy. In the end, we implement our algorithm on a practical cyber-physical application to demonstrate that the algorithm does converge to the same optimal policy that was obtained when dynamics of the game was known. We also analytically show the convergence of our algorithm to the MFE of the game. To the best of our knowledge, this is the first RL algorithm to learn optimal policies of non-stationary mean field games.

## REFERENCES

[1] D. Vasal, "Signaling equilibria in mean-field games," pp. 1–21, 2019.
[2] M. Huang, R. P. Malhamé, P. E. Caines *et al.*, "Large population stochastic dynamic games: closed-loop mckean-vlasov systems and the nash certainty equivalence principle," *Communications in Information & Systems*, vol. 6, no. 3, pp. 221–252, 2006.
[3] J.-M. Lasry, P.-L. Lions, J.-M. Lasry, and P.-L. Lions, "Mean field games," *Japan. J. Math*, vol. 2, pp. 229–260, 2007.
[4] X. Guo, A. Hu, R. Xu, and J. Zhang, "Learning mean-field games," in *Advances in Neural Information Processing Systems*, 2019, pp. 4967–4977.
[5] N. Tiwari, A. Ghosh, and V. Aggarwal, "Reinforcement learning for mean field game," *arXiv preprint arXiv:1905.13357*, 2019.
[6] Y. Yang, R. Luo, M. Li, M. Zhou, W. Zhang, and J. Wang, "Mean field multi-agent reinforcement learning," *arXiv preprint arXiv:1802.05438*, 2018.
[7] J. Subramanian, "Reinforcement learning in stationary mean-field games," *Proceedings of the 18th International Conference of Autonomous Agents and Multi-Agent Systems, AAMAS'19*, pp. 251–259, 2019.
[8] R. Elie, J. Pérolat, M. Laurière, M. Geist, and O. Pietquin, "Approximate fictitious play for mean field games," *arXiv preprint arXiv:1907.02633*, 2019.
[9] E. Maskin and J. Tirole, "Markov perfect equilibrium. I. Observable actions," *Journal of Economic Theory*, vol. 100, no. 2, pp. 191–219, 2001.
[10] A. Nayyar, A. Mahajan, and D. Teneketzis, "Decentralized Stochastic Control with Partial History Sharing: A Common Information Approach," *IEEE Transactions on Automation Control*, vol. 58, no. 7, 2013.
[11] H. Van Seijen, H. Van Hasselt, S. Whiteson, and M. Wiering, "A theoretical and empirical analysis of expected sarsa," in *2009 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning*. IEEE, 2009, pp. 177–184.
[12] M. Huang and Y. Ma, "Mean field stochastic games with binary actions: Stationary threshold policies," in *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*. IEEE, 2017, pp. 27–32.
[13] ——, "Mean field stochastic games: Monotone costs and threshold policies," in *2016 IEEE 55th Conference on Decision and Control (CDC)*. IEEE, 2016, pp. 7105–7110.
[14] ——, "Mean field stochastic games with binary action spaces and monotone costs," *arXiv preprint arXiv:1701.06661*, 2017.
[15] L. Jiang, V. Anantharam, and J. Walrand, "How bad are selfish investments in network security?" *IEEE/ACM Transactions on Networking*, vol. 19, no. 2, pp. 549–560, 2010.