

# MODEL-FREE MEAN-FIELD REINFORCEMENT LEARNING: MEAN-FIELD MDP AND MEAN-FIELD Q-LEARNING

RENÉ CARMONA, MATHIEU LAURIÈRE & ZONGJUN TAN

**ABSTRACT.** We study infinite horizon discounted Mean Field Control (MFC) problems with common noise through the lens of Mean Field Markov Decision Processes (MFMDP). We allow the agents to use actions that are randomized not only at the individual level but also at the level of the population. This common randomization allows us to establish connections between both closed-loop and open-loop policies for MFC and Markov policies for the MFMDP. In particular, we show that there exists an optimal closed-loop policy for the original MFC. Building on this framework and the notion of state-action value function, we then propose reinforcement learning (RL) methods for such problems, by adapting existing tabular and deep RL methods to the mean-field setting. The main difficulty is the treatment of the population state, which is an input of the policy and the value function. We provide convergence guarantees for tabular algorithms based on discretizations of the simplex. Neural network based algorithms are more suitable for continuous spaces and allow us to avoid discretizing the mean field state space. Numerical examples are provided.

**Key words.** Mean field reinforcement learning, Mean field Markov Decision Processes, McKean-Vlasov control

**AMS subject classification.** 65M12, 65M99, 93E20, 93E25

## CONTENTS

1. <b>Introduction</b>	2
2. <b>Model Description and Notations</b>	3
2.1. Probabilistic set-up of the MFC model	3
2.2. Probabilistic framework and classes of policies	6
2.3. Optimization and value functions	9
3. <b>Mean-Field MDP</b>	9
3.1. Mean-field MDP framework	10
3.2. Assumptions and optimization problem for MFMDP	11
3.3. Dynamic Programming principle for MFMDP	12
4. <b>Relations between the models</b>	13
4.1. Relations between MFC closed-loop policies and MFMDP policies	14
4.2. Relations between MFC closed-loop and open-loop policies	16
5. <b>Mean-Field Q-Learning</b>	18
5.1. State-action value function	18
5.2. Controls for finite state and action spaces	20
5.3. Simplex discretization and tabular MFQ-learning	21
5.4. Deep reinforcement learning for MFMDP	27
6. <b>Numerical Examples</b>	28

---

This work has been supported by NSF grant DMS-1716673 and ARO grant W911NF-17-1-0578.

6.1. Example 1: Cyber security model	28
6.2. Example 2: Discrete distribution planning	31
6.3. Example 3: Swarm motion	33
References	34
Appendix A. Auxiliary results for Section 2	40
Appendix B. Proofs for Section 4.1	42
Appendix C. Proofs for Section 4.2	43
Appendix D. Disintegration of kernels	43
Appendix E. Details on Q-learning, Section 5	43
E.1. Proof of Theorem 36	43
E.2. DDPG algorithm	45

## 1. Introduction

In today's highly connected world, important applications often involve very large number of interacting rational agents. Understanding how individual decisions aggregate to create global outcome and how, in turn, the agents react to those outcomes to adjust their behavior is a major challenge for numerous applications. Theoretical analyses distinguish between competitive and cooperative scenarios using game-theoretic notions. From a computational viewpoint, solving games with multiple players becomes infeasible when their number grows, in no small part because the number of pairwise interactions increases exponentially. To cope with this issue, mean field games (MFG) and mean field control (MFC) problems, also called McKean-Vlasov (MKV) control have been introduced (see e.g. [31, 28, 9, 12, 13]). Assuming that the population is homogeneous (the agents have the same transition and cost functions), and that the interactions are symmetric (these functions depend only on the empirical distribution of the other agents), the main idea is to use a mean-field approximation of the population's state in order to simplify the model. It is then sufficient to study the interactions between one representative player and the population distribution rather than the interactions between every pair of players. In the past decade, theoretical results and potential applications have received a growing level of interest. Numerical methods, which are a crucial tool for applications, have also been developed, mostly based on deterministic methods for partial differential equations (see e.g. [2, 1, 3]). Computational methods based on deep learning have recently been introduced and seem particularly suitable to tackle high-dimensional MFG and MFC problems (see e.g. [16, 15, 5, 24, 21, 35, 4]).

In the past few years, the question of learning solutions of mean field problems in a model-free way has gained momentum (see e.g., [36, 27, 22, 18, 25, 26, 7, 34, 33, 14]). Roughly speaking, the main point is to develop computational methods that can compute MFG or MFC solutions only by sampling realizations of trajectories and without having access to the model (i.e., the transition and cost functions). Model-free methods for a single agent control problem have been developed in the framework of reinforcement learning (RL), building on the formalism of Markov Decision Processes (MDP). One of the most well-known methods is the so-called Q-learning which exploits the dynamic programming principle satisfied by the state-action value function (also called Q-function). Multi-agent reinforcement learning (MARL) extends RL methods to situations in which several agents are simultaneously learning. Breakthrough results have been obtained in games with a small number of players (such as chess or go).

In this work, we focus on a setting with an infinite cooperative population modeled by an MFC problem in discrete time with an infinite horizon discounted cost. Not only is the dynamics subject to idiosyncratic and common noise, but the actions too are subject to both idiosyncratic and common randomness. The problem can be interpreted as one posed to a central planner who helps a very large population of agents to minimize its social cost and, to this end, can first sample a random policy for the whole population before letting each agent sample their action based on this common policy. This is a distinctive feature compared with the existing literature, in particular [25, 26] which has studied RL for MFC without common noise and [33] which has studied MFMDP in the presence of common noise but without common randomization. We use this extra source of randomness to connect open-loop and closed-loop policies for the MFC problem to a mean field MDP (MFMDP) whose state is the population distribution. Defining properly this MFMDP and dealing rigorously with common noise and common policy randomization leads us to carry out a careful probabilistic analysis of this type of problems. Besides the development of the theoretical framework, we also investigate how RL methods can be adapted to the MFC setting, using tabular methods or neural network based methods. Since policies are common to the whole population in the MFMDP, randomization used in RL methods translates into common randomization from the point of view of the MFC problem. We illustrate the performance of two methods on several numerical examples.

The main contributions are threefold. **(1)** We introduce a MFMDP and study its connection with the original MFC problem: we prove a DPP for the MFMDP value function (Theorem 19), on which we build to enable to prove equality of the open-loop and closed-loop value functions for the MFC problem (Theorem 27). Furthermore, we show existence of a stationary closed-loop policy (Proposition 25). **(2)** We study the state-action value function (or Q-function) of the MFMDP, for which we prove a DPP (Theorem 30). **(3)** We propose several RL methods: a tabular Q-learning relying on a discretization of the mean-field state simplex (Theorem 35), and a deep RL method to deal with continuous state or action spaces, which allows us to avoid simplex discretization or to deal with randomized actions.

The rest of the paper is organized as follows. Section 2 introduces the main concepts for the MFC problem, including the probabilistic framework and the notions of open-loop and closed-loop policies. Then, the corresponding MFMDP and its DPP are given in Section 3. The connections between the MFC and the MFMDP are developed in Section 4. We then turn our attention to the numerical aspects. In Section 5, we introduce the state-action value function, prove it satisfies a DPP, and propose computational methods based on RL. Several numerical examples are provided in Section 6 to illustrate original features of the MFMDP with common noise and randomized actions.

## 2. Model Description and Notations

Throughout the paper we work with Borel spaces, namely spaces homeomorphic to a non-empty Borel subset of some Polish space. If  $C$  is such a space, we denote by  $\mathcal{B}_C$  its Borel  $\sigma$ -field and  $\mathcal{P}(C)$  the space of probability measures on  $(C, \mathcal{B}_C)$  implicitly assumed to be equipped with the topology of the weak convergence and its corresponding Borel  $\sigma$ -field  $\mathcal{B}_{\mathcal{P}(C)}$ . For all the measurability issues we refer the reader to any of the textbooks [10] or [29].

### 2.1. Probabilistic set-up of the MFC model.

In this section, we specify what we mean by mean-field models with common noise. We first introduce the major building blocks, leaving the description of the dynamics for later on.

**Definition 1** (MFC model). *An infinite horizon discounted **mean-field control (MFC) model** with common noise is based on the following elements  $(S, A, E, E^0, F, f, \gamma)$ :*

- A Borel space  $(S, \mathcal{B}_S)$  for the state space.
- A Borel space  $(A, \mathcal{B}_A)$  for the action space.
- Two Borel spaces  $(E, \mathcal{B}_E)$  and  $(E^0, \mathcal{B}_{E^0})$  for the values of the idiosyncratic and common noise.
- A Borel measurable function  $F : S \times A \times \mathcal{P}(S \times A) \times E \times E^0 \rightarrow S$  called the system function.
- A bounded Borel measurable function  $f : S \times A \times \mathcal{P}(S \times A) \rightarrow \mathbb{R}$  called the one-stage cost function.
- A discount factor  $\gamma \in (0, 1)$ .

The system function  $F$  is used to describe the evolution of the state process based on a state, an action, a mean-field interaction term<sup>1</sup> and two noise terms.

Even if we are willing to postpone the regularity conditions on  $F$  and  $f$ , the above definition is still incomplete. It introduces the building blocks of the model, but does not explain how the actions are taken and how the system evolves over time. In order to motivate the nature of the assumptions we are about to introduce, we take an informal excursion in the world of finitely many actors. Our goal is to motivate the following important features of our model:

- (i) The mean-field interactions are conditioned on all shared information.
- (ii) The actions are randomized by additional sources of randomness.

**2.1.1. Motivation from Finitely Many Player Models.** Let us imagine that  $N$  robots with states  $X_n^1, \dots, X_n^N$ , take actions  $\alpha_n^1, \dots, \alpha_n^N$  at time  $n$  and that their next state is given by:

$$X_{n+1}^i = F\left(X_n^i, \alpha_n^i, \frac{1}{N} \sum_{j=1}^N \delta_{(X_n^j, \alpha_n^j)}, \epsilon_{n+1}^i\right), \quad n \geq 0, \quad i = 1, \dots, N$$

$\epsilon_n^1, \dots, \epsilon_n^N, \dots$  being independent and identically distributed (i.i.d. from now on) random shocks with distribution  $\nu \in \mathcal{P}(E)$ . Let us also assume that at each time  $n$ , a central unit collects the information about the states, the actions and the costs (or rewards) incurred by the individual robots, and minimizes the overall average cost to the system as given by:

$$\frac{1}{N} \sum_{j=1}^N \mathbb{E} \left[ \sum_{n=0}^{\infty} \gamma^n f(X_n^j, \alpha_n^j, \frac{1}{N} \sum_{i=1}^N \delta_{(X_n^i, \alpha_n^i)}) \right].$$

In the limit  $N \rightarrow \infty$ , using standard propagation of chaos arguments, we expect that the coupled evolutions of the states of the individual robots will become independent, and that each state evolves according to the dynamics:

$$X_{n+1} = F(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}, \epsilon_{n+1}), \quad n \geq 0,$$

and one can then imagine that the optimization of the central unit reduces to the minimization:

$$\inf_{\alpha = (\alpha_n)_{n \geq 0}} \mathbb{E} \left[ \sum_{n=0}^{\infty} \gamma^n f(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}) \right],$$

where  $\mathbb{P}_{(X_n, \alpha_n)}$  denotes the joint law of the state-action couple at time  $n$ . This is the formulation of (discrete time) MFC problem, which we will not call and MDP because of the lack of Markov property due to the presence of  $\mathbb{P}_{(X_n, \alpha_n)}$  in the equation, responsible for the McKean-Vlasov nature of the dynamics.

---

<sup>1</sup>The interactions are through the joint state-action distribution, which is sometimes referred to as “extended” MFC or MFC of controls.

Still, such a model does not account for the fact that the robots evolve in an environment which is most likely random. Since all the robots face the same random shocks due to the randomness of the common environment, we introduce a sequence  $(\epsilon_n^0)_{n \geq 1}$  of i.i.d. random elements in  $E^0$ , with common law  $\nu^0 \in \mathcal{P}(E^0)$ , independent of the idiosyncratic shocks  $(\epsilon_n^i)_{n \geq 1, i=1, \dots, N}$  of the individual robots, and with this addition to the model, the individual robot state dynamics become:

$$X_{n+1}^i = F\left(X_n^i, \alpha_n^i, \frac{1}{N} \sum_{j=1}^N \delta_{(X_n^j, \alpha_n^j)}, \epsilon_{n+1}^i, \epsilon_{n+1}^0\right), \quad n \geq 0, \quad i = 1, \dots, N.$$

Even though we added significant sources of coupling between the robots experiences, exchangeability among the robots persists, and in the limit  $N \rightarrow \infty$ , a conditional form of the propagation of chaos theory should lead to conditional independence of individual dynamics which should take the generic form:

$$X_{n+1} = F(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}^0, \epsilon_{n+1}, \epsilon_{n+1}^0), \quad n \geq 0,$$

$\mathbb{P}_{(X_n, \alpha_n)}^0$  being the conditional law of the state-action couple  $(X_n, \alpha_n)$  given the common noise  $(\epsilon_n^0)_{n \geq 1}$ . Accordingly, the optimization of the central unit should be:

$$\inf_{\alpha = (\alpha_n)_{n \geq 0}} \mathbb{E} \left[ \sum_{n=0}^{\infty} \gamma^n f(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}^0) \right].$$

This form of Mean Field discrete time Control problem with common noise is now very close to the model we investigate in this paper. The last question we would like to address before turning to the theoretical analysis of the model is: *Could the conditional distributions  $\mathbb{P}_{(X_n, \alpha_n)}^0$  depend upon some extra sources of randomness?* Indeed, if individual robots and the central unit are allowed to use mixed strategies, they need independent sources of randomness to randomize their actions and decisions at each time. So since the cost optimization and the choice of an action strategy are performed by the central unit,

- the central unit needs a source of randomness to *randomize* the choice of a mixed policy to be dispatched to the individual robots, and
- the individual robots need to sample their actions from the policy sampled for them by the central unit.

So in the limit  $N \rightarrow \infty$ , the *idiosyncratic* randomizations of the individual robots average out, while the *central unit randomization* remains.

As a result, the conditioning in  $\mathbb{P}_{(X_n, \alpha_n)}^0$  should be with respect to the common noise  $\epsilon_1^0, \dots, \epsilon_n^0$  as well as the sources of randomness used by the central unit to randomize the policy handed out to the individual robots for implementation.

**2.1.2. Back to the MFC problem formulation.** Motivated by the previous  $N$ -agent model discussion when the size of the population  $N$  tends to infinity, we propose the following set-up to disentangle clearly the various sources of randomness.

We assume that all the sources of randomness are from a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  supporting

- (i) an i.i.d. sequence  $(\epsilon_n)_{n \geq 0}$  with distribution  $\nu \in \mathcal{P}(E)$  modeling the idiosyncratic random shocks;
- (ii) an i.i.d. sequence  $(\epsilon_n^0)_{n \geq 0}$  with distribution  $\nu^0 \in \mathcal{P}(E^0)$  modeling the common noise;
- (iii) a random variable  $\mathcal{U}$  with distribution  $\mathbb{P}_{\mathcal{U}}$  in a Borel space  $(\Upsilon, \mathcal{B}_{\Upsilon})$  providing the randomization for the initial state;

- (iv) an i.i.d. sequence  $(\vartheta_n)_{n \geq 0}$  of random variables in a Borel space  $(\Theta, \mathcal{B}_\Theta)$  with distribution  $\mathbb{P}_\vartheta$  providing the generic robot with a source of randomization for their action choices;
- (v) an i.i.d. sequence  $(\vartheta_n^0)_{n \geq 0}$  of random variables in a Borel space  $(\Theta^0, \mathcal{B}_{\Theta^0})$  with distribution  $\mathbb{P}_{\vartheta^0}$  providing the central unit with a source of randomization for the choices of policies.

We assume that all these random sequences are independent of each other. We also assume that  $\mathbb{P}_\vartheta$  and  $\mathbb{P}_{\vartheta^0}$  are both atomless. This guarantees the existence of Borel measurable functions  $h : \Theta \mapsto [0, 1]$  and  $h^0 : \Theta^0 \mapsto [0, 1]$  which are uniformly distributed when viewed as random variables on the probability spaces  $(\Theta, \mathcal{B}_\Theta, \mathbb{P}_\vartheta)$  and  $(\Theta^0, \mathcal{B}_{\Theta^0}, \mathbb{P}_{\vartheta^0})$  respectively. The uniform random variables constructed with the functions  $h$  and  $h^0$  will be used repeatedly with the following classical result from measure theory which we state for the sake of later reference.

**Lemma 2** (Blackwell-Dubins Lemma). *For any Polish space  $B$ , there exists a measurable function  $\rho_B : \mathcal{P}(B) \times [0, 1] \mapsto B$ , which we shall call the Blackwell-Dubins function of the space  $B$ , satisfying*

- i) for each  $\nu \in \mathcal{P}(B)$  and each uniform random variable  $U \sim U(0, 1)$ , the  $B$ -valued random variable  $\rho_B(\nu, U)$  has distribution  $\nu$ ;*
- ii) for almost every  $u \in [0, 1]$ , the function  $\nu \mapsto \rho_B(\nu, u)$  is continuous for the weak topology of  $\mathcal{P}(B)$ .*

For the sake of illustration, some of the computations will be done in the canonical probability space  $(\Omega^c, \mathcal{F}^c, \mathbb{P}^c)$  defined by:

$$\Omega^c = \Upsilon \times \Theta \times \Theta^0 \times (E \times E^0 \times \Theta \times \Theta^0)^\infty, \quad \mathcal{F}^c = \mathcal{B}_\Upsilon \times \mathcal{B}_\Theta \times \mathcal{B}_{\Theta^0} \times (\mathcal{B}_E \times \mathcal{B}_{E^0} \times \mathcal{B}_\Theta \times \mathcal{B}_{\Theta^0})^\infty,$$

and the product probability  $\mathbb{P}^c$  given by:

$$\mathbb{P}^c = \mathbb{P}_\mathcal{U} \otimes \mathbb{P}_\vartheta \otimes \mathbb{P}_{\vartheta^0} \otimes (\nu \otimes \nu^0 \otimes \mathbb{P}_\vartheta \otimes \mathbb{P}_{\vartheta^0})^\infty.$$

A generic element  $\omega \in \Omega$  reads

$$\omega = (u, \theta_0, \theta_0^0, e_1, e_1^0, \theta_1, \theta_1^0, e_2, e_2^0, \dots, e_n, e_n^0, \theta_n, \theta_n^0, \dots),$$

and we realize the random variables as coordinate mappings,  $\mathcal{U}(\omega) = u$ , and for  $n \geq 0$ ,

$$\vartheta_n(\omega) = \theta_n, \quad \vartheta_n^0(\omega) = \theta_n^0, \quad \varepsilon_{n+1}(\omega) = e_{n+1}, \quad \varepsilon_{n+1}^0(\omega) = e_{n+1}^0.$$

We will use the short hand notations  $\underline{\vartheta}_n = (\vartheta_0, \dots, \vartheta_n)$ ,  $\underline{\vartheta}_n^0 = (\vartheta_0^0, \dots, \vartheta_n^0)$  for every  $n \geq 0$ , and  $\underline{\varepsilon}_n = (\varepsilon_1, \dots, \varepsilon_n)$ ,  $\underline{\varepsilon}_n^0 = (\varepsilon_1^0, \dots, \varepsilon_n^0)$  for every  $n \geq 1$ .

## 2.2. Probabilistic framework and classes of policies.

**2.2.1. Filtrations, action processes, and control processes.** We now introduce the framework that will be used to rigorously study the MFC problem. As explained intuitively in the introduction, we will distinguish several types of randomness. We will also distinguish between actions (elements of  $A$ ), and controls (probability measures on  $A$ ). These are the building blocks to define later the notion of policy (see § 2.2.3).

The  $\sigma$ -field for the initial state is denoted by  $\mathcal{F}_{x_0} = \sigma\{\mathcal{U}\}$ . We introduce four filtrations to define the action and control processes. The filtrations of the idiosyncratic and common noises are:

$$\mathcal{F}_0^\varepsilon = \mathcal{F}_0^{\varepsilon^0} = \{\emptyset, \Omega\}, \quad \mathcal{F}_n^\varepsilon = \sigma\{\underline{\varepsilon}_n\}, \quad \mathcal{F}_n^{\varepsilon^0} = \sigma\{\underline{\varepsilon}_n^0\}, \quad n \geq 1.$$

The filtration of the (idiosyncratic) action randomization is:

$$\mathcal{F}_n^\Theta = \sigma\{\Theta_0, \dots, \Theta_n\}, \quad n \geq 0.$$

The filtration of the (common) policy randomization is:

$$\mathcal{F}_n^{\Theta^0} = \sigma\{\Theta_0^0, \dots, \Theta_n^0\}, \quad n \geq 0.$$

We also introduce three new filtrations  $\mathbb{F}^0 = (\mathcal{F}_n^0)_{n \geq 0}$ ,  $\mathbb{G}^c = (\mathcal{G}_n^c)_{n \geq 0}$  and  $\mathbb{G}^a = (\mathcal{G}_n^a)_{n \geq 0}$  defined by:

$$\begin{aligned} \mathcal{F}_0^0 &= \sigma\{\vartheta_0^0\} & \mathcal{F}_n^0 &= \mathcal{F}_n^{\varepsilon^0} \vee \mathcal{F}_n^{\Theta^0}, = \sigma\{\vartheta_n^0, \varepsilon_n^0\}, \quad n \geq 1, \\ \mathcal{G}_0^c &= \sigma\{\mathcal{U}, \vartheta_0^0\}, & \mathcal{G}_n^c &= \mathcal{F}_{x_0} \vee \mathcal{F}_n^{\varepsilon} \vee \mathcal{F}_n^{\varepsilon^0} \vee \mathcal{F}_n^{\Theta^0} \vee \mathcal{F}_{n-1}^{\Theta} = \sigma\{\mathcal{U}, \vartheta_{n-1}, \vartheta_n^0, \varepsilon_n, \varepsilon_n^0\}, \quad n \geq 1, \\ \mathcal{G}_0^a &= \sigma\{\mathcal{U}, \vartheta_0^0, \vartheta_0\}, & \mathcal{G}_n^a &= \mathcal{F}_{x_0} \vee \mathcal{F}_n^{\varepsilon} \vee \mathcal{F}_n^{\varepsilon^0} \vee \mathcal{F}_n^{\Theta^0} \vee \mathcal{F}_n^{\Theta} = \sigma\{\mathcal{U}, \vartheta_{n-1}, \vartheta_n^0, \varepsilon_n, \varepsilon_n^0, \vartheta_n\}, \quad n \geq 1. \end{aligned}$$

Next, we introduce the terminology which is going to help us characterize the information available to a generic robot and the central unit to make their choices of actions and policies. Incidentally, we start referring to the robots as agents and everything related to agents (e.g. actions, controls, policies, costs, etc) will be referred as level-0. This is in contrast with the lifted stochastic optimization model which we introduce in Section 3 whose elements will be referred to as level-1. For reasons which will become clear later, the central unit controlling the robots will be called the *level-1 controller*. This new terminology frees our presentation of the theoretical results from the gory details of the robotic application we used to motivate the set-up.

**Definition 3.** A **level-0 action** is an element of  $A$ . A **level-0 (mixed) control** is a random probability measure on  $(A, \mathcal{B}_A)$ , that is, any random variable with values in the Borel space  $(\mathcal{P}(A), \mathcal{B}_{\mathcal{P}(A)})$ .

Typically, a level-0 action is denoted by  $a$  and a level-0 control is denoted by  $\mathbf{a}$ . Unless specified otherwise, all the controls we consider are mixed. So for the sake of brevity, we will omit this term.

We now consider the notion of action and control processes for a representative agent. Intuitively, an action process is the realization of a control process, where the sampling is done using  $(\vartheta_n)_{n \geq 0}$ .

**Definition 4.** A **level-0 action process** is a sequence of random variables  $\boldsymbol{\alpha} = (\alpha_n)_{n \geq 0}$  with values in  $A$  which is adapted to the filtration  $\mathbb{G}^a$ . The set of such action processes is denoted by  $\mathbb{A}$ . A **level-0 control process** is a sequence  $\mathbf{a} = (\mathbf{a}_n)_{n \geq 0}$  of level-0 controls which is adapted to the filtration  $\mathbb{G}^c$ . Finally, an action process  $\boldsymbol{\alpha} = (\alpha_n)_{n \geq 0}$  is said to be a **realization** of a level-0 control process  $\mathbf{a} = (\mathbf{a}_n)_{n \geq 0}$  if  $\mathcal{L}(\alpha_n | \sigma\{\mathcal{U}, \vartheta_{n-1}, \vartheta_n^0, \varepsilon_n, \varepsilon_n^0\}) = \mathbf{a}_n$ ,  $\mathbb{P}$ -a.s., for every  $n \geq 0$ .

Here and in the following, we use the notations  $\mathbb{P}_\xi$  and  $\mathcal{L}(\xi)$  interchangeably for the distribution of a random element  $\xi$ , and we use natural extensions to denote conditional distributions.

It can be shown (see Lemma 39 in the appendix) that, for any level-0 control process  $\mathbf{a}$  and any two realizations  $\boldsymbol{\alpha}, \boldsymbol{\alpha}'$  of  $\mathbf{a}$ , every bounded Borel measurable function  $h$ ,  $\mathbb{E}[h(\alpha'_n) | \mathcal{G}_n^c] = \int_A h(\alpha) \mathbf{a}_n(d\alpha) = \mathbb{E}[h(\alpha_n) | \mathcal{G}_n^c]$ ,  $\mathbb{P}$ -a.s.,  $n \geq 0$ .

**2.2.2. Conditional distribution and state process.** We are now in a position to describe precisely the mean-field interactions in the system function, and provide a clear definition of the state process driven by a mixed control process in the mean-field model with common noise.

**Definition 5.** For any initial distribution  $\mu_0 \in \mathcal{P}(S)$  and level-0 action process  $\boldsymbol{\alpha} = (\alpha_n)_{n \geq 0}$ , we say that a process  $\mathbf{X}^{\boldsymbol{\alpha}, \mu_0} = (X_n^{\boldsymbol{\alpha}, \mu_0})_{n \geq 0}$  is a **state process associated to  $(\boldsymbol{\alpha}, \mu_0)$**  for the MFC model if:  $X_0^{\boldsymbol{\alpha}, \mu_0}$  is an  $S$ -valued  $\mathcal{F}_{x_0}$ -random variable with distribution  $\mu_0$ , and for every  $n \geq 0$ ,

$$(1) \quad X_{n+1}^{\boldsymbol{\alpha}, \mu_0} = F(X_n^{\boldsymbol{\alpha}, \mu_0}, \alpha_n, \mathbb{P}_{(X_n^{\boldsymbol{\alpha}, \mu_0}, \alpha_n)}^0, \varepsilon_{n+1}, \varepsilon_{n+1}^0),$$

where  $\mathbb{P}_{(X_n^{\alpha, \mu_0}, \alpha_n)}^0$  is a regular version of  $\mathcal{L}((X_n^{\alpha, \mu_0}, \alpha_n) | \mathcal{F}_n^0)$ , the conditional joint distribution of state-action at time  $n$  with respect to common noise and common randomization up to current time.

Such a state process  $\mathbf{X}^{\alpha, \mu_0}$  is adapted to the filtration  $\mathbb{G}^x = (\mathcal{G}_n^x)_{n \geq 0}$ , defined by

$$(2) \quad \mathcal{G}_0^x = \sigma\{\mathcal{W}\}, \quad \mathcal{G}_n^x = \sigma\{\mathcal{W}, (\vartheta_k, \vartheta_k^0, \varepsilon_{k+1}, \varepsilon_{k+1}^0)_{k=0, \dots, n-1}\}, \quad n \geq 1.$$

For each level-0 action process  $\alpha$  and each  $n \geq 0$ , we denote by  $\mathbb{P}_{X_n^{\alpha, \mu_0}}^0$  a regular version of the conditional distribution  $\mathcal{L}(X_n^{\alpha, \mu_0} | \mathcal{F}_n^0)$ . It holds:

$$(3) \quad \mathbb{P}_{X_n^{\alpha, \mu_0}}^0 = \mathcal{L}(X_n^{\alpha, \mu_0} | \sigma\{\varepsilon_n^0, \vartheta_{n-1}^0\}), \quad \mathbb{P} - a.s..$$

this is due to the fact that  $X_n^{\alpha, \mu_0}$  is  $\mathcal{G}_n^x$ -measurable, hence  $X_n^{\alpha, \mu_0} \perp_{\mathcal{F}_n^0 \vee \mathcal{F}_{n-1}^0} \vartheta_n^0$ .

**2.2.3. Open-loop and closed-loop policies.** We now introduce two concepts of policies: open-loop and closed-loop ones.

We first consider open-loop policies. For each  $n \geq 0$ , let  $\Xi_n = \Upsilon \times (\Theta \times \Theta^0 \times E \times E^0)^n \times \Theta^0$ , with the convention that  $\Xi_0 = \Upsilon \times \Theta^0$ , and let us define  $\xi = (\xi_n)_{n \geq 0}$  by:

$$\xi_0 = (\mathcal{W}, \vartheta_n^0), \quad \text{and} \quad \xi_n = (\mathcal{W}, (\vartheta_k, \vartheta_k^0, \varepsilon_{k+1}, \varepsilon_{k+1}^0)_{k=0, \dots, n-1}, \vartheta_n^0), \quad n \geq 1.$$

**Definition 6.** An **level-0 open-loop policy** is a sequence  $\pi = (\pi_n)_{n \geq 0}$  of deterministic measurable functions  $\pi_n : \Xi_n \rightarrow \mathcal{P}(A)$ , called the **open-loop strategy functions** at time  $n$  of policy  $\pi$ . The set of all open-loop policies is denoted by  $\Pi^{OL}$ . A level-0 control process  $\mathbf{a} = (a_n)_{n \geq 0}$  is said to be **generated by the open-loop policy  $\pi$**  if:

$$(4) \quad a_n = \pi_n(\xi_n), \quad \mathbb{P} - a.s., \quad n \geq 0.$$

Notice that because of measurability restrictions ( $a_n$  must be adapted to the filtration  $\mathbb{G}^c$ ), there is a one-to-one correspondence between level-0 control processes  $\mathbf{a}$  and open-loop policies  $\pi$ , and both objects can be identified through equation (4). We shall use one object or the other depending on whether we want to emphasize the stochastic process or the sequence of deterministic functions defining the process. To be consistent with Definition 4, we shall often short-circuit the control process  $\mathbf{a}$  and say that an action process  $\alpha = (\alpha_n)_{n \geq 0}$  is (a realization of the control process) **generated by  $\pi$**  if  $\mathcal{L}(\alpha_n | \mathcal{G}_n^c) = \pi_n(\xi_n)$ ,  $\mathbb{P} - a.s.$ , for all  $n \geq 0$ .

We now consider closed-loop policies.

**Definition 7.** A **closed-loop Markov strategy function** is a measurable function from  $S \times \mathcal{P}(S) \times \Theta^0$  into  $\mathcal{P}(A)$ . A **closed-loop Markov policy  $\pi = (\pi_n)_{n \geq 0}$**  is a sequence of such functions. The set of all closed-loop Markov policies is denoted by  $\Pi^{CL}$ .

The choice of this form of closed-loop Markov strategy function is suggested by the mean-field nature of the dynamics (1) and the form of the costs (6). Taking values in  $\mathcal{P}(A)$  instead of  $A$  indicates that the strategy functions are mixed. Typically, they suggest that at each time  $n \geq 0$ , the action  $\alpha_n \in A$  taken according to such a policy should be sampled from a probability measure depending directly on the values of  $X_n^{\alpha, \mu_0}$  and  $\mathbb{P}_{X_n^{\alpha, \mu_0}}^0$ , and the random variable  $\vartheta_n^0$  used by the level-1 controller to randomize their choice. We formalize this procedure in Definition 8 below.

**Definition 8.** For a closed-loop Markov policy  $\pi \in \Pi^{CL}$  and an initial distribution  $\mu_0 \in \mathcal{P}(S)$ , a pair of state and action processes  $(\mathbf{X}, \alpha) = (X_n, \alpha_n)_{n \geq 0}$  is said to be **generated by  $(\pi, \mu_0)$**  if

i)  $\mathbf{X}$  is a state process associated to  $(\alpha, \mu_0)$  in the sense of Definition 5.

ii) The action process  $\alpha$  is adapted to  $\mathbb{G}^a$  and satisfies

$$(5) \quad \mathcal{L}(\alpha_n | \mathcal{G}_n^c) = \pi_n(X_n, \mathbb{P}_{X_n}^0, \vartheta_n^0), \quad \mathbb{P} - a.s., \quad n \geq 0.$$

In Definition 8, we can view the state and action processes as constructed simultaneously by alternatively invoking the system dynamics (1) and the sampling procedure consistent with (5). A convenient way to construct an action process  $\alpha$  satisfying (5) is to use the Blackwell-Dubin's Lemma 2. Indeed, if  $\rho_A$  is the Blackwell-Dubin's function of  $A$  and the uniformly distributed random variables  $U_n$  is given by  $U_n = h(\vartheta_n)$ , we can choose  $\alpha_n = \rho_A(\pi_n(X_n, \mathbb{P}_{X_n}^0, \vartheta_n^0), U_n)$ ,  $\mathbb{P}$ -a.s.,  $n \geq 0$ .

**Remark 9.** *Even though we call a policy  $\pi \in \Pi^{CL}$  a “Markov” policy, it does not imply any Markov property for the state process  $\mathbf{X}$  associated to such a policy. This abuse of terminology can be explained by our intention to work with level-1 Markov policies which will imply the Markov property for a lifted measure-valued state process constructed in the next section. Also, since we only use the term “Markov policy” in the closed-loop setting, we shall most often drop the term closed-loop hereafter and only call them simply Markov policies.*

### 2.3. Optimization and value functions.

We now move to the definition of the optimization problems for an MFC model with open-loop and closed-loop Markov policies.

We now introduce the value function associated to an action process.

**Definition 10.** *For every level-0 action process  $\alpha$ , the value function  $J^\alpha : \mathcal{P}(S) \rightarrow \mathbb{R}$  is defined for every  $\mu_0 \in \mathcal{P}(S)$  by:*

$$(6) \quad J^\alpha(\mu_0) := \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n f \left( X_n^{\alpha, \mu_0}, \alpha_n, \mathbb{P}_{(X_n^{\alpha, \mu_0}, \alpha_n)}^0 \right) \right],$$

where the state process  $\mathbf{X}^{\alpha, \mu_0}$  is associated to  $(\alpha, \mu_0)$  according to the dynamics (1).

The value of  $J^\alpha(\mu)$  is well-defined for every  $\mu \in \mathcal{P}(S)$  because  $f$  is measurable and bounded. Furthermore, this value depends only upon the sequence of one-dimensional marginal distributions of the  $\mathcal{P}(S \times A)$ -valued process  $(\mathbb{P}_{(X_n^{\alpha, \mu_0}, \alpha_n)}^0)_{n \geq 0}$ .

For any open-loop or closed-loop policy  $\pi$ , we can show that  $\mathbb{P}_{(X_n, \alpha_n)}^0$  depends on the action process only through the policy  $\pi$  provided  $(\mathbf{X}, \alpha)$  is generated by  $\pi$ . See Lemma 40 in the appendix. As a consequence, we can define, for any level-0 action process  $\alpha$  generated by  $\pi$ ,

$$J^\pi(\mu) = J^\alpha(\mu), \quad \mu \in \mathcal{P}(S).$$

Accordingly, we define the optimal open-loop value function and the optimal closed-loop value function as:

$$J^{OL,*}(\mu) = \inf_{\pi \in \Pi^{OL}} J^\pi(\mu), \quad J^{CL,*}(\mu) = \inf_{\pi \in \Pi^{CL}} J^\pi(\mu), \quad \mu \in \mathcal{P}(S),$$

which are finite because we assume that the one-stage cost function  $f$  is bounded and  $\gamma \in (0, 1)$ .

## 3. Mean-Field MDP

In this section, we introduce a Markov Decision Process (MDP) which we use to identify an optimal closed-loop Markov policy for our original MFC model.

### 3.1. Mean-field MDP framework.

The key observation is that, for a mixed Markov closed-loop policy  $\pi \in \Pi^{CL}$ , the associated value function  $J^\pi$  can be viewed as the value function of an MDP with state space  $\mathcal{P}(S)$ , state process  $(\mathbb{P}_{X_n}^0)_{n \geq 0}$  and action process  $(\mathbb{P}_{(X_n, \alpha_n)}^0)_{n \geq 0}$  with values in the new action space  $\mathcal{P}(S \times A)$ . Actions need to be consistent with the state in the sense that the first marginal of any action which can be taken while in a given state, has to be equal to the state itself. We provide a rigorous definition of this MDP, and we show its connection to the original MFC model.

A **mean-field MDP (MFMDP)** consists of a six-tuple  $(\bar{S}, \bar{A}, \bar{\Gamma}, P, \bar{f}, \gamma)$  as described below:

- The state space is the Borel space  $\bar{S} := \mathcal{P}(S)$ ; a generic element in  $\bar{S}$  is denoted by  $\mu$ .
- The action space is the Borel space  $\bar{A} := \mathcal{P}(S \times A)$ ; a generic element in  $\bar{A}$  is denoted by  $\bar{a}$ .
- The control constraint is a set-valued function  $\bar{U}$  from  $\bar{S}$  into the set of non-empty subsets of  $\bar{A}$  defined by:

$$(7) \quad \bar{U}(\mu) := \{\bar{a} \in \bar{A}; \text{pr}_1(\bar{a}) = \mu\}, \quad \forall \mu \in \bar{S};$$

where  $\text{pr}_1 : \bar{A} \rightarrow \bar{S}$  is the projection function that maps  $\bar{a} \in \bar{A}$  onto its first marginal distribution on  $S$ . We shall also use the notation

$$(8) \quad \bar{\Gamma} := \{(\mu, \bar{a}) \in \bar{S} \times \bar{A}; \bar{a} \in \bar{U}(\mu)\}.$$

- The transition probability kernel  $P : \bar{\Gamma} \rightarrow \mathcal{P}(\bar{S})$ , which is a Borel measurable function.
- The one-stage cost function  $\bar{f} : \bar{\Gamma} \rightarrow \mathbb{R}$ , which is a bounded measurable function.
- The discount coefficient  $\gamma \in (0, 1)$ .

**Remark 11.** The projection map  $\text{pr}_1$  is continuous, so the constraint set  $\bar{U}(\mu)$  is closed in  $\bar{A}$  for every  $\mu \in \bar{S}$ . The graph  $\text{Gr}(\text{pr}_1) := \{(\bar{a}, \mu) : \text{pr}_1(\bar{a}) = \mu\} \subset \bar{A} \times \bar{S}$  is closed, so  $\bar{\Gamma}$  is also closed in  $\bar{S} \times \bar{A}$ . Hence  $\bar{\Gamma}$  is an analytic subset of  $\bar{S} \times \bar{A}$ , and a Polish space on its own. We assume that  $\bar{\Gamma}$  is endowed with the induced topology as well as the trace  $\sigma$ -field inherited from  $\bar{S} \times \bar{A}$ .

**Definition 12.** The six-tuple  $(\bar{S}, \bar{A}, \bar{\Gamma}, P, \bar{f}, \gamma)$  is said to be the **MFMDP lifted from the MFC model**  $(S, A, E, E^0, F, f, \gamma)$  of Definition 1 if it satisfies:

- The transition kernel  $P$  is given by

$$(9) \quad P(\mu, \bar{a})(d\mu') = (\nu^0 \circ \bar{F}(\mu, \bar{a}, \cdot)^{-1})(d\mu'), \quad (\mu, \bar{a}) \in \bar{\Gamma},$$

where  $\bar{F} : \bar{\Gamma} \times E^0 \rightarrow \bar{S}$  is the system function defined in terms of  $F$  by:

$$(10) \quad \bar{F}(\mu, \bar{a}, e^0) = (\bar{a} \otimes \nu) \circ F(\cdot, \cdot, \bar{a}, \cdot, e^0)^{-1}, \quad (\mu, \bar{a}, e^0) \in \bar{\Gamma} \times E^0.$$

- The one-stage cost function  $\bar{f} : \bar{\Gamma} \rightarrow \mathbb{R}$  of the MFMDP satisfies:

$$(11) \quad \bar{f}(\mu, \bar{a}) = \int_{S \times A} f(x, \alpha, \bar{a}) \bar{a}(dx, d\alpha), \quad (\mu, \bar{a}) \in \bar{\Gamma}.$$

Here and in the following we denote by  $\nu \circ g^{-1}$  the push-forward of a measure  $\nu$  by a measurable function  $g$ . We defined the dynamics using a transition kernel  $P$  for two reasons: 1) to conform with the standard literature on MDPs which seems to prefer transition kernels to system functions; 2) we shall restrict ourselves to Markovian policies and control processes given by feedback functions of the state for the analysis of the lifted MDP.

We can check that  $\bar{F}$  is Borel measurable; see e.g. [10, Proposition 7.29] and Remark 41 in the appendix.

**Remark 13.** *In anticipation for what is going to come next, we want to emphasize that the above MFMDP satisfies the assumptions of [10, Chapter 8-9]. We will use the results therein to derive a form of Dynamic Programming Principle (DPP) for the MFMDP.*

To highlight the tight connections between the lifted MFMDP and the original MFC, we define the notion of mixed strategy and mixed Markov policy for MFMDP.

**Definition 14.** *We call **level-1 pure strategy function** any Borel measurable function from  $\bar{S}$  into  $\bar{A}$  whose graph is contained in  $\bar{\Gamma}$ . We call **level-1 mixed strategy function** any Borel measurable function  $\bar{\pi}$  from  $\bar{S}$  into  $\mathcal{P}(A)$  satisfying  $\bar{\pi}(\mu)(\bar{U}(\mu)) = 1$ ,  $\mu \in \bar{S}$ . We denote by  $\bar{\Pi}^p$  (resp.  $\bar{\Pi}$ ) the set of pure (resp. mixed) strategy functions.*

We shall identify every  $\beta \in \bar{\Pi}^p$  with the corresponding level-1 mixed strategy  $\bar{\pi}$  defined by  $\bar{\pi}(\mu) = \delta_{\beta(\mu)}$  for  $\mu \in \bar{S}$ . For the sake of brevity, we sometimes omit the term “mixed” or “randomized” when  $\bar{\pi} \in \bar{\Pi}$ , but we always keep the term “pure” or “non-randomized” when  $\bar{\pi} \in \bar{\Pi}^p$ . We now define policies as sequences of strategy functions.

**Definition 15.** *A **mixed Markov policy** for MFMDP is an element of  $\bar{\Pi} := (\bar{\Pi})^{\mathbb{N}}$ . Similarly, a **pure policy** is an element of  $\bar{\Pi}^p := (\bar{\Pi}^p)^{\mathbb{N}}$ . We say that a policy  $\bar{\pi} = (\bar{\pi}_n)_{n \geq 0}$  is **stationary** if the strategy functions  $\bar{\pi}_n$  are equal for all  $n$ .*

Keeping with the spirit of the previous section, these policies should be called “Markov” policies. We restrict ourselves to these policies and refrain from using history dependent policies because we are mostly interested in optimizing value functions, and we know that for MDPs like our lifted MFMDP, for each history dependent policy, there exists a Markov policy with the same value function (as defined below in Definition 17). See for example [10, Proposition 9.1].

**Definition 16.** *A pair of state and action processes  $(\mu, \bar{a}) = (\mu_n, \bar{a}_n)_{n \geq 0}$  is said to be **generated by**  $(\bar{\pi}, \mu) \in \bar{\Pi} \times \bar{S}$  if the following conditions are satisfied:  $\mu_0 = \mu$ ,*

$$(12) \quad \mu_{n+1} = \bar{F}(\mu_n, \bar{a}_n, \varepsilon_{n+1}^0), \quad \mathbb{P} - a.s. \quad n \geq 0,$$

*and  $\bar{a}$  is an  $\bar{A}$ -valued process adapted to  $\mathbb{F}^0$  satisfying:*

$$(13) \quad \mathcal{L}(\bar{a}_n | \mu_n) = \bar{\pi}_n(\mu_n), \quad \mathbb{P} - a.s. \quad n \geq 0.$$

### 3.2. Assumptions and optimization problem for MFMDP.

We will sometimes rely on the following assumptions in the analysis of the model.

**Assumption (H1).** • **System function  $F$ :** *For every  $(e, e^0) \in E \times E^0$ , the function  $F(\cdot, \cdot, \cdot, e, e^0)$  is continuous in its remaining variables.*

• **One-stage cost function  $f$ :**  *$f : S \times A \times \mathcal{P}(S \times A) \rightarrow \mathbb{R}$  is continuous.*

To show existence of an optimal policy, we will make use of the following extra assumption:

**Assumption (H2). Compactness:** *The state space  $S$  and the action space  $A$  are compact metric spaces.*

In order to obtain a dynamic programming principle (DPP) with Borel measurable mixed Markov policies, measurability issues lead us to work under Assumption (H1) for the MFC model. It can be shown that, under this assumption,  $\bar{F}$  is Borel measurable and for every  $e^0 \in E^0$ ,  $\bar{F}(\cdot, \cdot, e^0)$  is continuous in its remaining variables, and  $\bar{f}$  is bounded and lower semi-continuous. See Lemma 42 in the appendix.

**Definition 17.** For every  $\bar{\pi} \in \bar{\Pi}$ , we define the **value function**  $J^{\bar{\pi}}$  by:

$$(14) \quad \bar{J}^{\bar{\pi}}(\mu) := \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n \bar{f}(\mu_n, \bar{a}_n) \right], \quad \mu \in \bar{S},$$

where  $(\mu, \bar{a}) = (\mu_n, \bar{a}_n)_{n \geq 0}$  is any pair of state and action processes generated by  $(\bar{\pi}, \mu)$ . If  $\bar{\pi} \in \bar{\Pi}$  is stationary with a strategy function  $\bar{\pi} \in \bar{\Pi}$ , then we let  $\bar{J}^{\pi} = \bar{J}^{\bar{\pi}}$ .

It can be shown that the value function  $\bar{J}^{\bar{\pi}}$  given in (14) is well defined because the expectation in (14) does not depend upon the particular choice of the pair of state action processes  $(\mu, \bar{a})$  generated by  $(\bar{\pi}, \mu)$ . See Lemma 43 in the appendix.

**Remark 18.** Using measurability arguments found for example in [10, Chapter 7], one can check that  $\bar{J}^{\bar{\pi}}$  is Borel measurable when  $\bar{\pi} \in \bar{\Pi}$ .

With the value function for MFMDP at hand, we define the **optimal value function of the MFMDP** as:

$$\bar{J}^*(\mu) = \inf_{\bar{\pi} \in \bar{\Pi}} \bar{J}^{\bar{\pi}}(\mu), \quad \mu \in \bar{S}.$$

**3.3. Dynamic Programming principle for MFMDP.** We state and prove the Dynamic Programming principle for the optimal value function with Borel measurable mixed Markov policies.

**Theorem 19.** Assume that **(H1)** and **(H2)** hold. Then, the function  $\bar{J}^*$  is bounded and lower semi-continuous, and moreover it is the unique bounded and lower semi-continuous function satisfying the following dynamic programming equation with unknown  $\bar{J}$ :

$$(15) \quad \bar{J}(\mu) = \inf_{\bar{a} \in \bar{U}(\mu)} \left\{ \bar{f}(\mu, \bar{a}) + \gamma \mathbb{E} \left[ \bar{J}(\bar{F}(\mu, \bar{a}, \varepsilon^0)) \right] \right\}, \quad \mu \in \bar{S}.$$

Furthermore, there exists a pure stationary  $\bar{\pi}^* = (\bar{\pi}^*, \bar{\pi}^*, \dots) \in \bar{\Pi}^p$  that is optimal, i.e.,  $\bar{J}^{\bar{\pi}^*} = \bar{J}^*$ .

The Dynamic Programming Principle (15) is known to hold for universally measurable policies. See for example [10, Proposition 9.8].<sup>2</sup> The gist of the above theorem is to show that it also holds for Borel measurable policies.

*Proof. Step 1: Bellman operator and fixed point with universal measurability.* We define the Bellman operator  $\bar{T}$  by:

$$(16) \quad [\bar{T}\bar{J}](\mu) = \inf_{\bar{a} \in \bar{U}(\mu)} \left\{ \bar{f}(\mu, \bar{a}) + \gamma \mathbb{E} \left[ \bar{J}(\bar{F}(\mu, \bar{a}, \varepsilon^0)) \right] \right\}, \quad \mu \in \bar{S}.$$

Then, by [10, Proposition 9.8],  $\bar{T}$  is a strict contraction on the space of bounded universally measurable functions on  $\bar{S}$  and the fixed point coincides with the optimal value function for universally measurable policies, namely,  $\bar{J}^{*, Univ}$  defined as:

$$\bar{J}^{*, Univ}(\mu) = \inf_{\bar{\pi} \in \bar{\Pi}^{Univ}} \bar{J}^{\bar{\pi}}, \quad \mu \in \bar{S}$$

where  $\bar{\Pi}^{Univ}$  is the set of universally measurable mixed strategy functions, and for every  $\bar{\pi} \in \bar{\Pi}^{Univ}$ ,  $\bar{J}^{\bar{\pi}}$  is defined as in (14) for the Boreal measurable case.

**Step 2: Fixed point with Borel measurability.** We will apply the Banach fixed point theorem for  $\bar{T}$  on  $\mathcal{L}sc(\bar{S})$ , which denotes the set of real valued, bounded and lower semi-continuous functions on

<sup>2</sup>Beware that the MFMDP setting is denoted with overlines in the notation of the present paper, but due to the common noise it corresponds to a stochastic model in [10], which is denoted without overlines.

$\bar{S}$ . This set is a closed subset of the Banach space of real valued bounded functions on  $\bar{S}$  endowed with the sup norm.

The key point is to show that  $\bar{T}$  leaves  $\mathcal{L}sc(\bar{S})$  invariant. This follows by the measurable selection theorem for lower semi-continuous functions given in [10, Proposition 7.33], since we can show that the content of the curly bracket in (16) is lower semi-continuous if  $\bar{J} \in \mathcal{L}sc(\bar{S})$ . Indeed,  $\bar{f}$  is lower semi-continuous (see Lemma 42 in the Appendix). Moreover, since  $\bar{F}$  is continuous for  $\epsilon^0$  fixed (again by Lemma 42), the expectation is a continuous function of  $(\mu, \bar{a})$  whenever  $\bar{J}$  is continuous by the dominated convergence theorem. Now since a function is lower semi-continuous if and only if it is the pointwise limit of an increasing sequence of continuous functions, we can use the monotone convergence theorem to show that the expectation is the limit of an increasing sequence of continuous functions, hence that it is lower semi-continuous.

Last,  $\bar{T}$  is a strict contraction for the sup norm. We thus conclude by the Banach fixed point theorem that  $\bar{T}$  has a unique fixed point in  $\mathcal{L}sc(\bar{S})$ , which we denote by  $\bar{J}^{*,bd,lsc}$ . In other words,  $\bar{J}^{*,bd,lsc}$  is the unique bounded lower semi-continuous function satisfying the dynamic programming equation (15).

**Step 3:**  $\bar{J}^{*,bd,lsc} = \bar{J}^*$ . Being lower semi-continuous,  $\bar{J}^{*,bd,lsc}$  is universally measurable, so it is also a fixed point in the space of universally measurable functions. Hence, by uniqueness, it coincides with  $\bar{J}^{*,Univ}$ . Furthermore  $\bar{\Pi} \subseteq \bar{\Pi}^{Univ}$ . So we deduce:

$$(17) \quad \bar{J}^{*,bd,lsc}(\mu) = \bar{J}^{*,Univ}(\mu) \leq \bar{J}^*(\mu), \quad \mu \in \bar{S}.$$

Assumption **(H2)** implies that the lifted MFMDP satisfies the assumptions of [10, Corollary 9.17.2] so there exists a stationary pure (at the level-1) Borel measurable policy  $\bar{\pi}^* = (\bar{\pi}^*, \bar{\pi}^*, \dots)$  which is optimal in the sense that:

$$\bar{J}^{\bar{\pi}^*}(\mu) = \bar{J}^{*,Univ}(\mu), \quad \mu \in \bar{S}.$$

Consequently:

$$\bar{J}^*(\mu) \leq \bar{J}^{\bar{\pi}^*}(\mu) = \bar{J}^{*,Univ}(\mu) = \bar{J}^{*,bd,lsc}(\mu), \quad \mu \in \bar{S},$$

which together with (17) gives the equality  $\bar{J}^* = \bar{J}^{*,Univ}$ .

Hence  $\bar{J}^*$  satisfies the dynamic programming principle (15) and by Step 4, it is the unique bounded lower semi-continuous function satisfying it.  $\square$

**Remark 20.** When the mixed strategy function  $\bar{\pi}_n \in \bar{\Pi}^{Univ}$  is only universally measurable, for each time  $n$ , as in the above proof of Theorem 19, the understanding of condition (13) requires a modicum of care. Let  $q_n = \mathcal{L}(\mu_n) \in \mathcal{P}(\bar{S})$  be the distribution of the random state  $\mu_n$  with values in  $\bar{S}$ . We consider a Borel measurable kernel,  $\bar{\pi}_{q_n} : (\bar{S}, \mathcal{B}_{\bar{S}}) \rightarrow (\mathcal{P}(\bar{A}), \mathcal{B}_{\mathcal{P}(\bar{A})})$ , such that  $\bar{\pi}_{q_n}(\mu) = \bar{\pi}_n(\mu)$  for  $q_n$ -almost every  $\mu \in \bar{S}$  (see [10, Lemma 7.28 (c)] for existence). Then condition (13) says that  $\bar{\pi}_{q_n}$  is a regular version of the conditional probability of  $\bar{a}_n$  given  $\mu_n$ . Furthermore, the integration of a function  $\phi(\mu_n, \cdot)$  with respect to  $\bar{\pi}_n(\mu_n)$  should be understood in the following sense:

$$\int_{\bar{A}} \phi(\mu, \bar{a}) \bar{\pi}_n(\mu)(d\bar{a}) = \int_{\bar{A}} \phi(\mu, \bar{a}) \bar{\pi}_{q_n}(\mu)(d\bar{a}),$$

for  $q_n$ -almost every  $\mu \in \bar{S}$ , where  $q_n = \mathcal{L}(\mu_n)$ .

#### 4. Relations between the models

#### 4.1. Relations between MFC closed-loop policies and MFMDP policies.

In this section, we discuss some of the connections between the original level-0 MFC model and the lifted MFMDP model.

We start by highlighting that, intuitively, a closed-loop policy for the MFC can be viewed as sampled from a policy for the MFMDP by picking the common randomness. The following definition formalizes this idea.

**Definition 21.** *Let  $\pi \in \Pi^{CL}$  and  $\bar{\pi} \in \bar{\Pi}$ . We say that they correspond to each other if for each  $\mu \in \bar{S}$  and  $n \geq 0$ ,  $\bar{\pi}_n(\mu) \in \mathcal{P}(\bar{A})$  is equal to the push forward of  $\mathbb{P}_{\vartheta^0}$  by the map:*

$$\Theta^0 \ni \theta^0 \mapsto \mu(dx)\pi_n(x, \mu, \theta^0)(d\alpha) \in \bar{A}.$$

Note that if  $\pi$  and  $\bar{\pi}$  correspond to each other, then one is stationary if and only if the other one is. The main result of this section is the following.

**Theorem 22.** *Assume (H1) holds. Then for every  $\mu \in \mathcal{P}(S)$ ,  $\{J^\pi(\mu) : \pi \in \Pi^{CL}\} = \{\bar{J}^{\bar{\pi}}(\mu) : \bar{\pi} \in \bar{\Pi}\}$ . Similarly, for stationary policies, we have: for every  $\mu \in \mathcal{P}(S)$ ,  $\{J^\pi(\mu) : \pi \in \Pi^{CL} \text{ stationary}\} = \{\bar{J}^{\bar{\pi}}(\mu) : \bar{\pi} \in \bar{\Pi} \text{ stationary}\}$ .*

This result means that for every  $\pi \in \Pi^{CL}$ , there exists  $\bar{\pi} \in \bar{\Pi}$  such that  $J^\pi = \bar{J}^{\bar{\pi}}$  and conversely, for every  $\bar{\pi} \in \bar{\Pi}$ , there exists  $\pi \in \Pi^{CL}$  such that this equality holds, with the same result in the stationary case.

In preparation for the proof of this result, we give two useful technical lemmas whose proofs are deferred to the appendix. These lemmas describe the properties of conditional distributions of the level-0 state and action processes.

**Lemma 23.** *Assume (H1) holds. Let  $\alpha \in \mathbb{A}$ ,  $\mu_0 \in \mathcal{P}(S)$ , and let  $\mathbf{X}$  be the associated state process. Then:*

$$(18) \quad \mathbb{P}_{X_{n+1}}^0 = \bar{F}(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0, \varepsilon_{n+1}^0), \quad \mathbb{P} - a.s. \quad n \geq 0.$$

So  $\mathcal{L}(\mathbb{P}_{X_{n+1}}^0) = P(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0)$ ,  $n \geq 0$ , where the transition kernel  $P$  was defined in (9).

**Lemma 24.** *Assume (H1) holds. Let  $\alpha \in \mathbb{A}$ ,  $\mu_0 \in \mathcal{P}(S)$ , and let  $\mathbf{X}$  be the associated state process. For every  $n \geq 0$ , let  $\kappa_n : \bar{S} \rightarrow \mathcal{P}(\bar{A})$  be the Borel measurable disintegration kernel of  $\mathcal{L}(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0)$  along its first marginal. Then, if  $(\zeta, \bar{\eta})$  is an  $(\bar{S} \times \bar{A})$ -valued pair of stochastic processes which are  $\mathbb{F}^0$ -adapted, and satisfy:  $\zeta_0 = \mu_0$ ,  $\mathbb{P} - a.s.$ ,  $\zeta_{n+1} = \bar{F}(\zeta_n, \bar{\eta}_n, \varepsilon_{n+1}^0)$ ,  $\mathbb{P} - a.s.$ ,  $n \geq 0$ , and if  $\mathcal{L}(\bar{\eta}_n | \zeta_n) = \kappa_n(\zeta_n)$ ,  $\mathbb{P} - a.s.$   $n \geq 0$ , we have:*

$$(19) \quad \mathcal{L}(\zeta_n, \bar{\eta}_n) = \mathcal{L}(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0), \quad n \geq 0.$$

Intuitively,  $(\kappa_n)_{n \geq 0}$  plays the role of the conditional law of the action process. We will come back to this interpretation in the proof of Lemma 29. We are now ready to prove Theorem 22.

*Proof of Theorem 22. Step 1:* Let  $\pi \in \Pi^{CL}$ . Let  $\bar{\pi}$  corresponding to  $\pi$  in the sense of Definition 21. We now check the equality of the value functions  $J^\pi$  and  $\bar{J}^{\bar{\pi}}$ . Note that if  $\bar{\pi}$  is stationary, then so is  $\pi$ . Let  $(\mathbf{X}, \alpha)$  be a pair of state and action processes generated by  $(\pi, \mu_0)$ . Then

$$J^\pi(\mu_0) = \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n f(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}^0) \right] = \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n \bar{f}(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0) \right].$$

We now check that  $\mu_n = \mathbb{P}_{X_n}^0$  and  $\bar{a}_n = \mathbb{P}_{(X_n, \alpha_n)}^0$  form a pair of state and action processes generated by  $(\bar{\pi}, \mu_0)$ . This is indeed the case because (12) is implied by Lemma 23 and (13) is implied by the definition of  $\bar{\pi}_n$ . Consequently,

$$\bar{J}^{\bar{\pi}}(\mu_0) = \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n \bar{f}(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0) \right] = J^{\pi}(\mu_0).$$

**Step 2:** Conversely, let  $\bar{\pi} = (\bar{\pi}_n)_{n \geq 0}$  in  $\bar{\Pi}$ . For every  $n \geq 0$ ,  $\bar{\pi}_n : \bar{S} \rightarrow \mathcal{P}(\bar{A})$  is a Borel measurable map such that for every  $\mu \in \bar{S}$  we have  $\text{pr}_1(\bar{\pi}_n(\mu)) = \mu$ . According to the universal disintegration theorem [29, Corollary 1.26], there exists a Borel measurable probability kernel  $K : S \times \mathcal{P}(S \times A) \times \mathcal{P}(S) \rightarrow \mathcal{P}(A)$  such that for every  $\rho \in \mathcal{P}(S \times A)$  and  $\mu \in \mathcal{P}(S)$  such that  $\text{pr}_1(\rho) = \mu$ , we have  $\rho = \mu \hat{\otimes} K(\cdot, \rho, \mu)$ , where  $\hat{\otimes}$  denotes the product of a measure and a kernel. So for every integer  $n \geq 0$ ,  $x \in S$ ,  $\mu \in \bar{S}$  and  $\theta^0 \in \Theta^0$ , we define:

$$(20) \quad \pi_n(x, \mu, \theta^0) := K(x, \rho_{\bar{A}}(\bar{\pi}_n(\mu), h^0(\theta^0)), \mu),$$

where  $\rho_{\bar{A}}$  is the Blackwell-Dubins function of  $\bar{A}$ . Note that if  $\bar{\pi}$  is stationary, then so is  $\pi$ . Because the functions  $K$ ,  $h^0$  and  $\rho_{\bar{A}}$  are Borel measurable, so is the strategy function  $\pi_n$  for every  $n \geq 0$ . Hence  $\pi = (\pi_n)_{n \geq 0} \in \Pi^{CL}$ . Recall that the function  $h^0$  was introduced in Section 2.1.2, and that  $h^0(\vartheta^0)$  is uniformly distributed on  $[0, 1]$  by construction. Notice that for every  $\mu \in \bar{S}$  and for almost every  $\theta^0 \in [0, 1]$ , the definition of the universal disintegration kernel  $K$  implies that:

$$(21) \quad \left( \rho_{\bar{A}}(\bar{\pi}_n(\mu), \theta^0) \right)(dx, d\alpha) = \mu(dx) K(x, \rho_{\bar{A}}(\bar{\pi}_n(\mu), \theta^0), \mu)(d\alpha),$$

and as a result, we have:

$$(22) \quad \left( \rho_{\bar{A}}(\bar{\pi}_n(\mu), h^0(\theta^0)) \right)(dx, d\alpha) = \mu(dx) \pi_n(x, \mu, \theta^0)(d\alpha).$$

When  $\theta^0$  is replaced by  $\vartheta_n^0$ , by Blackwell-Dubins lemma (see first point in Lemma 2), the left hand side of (22) is a random variable with values in  $\bar{A} = \mathcal{P}(S \times A)$  with distribution  $\bar{\pi}_n(\mu)$ .

Next, we show that  $J^{\pi} = \bar{J}^{\bar{\pi}}$ . Let  $\mu_0 \in \bar{S}$ . Let  $(\zeta, \bar{\eta})$  be state and action processes generated by  $(\bar{\pi}, \mu_0)$  (see Definition 16). Let  $(\mathbf{X}, \alpha)$  be a pair of state and action processes generated by  $\pi$  and  $\mu_0$ . Using the fact that  $\mathbb{P}_{(X_n, \alpha_n)}^0 = \mathbb{P}_{X_n}^0 \hat{\otimes} \pi_n(\cdot, \mathbb{P}_{X_n}^0, \vartheta_n^0)$ , and the fact that  $\vartheta_n^0$  is independent of  $\mathbb{P}_{X_n}^0$  by equation (3), we have:

$$\begin{aligned} J^{\pi}(\mu_0) &= \sum_{n \geq 0} \gamma^n \mathbb{E} \left[ \int_{S \times A} f(x, \alpha, \mathbb{P}_{(X_n, \alpha_n)}^0) \mathbb{P}_{(X_n, \alpha_n)}^0(dx, d\alpha) \right] \\ &= \sum_{n \geq 0} \gamma^n \mathbb{E} \left[ \int_{S \times A} f(x, \alpha, \mathbb{P}_{X_n}^0 \hat{\otimes} \pi_n(\cdot, \mathbb{P}_{X_n}^0, \vartheta_n^0)) \mathbb{P}_{X_n}^0(dx) \pi_n(x, \mathbb{P}_{X_n}^0, \vartheta_n^0)(d\alpha) \right] \\ &= \sum_{n \geq 0} \gamma^n \mathbb{E} \left[ \int_{S \times A} f(x, \alpha, \bar{a}) \bar{a}(dx, d\alpha) \bar{\pi}_n(\mathbb{P}_{X_n}^0)(d\bar{a}) \right]. \end{aligned}$$

The last equality holds by the fact that both sides of (22) with  $\theta^0 = \vartheta_n^0$  are random variables with values in  $\bar{A} = \mathcal{P}(S \times A)$  with distribution  $\bar{\pi}_n(\mu)$ . On the other hand:

$$\begin{aligned} \bar{J}^{\bar{\pi}}(\mu_0) &= \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n \bar{f}(\zeta_n, \bar{\eta}_n) \right] \\ &= \sum_{n \geq 0} \gamma^n \mathbb{E} \left[ \int_{\bar{A}} \bar{f}(\zeta_n, \bar{a}) \mathcal{L}(\bar{\eta}_n | \zeta_n)(d\bar{a}) \right] \\ &= \sum_{n \geq 0} \gamma^n \mathbb{E} \left[ \int_{\bar{A}} \bar{f}(\mathbb{P}_{X_n}^0, \bar{a}) \bar{\pi}_n(\mathbb{P}_{X_n}^0)(d\bar{a}) \right], \end{aligned}$$

where the last equality holds by (13) because  $(\zeta, \bar{\eta})$  are generated by  $(\bar{\pi}, \mu_0)$ . This completes the proof.  $\square$

At this stage, we need to emphasize the crucial role played by the common randomization provided by the sequence  $(\vartheta_n^0)_{n \geq 0}$ . Its presence is what allowed us to prove that the value of a policy for the lifted MDP can always be achieved by a closed loop policy of the original MFC.

Comparing to [33], while they prove equality of the optimal value functions without the central randomization, the latter allows us to prove the identity of the value functions, policy by policy, even before taking the optimum values.

#### 4.2. Relations between MFC closed-loop and open-loop policies.

We first prove existence of optimal closed-loop Markov policies and then we prove equality of the open loop and closed loop value functions.

**Proposition 25.** *Assume (H1) and (H2) hold. There exists a stationary closed loop Markov policy for the original MFC that is optimal, i.e.,  $\pi^* = (\pi^*, \pi^*, \dots) \in \Pi^{CL}$  such that:  $J^{\pi^*} = J^{CL,*}$ .*

*Proof.* Let  $\bar{\pi}^*$  be an optimal pure stationary Markov policy for MFMDP whose existence is given in Theorem 19, and let  $\pi^* \in \Pi^{CL}$  be a closed-loop Markov policy whose value function is the same and whose existence is given in Theorem 22 (for the case of stationary policies). We have:

$$J^{\pi^*}(\mu) = \bar{J}^{\bar{\pi}^*}(\mu) = \inf_{\bar{\pi} \in \bar{\Pi}} \bar{J}^{\bar{\pi}}(\mu), \quad \mu \in \mathcal{P}(S),$$

Using Theorem 22 again, for every  $\pi \in \Pi^{CL}$  for MFC, there exists  $\bar{\pi} \in \bar{\Pi}$  for MFMDP such that  $\bar{J}^{\bar{\pi}} = J^\pi$ . So, for every  $\pi \in \Pi^{CL}$ ,

$$J^\pi(\mu) \geq \inf_{\bar{\pi} \in \bar{\Pi}} \bar{J}^{\bar{\pi}}(\mu) = J^{\pi^*}(\mu), \quad \mu \in \mathcal{P}(S),$$

which concludes the proof.  $\square$

**Remark 26.** *Notice that because  $\bar{\pi}^*$  is pure, since the Blackwell-Dubins function  $\rho_{\bar{A}}$  does not depend upon its second argument when the first is a point mass, we can conclude that  $\pi_n$  does not depend upon the common randomization as given by  $\vartheta_0^0$ . In other words, the above result would still hold even if we did not have the common randomization.*

We now show the equality of the open-loop and closed-loop optimal value functions of the MFC.

**Theorem 27.** *Assume (H1) holds. Then  $J^{OL,*} = J^{CL,*}$ .*

This result is a direct consequence of the following Lemma 28 and Lemma 29, together with Theorem 22. Indeed, by Lemma 28 and Lemma 29,

$$J^{CL,*} \geq J^{OL,*} \geq \bar{J}^*.$$

Moreover, by Theorem 22,

$$\bar{J}^* = J^{CL,*}.$$

Hence the above inequalities are equalities, which proves the first part of Theorem 27. The existence of an optimal closed-loop policy stems from Theorem 22, which entails the existence of an optimal open-loop policy by Lemma 28. Again, while this type of equality between the optimal value functions could be expected to hold under different assumptions and without the central randomization, we prove it here by leveraging the equalities proven in Lemma 28 and Lemma 29 policy by policy, before computing optima over sets of policies.

First, it is expected that every closed-loop policy for the MFC can be viewed as an open-loop policy for the MFC, which leads to the following result in terms of value functions.

**Lemma 28.** *Assume (H1) holds. For every  $\tilde{\pi} \in \Pi^{CL}$ , there exists  $\pi \in \Pi^{OL}$  such that:  $J^\pi = J^{\tilde{\pi}}$ .*

The proof is deferred to the Appendix C. Next, every open-loop policy for the MFC corresponds to a policy for the MFMDP, as we show in the following result.

**Lemma 29.** *Assume (H1) holds. For every  $\pi \in \Pi^{OL}$ , there exists  $\bar{\pi} \in \bar{\Pi}$  such that:  $\bar{J}^{\bar{\pi}} = J^\pi$ .*

*Proof.* Let  $\pi \in \Pi^{OL}$ . Let us fix an initial distribution  $\mu_0 \in \mathcal{P}(S)$ , let  $\alpha$  be an action process generated by  $\pi$ , and  $\mathbf{X}$  be the state process associated with  $(\alpha, \mu_0)$  (recall Definition 5). For each  $n \geq 0$ , we consider the probability kernel  $\kappa_n : \bar{S} \rightarrow \mathcal{P}(\bar{A})$  defined in the statement of Lemma 24. We construct by induction an  $(\bar{S} \times \bar{A})$ -valued pair of processes  $(\zeta, \bar{\eta})$  in the following way. For  $n = 0$  we set  $\zeta_0 = \mu_0$  and  $\bar{\eta}_0 = \rho_{\bar{A}}(\kappa_0(\zeta_0), h^0(\vartheta_0^0))$  where  $\rho_{\bar{A}}$  is the Blackwell-Dubins function of the space  $\bar{A}$  introduced in Lemma 2. Then for any  $n \geq 0$  we define: Obviously, each time we involve  $\vartheta_0^0$ , we rely on the central randomization. Still, the conclusion of this lemma should not be considered as obvious because it proves that we can pack all the dependence on the past carried by the open loop controls at level 0 into  $\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0$  and a probability measure on the space of actions at level 1.

$$\zeta_{n+1} = \bar{F}(\zeta_n, \bar{\eta}_n, \varepsilon_{n+1}^0), \quad \text{and} \quad \bar{\eta}_{n+1} = \rho_{\bar{A}}(\kappa_{n+1}(\zeta_{n+1}), h^0(\vartheta_{n+1}^0)).$$

The process  $\bar{\eta}$  is adapted to  $\mathbb{F}^0$  and by Lemma 2 it satisfies

$$\mathcal{L}(\bar{\eta}_n | \zeta_n) = \kappa_n(\zeta_n), \quad \mathbb{P} - a.s., \quad n \geq 0,$$

because  $\vartheta_{n+1}^0$  is independent of  $\zeta_n$ . Thus, by Lemma 24:

$$(23) \quad \mathcal{L}(\zeta_n, \bar{\eta}_n) = \mathcal{L}\left(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0\right), \quad n \geq 0.$$

Let  $\bar{\pi} = (\kappa_n)_{n \geq 0}$ . Since  $\kappa_n(\bar{U}(\mu)) = 1$  for every  $n \geq 0$ , we see that  $\bar{\pi} \in \bar{\Pi}$ . We conclude by noting that:

$$\bar{J}^{\bar{\pi}}(\mu_0) = \sum_{n=0}^{\infty} \gamma^n \mathbb{E}[\bar{f}(\zeta_n, \bar{\eta}_n)] = \sum_{n=0}^{\infty} \gamma^n \mathbb{E}\left[\bar{f}\left(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0\right)\right] = J^\pi(\mu_0),$$

where the second equality holds by (23). □

## 5. Mean-Field Q-Learning

### 5.1. State-action value function.

We now turn our attention to the question of *learning* the solution of the MFC problem in a model-free setting, i.e., assuming the model is unknown while still having access to sample realizations of state trajectories and associated rewards. Before considering algorithms, we first study the so-called state-action value function.

In order to take advantage of the strongest results proven so far, we now assume both **(H1)** and **(H2)** hold. Under these assumptions, recall that the DPP given by Theorem 19 holds. In this section, we restrict ourselves to non-randomized stationary policies. When  $\bar{\pi} = (\bar{\pi}, \bar{\pi}, \dots)$ , we use freely the notation  $\bar{J}^{\bar{\pi}} := \bar{J}^{\bar{\pi}}$ . Note however that we will use common randomization in the numerical section for the purpose of exploration.

In this section, without any loss of generality, we restrict the search for optimal policies to the set:

$$\bar{U}_B(\bar{A}|\bar{S}) = \{\bar{\pi} \in \bar{\Pi}^p \mid \forall \mu \in \bar{S}, \bar{\pi}(\mu) \in \bar{U}(\mu)\}.$$

For each  $\bar{\pi} \in \bar{U}_B(\bar{A}|\bar{S})$ , the mapping  $\bar{S} \ni \mu \mapsto \delta_{\bar{\pi}(\mu)} \in \mathcal{P}(\bar{A})$  which assigns to each  $\mu \in \bar{S}$  the Dirac point mass at the point  $\bar{\pi}(\mu) \in \bar{A}$  is a Borel measurable function by definition of the Borel  $\sigma$ -field of  $\mathcal{P}(\bar{A})$ .

Now, for each  $\bar{\pi} \in \bar{\Pi}$ , we introduce the state-action value function  $\bar{Q}^{\bar{\pi}} : \bar{\Gamma} \rightarrow \mathbb{R}$  defined by:

$$(24) \quad \bar{Q}^{\bar{\pi}}(\mu, \bar{a}) := \bar{f}(\mu, \bar{a}) + \sum_{n \geq 1} \gamma^n \mathbb{E}[\bar{f}(\mu_n, \bar{\pi}(\mu_n))], \quad (\mu, \bar{a}) \in \bar{\Gamma},$$

where the process  $(\mu_n)_{n \geq 0}$  starting at  $\mu_0 = \mu$  satisfies  $\mu_1 = \bar{F}(\mu, \bar{a}, \varepsilon_1^0)$ , and for every  $n \geq 1$ :  $\mu_{n+1} = \bar{F}(\mu_n, \bar{\pi}(\mu_n), \varepsilon_{n+1}^0)$ . Next we define the optimal state-action value function by:

$$(25) \quad \bar{Q}^*(\mu, \bar{a}) := \inf_{\bar{\pi} \in \bar{U}_B(\bar{A}|\bar{S})} \bar{Q}^{\bar{\pi}}(\mu, \bar{a}), \quad (\mu, \bar{a}) \in \bar{\Gamma}.$$

The main goal of this section is to prove the following dynamic programming principle for  $\bar{Q}^*$ .

**Theorem 30.** *Assume **(H1)** and **(H2)** hold. The optimal state-action value function  $\bar{Q}^*$  satisfies the so-called **Bellman equation for state-action value function**:*

$$(26) \quad \bar{Q}^*(\mu, \bar{a}) = \bar{f}(\mu, \bar{a}) + \gamma \mathbb{E} \left[ \inf_{\bar{a}' \in \bar{U}(\bar{F}(\mu, \bar{a}, \varepsilon^0))} \bar{Q}^*(\bar{F}(\mu, \bar{a}, \varepsilon^0), \bar{a}') \right], \quad (\mu, \bar{a}) \in \bar{\Gamma}.$$

We will prove this result by showing that  $\bar{Q}^*$  is the unique fixed point of state-action Bellman operator  $T$  defined on the set  $\mathcal{Lsc}(\bar{\Gamma})$  of bounded lower semi-continuous functions on  $\bar{\Gamma}$ , by:

$$(27) \quad [T\bar{Q}](\mu, \bar{a}) := \bar{f}(\mu, \bar{a}) + \gamma \mathbb{E} \left[ \inf_{\bar{a}' \in \bar{U}(\bar{F}(\mu, \bar{a}, \varepsilon^0))} \bar{Q}(\bar{F}(\mu, \bar{a}, \varepsilon^0), \bar{a}') \right], \quad (\mu, \bar{a}) \in \bar{\Gamma}.$$

We first justify in Lemma 31 the fact that the operator  $T$  is well-defined.

Since  $\bar{\Gamma}$  is a closed subset of the Polish space  $\bar{S} \times \bar{A}$ , it is a Borel space, and the space  $\mathcal{BD}_u(\bar{\Gamma})$  of bounded real-valued universally measurable functions on  $\bar{\Gamma}$  endowed with the sup norm  $\|f\|_\infty = \sup_{(\mu, \bar{a}) \in \bar{\Gamma}} |f(\mu, \bar{a})|$  is a Banach space. While the set  $\mathcal{Lsc}(\bar{\Gamma})$  is not a vector space, it is a closed subset of  $\mathcal{BD}_u(\bar{\Gamma})$ , hence a complete metric space for the metric  $d_\infty(f, f') = \|f - f'\|_\infty$ .

**Lemma 31.** *Assume **(H1)** and **(H2)** hold. The set  $\mathcal{Lsc}(\bar{\Gamma})$  is invariant under the state-action Bellman operator  $T$ , which is a strict contraction on this metric space.*

*Proof.* We first claim that the set  $\mathcal{Lsc}(\bar{\Gamma})$  is invariant under  $T$ . We need to show that  $T\bar{Q}$  is lower semi-continuous whenever  $\bar{Q}$  is. To wit, by the projection property for infima of lower semi-continuous functions (see for example [10, Proposition 7.33]), the function  $\bar{S} \ni \mu' \mapsto \inf_{\bar{a}' \in \bar{U}(\mu')} \bar{Q}(\mu', \bar{a}')$  is lower semi-continuous. Since  $\bar{\Gamma} \ni (\mu, \bar{a}) \mapsto \mu' = \bar{F}(\mu, \bar{a}, e^0) \in \bar{S}$  is continuous for  $e^0 \in E^0$  fixed, the infimum in formula (27) is then a lower semi-continuous function of  $(\mu, \bar{a})$  for fixed  $e^0 \in E^0$ . Finally, Fatou's theorem implies that the expectation in (27) is a lower semi-continuous function of  $(\mu, \bar{a}) \in \bar{\Gamma}$ . Furthermore  $\bar{f}$  is continuous. Consequently,  $T\bar{Q}$  is also lower semi-continuous. Now if  $\bar{Q}_1$  and  $\bar{Q}_2$  are elements of  $\mathcal{Lsc}(\bar{\Gamma})$  we have:

$$\begin{aligned} \|T\bar{Q}_1 - T\bar{Q}_2\|_\infty &\leq \gamma \mathbb{E} \left[ \sup_{(\mu, \bar{a}) \in \bar{\Gamma}} \left| \inf_{\bar{a}' \in \bar{U}(\bar{F}(\mu, \bar{a}, \varepsilon^0))} \bar{Q}_1(\bar{F}(\mu, \bar{a}, \varepsilon^0), \bar{a}') - \inf_{\bar{a}' \in \bar{U}(\bar{F}(\mu, \bar{a}, \varepsilon^0))} \bar{Q}_2(\bar{F}(\mu, \bar{a}, \varepsilon^0), \bar{a}') \right| \right] \\ &\leq \gamma \mathbb{E} \left[ \sup_{(\mu, \bar{a}) \in \bar{\Gamma}} \sup_{\bar{a}' \in \bar{U}(\bar{F}(\mu, \bar{a}, \varepsilon^0))} \left| \bar{Q}_1(\bar{F}(\mu, \bar{a}, \varepsilon^0), \bar{a}') - \bar{Q}_2(\bar{F}(\mu, \bar{a}, \varepsilon^0), \bar{a}') \right| \right] \\ &\leq \gamma \|\bar{Q}_1 - \bar{Q}_2\|_\infty. \end{aligned}$$

Since  $\gamma < 1$ , this proves that  $T$  is a strict contraction on  $\mathcal{Lsc}(\bar{\Gamma})$ . We conclude the proof of the result using the Banach fixed point theorem.  $\square$

Using the Markov property, we can rewrite the state-action value function  $\bar{Q}^{\bar{\pi}}$  in terms of the state value function  $\bar{J}^{\bar{\pi}}$ :

$$(28) \quad \bar{J}^{\bar{\pi}}(\mu) = \bar{Q}^{\bar{\pi}}(\mu, \bar{\pi}(\mu)), \quad \bar{\pi} \in \bar{\Pi}^p, \mu \in \bar{S}.$$

Now, for the optimal value functions, we have the following.

**Lemma 32.** *Assume (H1) and (H2) hold. For all  $(\mu, \bar{a}) \in \bar{\Gamma}$ :  $\bar{Q}^*(\mu, \bar{a}) = \bar{f}(\mu, \bar{a}) + \gamma \mathbb{E}[\bar{J}^*(\bar{F}(\mu, \bar{a}, \varepsilon^0))]$ .*

*Proof.* We show the inequalities in both directions. First, for  $(\mu, \bar{a}) \in \bar{\Gamma}$  given, let  $q = \mathcal{L}(\mu_1) \in \mathcal{P}(\bar{S})$  be the distribution of the random measure  $\mu_1 = \bar{F}(\mu, \bar{a}, \varepsilon^0)$ . By [10, Corollary 9.5.2] and the fact that  $\{(\varphi, \varphi, \dots) \text{ with } \varphi \in \bar{U}_B(\bar{A}|\bar{S})\} \subseteq \bar{\Pi}$ , we have:

$$\int_{\bar{S}} \bar{J}^*(\mu) q(d\mu) = \inf_{\bar{\pi} \in \bar{\Pi}} \int_{\bar{S}} \bar{J}^{\bar{\pi}}(\mu) q(d\mu) \leq \inf_{\substack{\bar{\pi} = (\varphi, \varphi, \dots); \\ \varphi \in \bar{U}_B(\bar{A}|\bar{S})}} \int_{\bar{S}} \bar{J}^{\bar{\pi}}(\mu) q(d\mu).$$

Hence  $f(\mu, \bar{a}) + \gamma \mathbb{E}[\bar{J}^*(\mu_1)] \leq \bar{Q}^*(\mu, \bar{a})$ .

Conversely, under the standing assumptions, by Theorem 19 there exists  $\bar{\pi}^* \in \bar{U}_B(\bar{A}|\bar{S})$  that is optimal. So we have:

$$\bar{Q}^*(\mu, \bar{a}) = \inf_{\bar{\pi} \in \bar{U}_B(\bar{A}|\bar{S})} \{f(\mu, \bar{a}) + \gamma \mathbb{E}[\bar{J}^{\bar{\pi}}(\mu_1)]\} \leq f(\mu, \bar{a}) + \gamma \mathbb{E}[\bar{J}^{\bar{\pi}^*}(\mu_1)] = f(\mu, \bar{a}) + \gamma \mathbb{E}[\bar{J}^*(\mu_1)],$$

for every  $(\mu, \bar{a}) \in \bar{\Gamma}$ , which concludes the proof.  $\square$

**Lemma 33.** *Assume (H1) and (H2) hold.  $\bar{Q}^*$  is lower semi-continuous and, as a result, there exists  $\bar{\pi} \in \bar{U}_B(\bar{A}|\bar{S})$  such that for every  $\mu \in \bar{S}$ ,  $\bar{\pi}(\mu) \in \arg \inf_{\bar{a} \in \bar{U}(\mu)} \bar{Q}^*(\mu, \bar{a})$ .*

*Proof.* For each fixed  $e^0 \in E^0$ ,  $(\mu, \bar{a}) \mapsto \bar{J}^*(\bar{F}(\mu, \bar{a}, e^0))$  is lower semi-continuous by lower semi-continuity of  $\bar{J}^*$  (see Theorem 19) and continuity of  $\bar{F}(\cdot, \cdot, e^0)$ . As in the proof of Lemma 31, Fatou's theorem implies that the function  $(\mu, \bar{a}) \mapsto \mathbb{E}[\bar{J}^*(\bar{F}(\mu, \bar{a}, e^0))]$  is also lower semi-continuous. Furthermore,  $\bar{f}$  is continuous. So, by the expression in Lemma 32,  $\bar{Q}^*$  is a lower semi-continuous function on  $\bar{\Gamma}$ .

Since  $\bar{\Gamma}$  is a closed subset of  $\bar{S} \times \bar{A}$  and  $\bar{A}$  is compact, by applying a selection theorem for lower semi-continuous function [10, Proposition 7.33] on  $\bar{Q}^* : \bar{\Gamma} \rightarrow \mathbb{R}$ , we obtain that there exists a Borel measurable function  $\tilde{\pi} \in \bar{U}_B(\bar{A}|\bar{S})$  whose graph is contained in  $\bar{\Gamma}$  and  $\bar{Q}^*(\mu, \tilde{\pi}(\mu)) = \inf_{\bar{a} \in \bar{U}(\mu)} \bar{Q}^*(\mu, \bar{a})$ , for all  $\mu \in \bar{S}$ .  $\square$

**Lemma 34.** *Assume (H1) and (H2) hold. For all  $\mu \in \bar{S}$ ,  $\inf_{\bar{a} \in \bar{U}(\mu)} \bar{Q}^*(\mu, \bar{a}) = \bar{J}^*(\mu)$ .*

*Proof.* We first show the inequality  $\inf_{\bar{a} \in \bar{U}(\mu)} \bar{Q}^*(\mu, \bar{a}) \geq \bar{J}^*(\mu)$ . Let us denote by  $\tilde{\pi} \in \bar{U}_B(\bar{A}|\bar{S})$  the strategy function in Lemma 33. By definition,

$$\inf_{\bar{a} \in \bar{U}(\mu)} \bar{Q}^*(\mu, \bar{a}) = \bar{Q}^*(\mu, \tilde{\pi}(\mu)) = \inf_{\bar{\pi} \in \bar{U}_B(\bar{A}|\bar{S})} \bar{Q}^{\bar{\pi}}(\mu, \tilde{\pi}(\mu)), \quad \mu \in \bar{S}.$$

Then for each  $\bar{\pi} \in \bar{U}_B(\bar{A}|\bar{S})$ , we denote  $\bar{\pi}^{\bar{\pi}} = (\bar{\pi}, \bar{\pi}, \bar{\pi}, \dots) \in \bar{\Pi}$ . So,

$$\bar{Q}^{\bar{\pi}}(\mu, \tilde{\pi}(\mu)) = f(\mu, \tilde{\pi}(\mu)) + \mathbb{E} \left[ \sum_{n=1}^{\infty} \gamma^n \bar{f}(\mu_n, \bar{\pi}(\mu_n)) \right] = \bar{J}^{\bar{\pi}^{\bar{\pi}}}(\mu) \geq \bar{J}^*(\mu), \quad \mu \in \bar{S},$$

which provides the first inequality. To prove the converse inequality, let  $\bar{\pi}^* \in \bar{U}_B(\bar{A}|\bar{S})$  be an optimal non-randomized stationary Markov policy whose existence is given by Theorem 19, and notice that for every  $\mu \in \bar{S}$ ,

$$\bar{J}^*(\mu) = \bar{J}^{\bar{\pi}^*}(\mu) = \bar{Q}^{\bar{\pi}^*}(\mu, \bar{\pi}^*(\mu)) \geq \bar{Q}^*(\mu, \bar{\pi}^*(\mu)) \geq \inf_{\bar{a} \in \bar{U}(\mu)} \bar{Q}^*(\mu, \bar{a}),$$

by equation 28, Lemma 32. This concludes the proof.  $\square$

We can now complete the proof of Theorem 30.

*Proof of Theorem 30.* The Bellman equation (26) is a direct consequence of Lemma 32 and Lemma 34. Since  $T$  is a strict contraction mapping on  $\mathcal{Lsc}(\bar{\Gamma})$  by Lemma 31 and since  $\mathcal{Lsc}(\bar{\Gamma})$  is closed in the Banach space  $\mathcal{BD}_u(\bar{\Gamma})$ , by the Banach fixed point theorem we conclude that  $\bar{Q}^*$  is the unique fixed point of  $T$  on  $\mathcal{Lsc}(\bar{\Gamma})$ .  $\square$

Next, we build upon the previous results to propose reinforcement learning algorithms for the original MFC problem. From now on, we assume that the state and action spaces are finite, unless otherwise specified.

## 5.2. Controls for finite state and action spaces.

In this rest of this section, we assume that  $S$  and  $A$  are finite, we denote their numbers of elements by  $|S|$  and  $|A|$  respectively, and we denote by  $x^{(1)}, \dots, x^{(|S|)}$  and  $\alpha^{(1)}, \dots, \alpha^{(|A|)}$  their elements. We first revisit the description of the action space and then propose two reinforcement learning methods in this setting. We shall explain later how to adapt reinforcement learning techniques to the case of continuous spaces.

Before introducing the mean-field Q-learning algorithm, we first provide a representation of the set  $\bar{\Gamma} \subseteq \bar{S} \times \bar{A} = \mathcal{P}(S) \times \mathcal{P}(S \times A)$  on which the  $\bar{Q}^*$  function is defined.

Since we assume that  $S$  finite, its lifted space  $\mathcal{P}(S)$  can be identified with a simplex  $\mathfrak{S}$  in  $\mathbb{R}^{|S|}$ . In other words, we treat a distribution  $\mu \in \mathcal{P}(S)$  as an  $|S|$ -dimensional vector  $(\mu^{(i)})_{i=1, \dots, |S|}$  whose non-negative coordinates sum up to one. Similarly, since  $A$  is finite, we identify  $\mathcal{P}(A)$  to a simplex  $\mathfrak{A}$  in  $\mathbb{R}^{|A|}$ . However, representing admissible actions  $\bar{a} \in \bar{U}(\mu) \subseteq \mathcal{P}(S \times A)$  of the lifted MDP requires a modicum of care due to the constraint. A first approach is to identify  $\mathcal{P}(S \times A)$  with a simplex in  $\mathbb{R}^{|S| \times |A|}$  and to view a lifted action  $\bar{a}$  as a  $|S| \times |A|$  matrix  $(\bar{a}(x^{(i)}, \alpha^{(j)}))_{1 \leq i \leq |S|, 1 \leq j \leq |A|}$  of non-negative numbers

summing up to 1. Then a pair  $(\mu, \bar{a}) \in \bar{S} \times \bar{A}$  is in  $\bar{\Gamma}$  if and only if the following linear constraint is satisfied:  $\sum_{j=1}^{|A|} \bar{a}(\mu^{(i)}, \alpha^{(j)}) = \mu^{(i)}$  for all  $i = 1, \dots, |S|$ . The above transformation is straightforward but not sufficient for our purposes because it provides only a representation of the actions and controls of the central planner, and it does not address the strategy functions of non-randomized stationary mixed Markovian closed-loop policies for an individual agent in our original optimization problem.

For any pair  $(\mu, \bar{a}) \in \bar{\Gamma}$ , we can define the mapping  $k_\mu : S \rightarrow \mathcal{P}(A)$ : for  $i = 1, \dots, |S|$ ,

$$k_\mu(x^{(i)}) = \bar{a}(\mu^{(i)}, \alpha^{(j)})/\mu(x^{(i)}), \quad \text{if } \mu(x^{(i)}) > 0,$$

and any value otherwise. Note that here, there is no common randomization. As proved above (see Theorem 19), there exists a non-randomized stationary policy for the lifted MDP. So the central planner can look for strategy functions within the set:

$$(29) \quad \mathcal{A} := \{\tilde{a} : S \rightarrow \mathcal{P}(A) \mid \tilde{a} \text{ Borel measurable}\}.$$

Consider the function  $\tilde{Q}^* : \mathcal{P}(S) \times \mathcal{A} \rightarrow \mathbb{R}$  defined by:

$$(30) \quad \tilde{Q}^*(\mu, \tilde{a}) := \bar{Q}^*(\mu, \mu \hat{\otimes} \tilde{a}).$$

Then the Bellman equation (26) becomes:

$$(31) \quad \tilde{Q}^*(\mu, \tilde{a}) = \int_{S \times A} f(x, \alpha, \mu \hat{\otimes} \tilde{a}) \tilde{a}(x, d\alpha) \mu(dx) + \gamma \mathbb{E} \left[ \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\mu_1, \tilde{a}') \right], \quad (\mu, \tilde{a}) \in \mathcal{P}(S) \times \mathcal{A},$$

where  $\mu_1 = \bar{F}(\mu, \mu \hat{\otimes} \tilde{a}, \varepsilon^0)$ , keeping in mind that the integral over  $S \times A$  is in fact a finite sum. Even though  $S$  and  $A$  are finite, equation (31) still needs to be understood as a fixed point in the space of bounded lower semi-continuous functions on a closed subset of a finite dimensional Euclidean space, as the measurability issues addressed in deriving equation (26) still remain. We also introduce the function  $\tilde{f} : \mathcal{P}(S) \times \mathcal{A} \rightarrow \mathbb{R}$  such that:

$$\tilde{f}(\mu, \tilde{a}) := \bar{f}(\mu, \mu \hat{\otimes} \tilde{a}) = \int_{S \times A} f(x, \alpha, \mu \hat{\otimes} \tilde{a}) \tilde{a}(x, d\alpha) \mu(dx), \quad (\mu, \tilde{a}) \in \mathcal{P}(S) \times \mathcal{A}.$$

In the rest of this section, we propose two model-free algorithms relying on the optimal state-action value function  $\bar{Q}^* : \bar{\Gamma} \rightarrow \mathbb{R}$  or equivalently  $\tilde{Q}^* : \mathcal{P}(S) \times \mathcal{A} \rightarrow \mathbb{R}$ .

### 5.3. Simplex discretization and tabular MFQ-learning.

We consider two settings, depending on whether the controls at level-0 are mixed or pure. In both cases, we prove convergence of a tabular Q-learning algorithm, after suitable discretization of the simplexes. When using pure controls, we can prove not only convergence of the value function but also of the optimizer.

**5.3.1. Q-learning with controls that are mixed at level-0.** Since the simplexes  $\mathfrak{S}$  and  $\mathfrak{A}$  are not finite, it is not possible to directly apply a tabular version of Q-learning algorithm to approximate  $\bar{Q}^*$ . A possible workaround is to first replace these simplexes by finite subsets  $\check{\mathfrak{S}} \subset \mathfrak{S}$  and  $\check{\mathfrak{A}} \subset \mathfrak{A}$ . Let  $\check{\mathcal{A}} = \{\check{a} : S \rightarrow \check{\mathfrak{A}}\}$ . In particular,  $|\check{\mathcal{A}}| = |\check{\mathfrak{A}}|^{|S|}$  because we identify functions in  $\check{\mathcal{A}}$  with  $|S|$ -dimensional vectors whose entries take values in the finite set  $\check{\mathfrak{A}}$ . To ensure that the mean-field term takes values in the finite set  $\check{\mathfrak{S}}$ , we use a projection: at time  $n$ , given  $\mu_n \in \check{\mathfrak{S}}$ , we compute  $\mu_{n+1} = \bar{F}(\mu_n, \mu_n \hat{\otimes} \check{a}, \varepsilon_{n+1}^0)$ , and then we project  $\mu_{n+1}$  back on  $\check{\mathfrak{S}}$  using a projection operator  $\text{Proj}_{\check{\mathfrak{S}}} : \mathcal{P}(S) \rightarrow \check{\mathfrak{S}}$ . Precise definitions of the discretization and the projection are provided below, after introducing a discrete version of the original MFC problem.

More precisely, we consider the **projected MFC problem**:

$$\inf_{\tilde{\pi} \in \tilde{\Pi}} \check{J}_{\tilde{\pi}}(\mu_0), \quad \mu_0 \in \check{\mathfrak{S}},$$

where  $\tilde{\Pi} = \{\tilde{\pi} : S \times \check{\mathfrak{S}} \rightarrow \check{\mathfrak{A}}\}$ , and for every strategy function  $\tilde{\pi} : S \times \check{\mathfrak{S}} \rightarrow \check{\mathfrak{A}}$ ,  $\check{J}^{\tilde{\pi}} : \check{\mathfrak{S}} \rightarrow \mathbb{R}$  is defined by:

$$(32) \quad \check{J}^{\tilde{\pi}}(\mu_0) = \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n \tilde{f} \left( \mu_n^{\mu_0, \tilde{\pi}}, \tilde{\pi}(\cdot, \mu_n^{\mu_0, \tilde{\pi}}) \right) \right]$$

where

$$(33) \quad \mu_{n+1}^{\mu_0, \tilde{\pi}} = \text{Proj}_{\check{\mathfrak{S}}} \circ \bar{F} \left( \mu_n^{\mu_0, \tilde{\pi}}, \mu_n^{\mu_0, \tilde{\pi}} \hat{\otimes} \tilde{\pi}(\cdot, \mu_n^{\mu_0, \tilde{\pi}}), \varepsilon_{n+1}^0 \right) =: \check{\Phi}^{\tilde{\pi}, \varepsilon_{n+1}^0}(\mu_n^{\mu_0, \tilde{\pi}}).$$

We will denote by  $\check{J}^*$  and  $\check{Q}^*$  respectively the optimal state and state-action value functions of this projected MFC problem. Here  $\check{Q} : \check{\mathfrak{S}} \times \check{\mathcal{A}} \rightarrow \mathbb{R}$  can be represented by a matrix (also called a table) in  $\mathbb{R}^{|\check{\mathfrak{S}}| \times |\check{\mathcal{A}}|}$  and is viewed as an approximation of  $\tilde{Q}^* : \mathcal{P}(S) \times \mathcal{A} \rightarrow \mathbb{R}$  of the original MFC problem.

This problem can be viewed as an MDP with finite state and action spaces. In this case, a straightforward adaptation of the tabular Q-learning algorithm leads to Algorithm 1. Note that, even in the absence of common noise, this algorithm is possibly stochastic since at each episode, the order in which the state-action pairs are picked is potentially random. In practice, the order could be fixed in advance or stem from a sampled trajectory.

---

**Algorithm 1:** Mean-Field Q-learning (MFQ) with simplex discretization

---

**Data:** A number of episodes  $N_{\text{epi}}$ ; a sequence of learning rates  $(\eta_n)_{n=0, \dots, N_{\text{epi}}-1}$ ; a sequence of state-action pairs  $(\check{\mu}_n, \check{a}_n)_{n \geq 0} \in \check{\mathfrak{S}} \times \check{\mathcal{A}}$ .

**Result:**  $\check{Q}_{N_{\text{epi}}}$ , an approximation of  $\check{Q}^*$  on  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$ .

1 **begin**

2     Initialize table  $\check{Q}_0 \in \mathbb{R}^{|\check{\mathfrak{S}}| \times |\check{\mathcal{A}}|}$ ,  $\mu_0 \in \check{\mathfrak{S}}$  and  $a_0 \in \mathcal{A}$

3     **for**  $n = 0, 1, \dots, N_{\text{epi}} - 1$  **do**

4         Execute action  $\check{a}_n$ , observe  $\check{\mu}'_{n+1} = \text{Proj}_{\check{\mathfrak{S}}} \circ \bar{F}(\check{\mu}_n, \check{\mu}_n \hat{\otimes} \check{a}_n, \varepsilon_{n+1}^0)$  and cost  $\tilde{f}(\check{\mu}_n, \check{a}_n)$

5         Initialize  $\check{Q}_{n+1} = \check{Q}_n$  on  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$

6         Set  $\check{Q}_{n+1}(\check{\mu}_n, \check{a}_n) = (1 - \eta_n)\check{Q}_n(\check{\mu}_n, \check{a}_n) + \eta_n \left( \tilde{f}(\check{\mu}_n, \check{a}_n) + \gamma \min_{\check{a}' \in \check{\mathcal{A}}} \check{Q}_n(\check{\mu}'_{n+1}, \check{a}') \right)$

7     **return**  $\check{Q}_{N_{\text{epi}}}$

---

Algorithm 1 returns the table  $\check{Q}_{N_{\text{epi}}}$  after  $N_{\text{epi}}$  episodes. We prove below that this table converges to the optimal Q-function  $\check{Q}^*$  in a suitable sense. To keep the paper at a reasonable length, we will make the following simplifying assumptions.

We endow the simplexes  $\check{\mathfrak{S}}$  and  $\check{\mathfrak{A}}$  respectively with the Euclidean distances  $d_{\check{\mathfrak{S}}}$  and  $d_{\check{\mathfrak{A}}}$  of the spaces  $\mathbb{R}^{|\check{\mathfrak{S}}|}$  and  $\mathbb{R}^{|\check{\mathfrak{A}}|}$ . Because  $S$  is finite, we can identify  $\mathcal{A}$  defined in (29) with  $\mathcal{P}(A)^{|S|}$  and endow it with the distance  $d_{\mathcal{A}}(\tilde{a}, \tilde{a}') = \sup_{x \in S} d_{\mathfrak{A}}(\tilde{a}(x), \tilde{a}'(x))$  for  $\tilde{a}, \tilde{a}' \in \mathcal{A}$ . Furthermore, we consider the following discretizations of the simplexes. Let  $\varepsilon_{\check{\mathfrak{S}}} > 0$  satisfying: for all  $\mu \in \check{\mathfrak{S}}$ , there exists  $\check{\mu} \in \check{\mathfrak{S}}$  s.t.  $d_{\check{\mathfrak{S}}}(\mu, \check{\mu}) \leq \varepsilon_{\check{\mathfrak{S}}}$ . Similarly, let  $\varepsilon_{\check{\mathfrak{A}}} > 0$  satisfying: for all  $\nu \in \check{\mathfrak{A}}$ , there exists  $\check{\nu} \in \check{\mathfrak{A}}$  such that  $d_{\check{\mathfrak{A}}}(\nu, \check{\nu}) \leq \varepsilon_{\check{\mathfrak{A}}}$ . Because  $S$  is finite and the definition of the distance  $d_{\mathcal{A}}$ , we have for every  $\tilde{a} \in \mathcal{A}$ , there exists  $\check{a} \in \check{\mathcal{A}}$ , s.t.  $d_{\mathcal{A}}(\tilde{a}, \check{a}) \leq \varepsilon_{\check{\mathcal{A}}}$ .

**Assumption (H3). Regularity of the data:**  $\tilde{f}$  is bounded and Lipschitz continuous with respect to  $(\mu, \tilde{a})$  with constant  $L_{\tilde{f}}$ , namely for every  $(\mu, \tilde{a}), (\mu', \tilde{a}') \in \mathfrak{S} \times \mathcal{A}$ , we have

$$|\tilde{f}(\mu, \tilde{a}) - \tilde{f}(\mu', \tilde{a}')| \leq L_{\tilde{f}} (\|\mu - \mu'\|_{d_{\mathfrak{S}}} + d_{\mathcal{A}}(\tilde{a}, \tilde{a}')) \quad \text{and} \quad \tilde{f}(\mu, \tilde{a}) \leq L_{\tilde{f}}.$$

Also,  $\bar{F}$  is Lipschitz continuous with respect to  $\mu$  and  $\tilde{a}$  with constant  $L_{\bar{F}}$  in expectation over the randomness of the common noise, namely: for every  $(\mu, \tilde{a}), (\mu', \tilde{a}') \in \mathfrak{S} \times \mathcal{A}$ ,

$$\mathbb{E}_{\varepsilon^0} [\|\bar{F}(\mu, \mu \hat{\otimes} \tilde{a}, \varepsilon^0) - \bar{F}(\mu', \mu' \hat{\otimes} \tilde{a}', \varepsilon^0)\|_{d_{\mathfrak{S}}}] \leq L_{\bar{F}} (\|\mu - \mu'\|_{d_{\mathfrak{S}}} + d_{\mathcal{A}}(\tilde{a}, \tilde{a}'))$$

**Assumption (H4). Regularity of the value function:**  $\bar{J}^*$  is Lipschitz continuous w.r.t.  $\mu$  with constant  $L_{\bar{J}^*}$ .

**Assumption (H5). Covering time:** There exists a finite  $T_{cov}$  such that with probability 1/2 (over the randomness of the common noise and of Algorithm 1) the following holds: For every starting point in  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$ , every element of  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$  has been visited before time  $T_{cov}$  during the execution of Algorithm 1.

The regularity of  $\bar{J}^*$  in (H4) can typically be ensured through suitable conditions on the data of the problem, as e.g. in [17, 11, 13]. Assumption (H5) is similar to the covering time assumption in [19]. In practice, exploration can be enhanced by adjusting the greediness level and by using exploring starts (if the learner can query an oracle which simulates transitions from any  $(\mu, \tilde{a})$ ). Note that the boundedness of the one-stage cost  $\tilde{f}$  from Assumption (H3) together with the fact that  $\gamma \in (0, 1)$  ensures the existence of a finite bound  $\check{J}_{bound}$  for the state value function of the projected MFC problem. We denote by  $\beta = (1 - \gamma)/2$  the horizon of the MDP corresponding to the projected MFC problem, and for  $\delta \in (0, 1)$ , we let  $T_{cov}(\delta) = \lceil T_{cov} \log_2(1/(2\delta)) \rceil$ . We consider projection operators  $\text{Proj}_{\check{\mathfrak{S}}} : \mathfrak{S} \rightarrow \check{\mathfrak{S}}$  and  $\text{Proj}_{\check{\mathcal{A}}} : \mathcal{A} \rightarrow \check{\mathcal{A}}$  such that  $(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) := (\check{\mu}, \check{a})$  for every  $(\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}$  where  $(\check{\mu}, \check{a})$  is the closest point (or one of the closest points, in case of equality) in  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$  with respect to  $d_{\mathfrak{S}}$  and  $d_{\mathcal{A}}$ . Based on simplexes discretizations, this point satisfies  $\|\mu - \check{\mu}\|_{d_{\mathfrak{S}}} \leq \varepsilon_{\mathfrak{S}}$  and  $d_{\mathcal{A}}(\tilde{a}, \check{a}) \leq \varepsilon_{\mathcal{A}}$ .

**Theorem 35.** Let  $\delta \in (0, 1)$  and  $\varepsilon > 0$ . Assume Assumptions (H3)–(H5) hold. Consider learning rates  $(\eta_n)_n$  satisfying: There exists  $\kappa \in (1/2, 1)$  such that for every  $(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \check{\mathcal{A}}$ ,  $\eta_n := \eta_n(\check{\mu}, \check{a}) = 1/(1 + C(n, \check{\mu}, \check{a}))^\kappa$  for each  $n \geq 0$ , where  $C(n, \check{\mu}, \check{a})$  is the number of times up to  $n$  that the pair  $(\check{\mu}, \check{a})$  has been visited in Algorithm 1. If the number of episodes  $N_{\text{epi}}$  is of order

$$(34) \quad \Omega \left( \left( \frac{(T_{cov}(\delta))^{1+3\kappa} \check{J}_{bound}^2 \ln(|\check{\mathfrak{S}}| |\check{\mathcal{A}}|^{|\mathcal{S}|} \check{J}_{bound}/(2\delta\beta\varepsilon))}{\beta^2 \varepsilon^2} \right)^{\frac{1}{\kappa}} + \left( \frac{(T_{cov}(\delta))}{\beta} \ln \left( \frac{\check{J}_{bound}}{\varepsilon} \right) \right)^{\frac{1}{1-\kappa}} \right),$$

then with probability  $1 - \delta$ , for all  $(\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}$ ,

$$\left| \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) - \bar{Q}^*(\mu, \mu \hat{\otimes} \tilde{a}) \right| \leq \varepsilon',$$

where  $\varepsilon' = \varepsilon + \left( \frac{\gamma}{1-\gamma} L_{\bar{J}^*} + L_{\tilde{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}} \right) \varepsilon_{\mathfrak{S}} + \frac{1}{1-\gamma} \left( L_{\tilde{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}} \right) \varepsilon_{\mathcal{A}}.$

Note that  $\varepsilon$  can be chosen as small as desired provided  $N_{\text{epi}}$  is large enough. The second and the third terms in the error  $\varepsilon'$  are proportional to  $\varepsilon_{\mathfrak{S}}$  and  $\varepsilon_{\mathcal{A}}$ , which is somehow unavoidable in general due to the projection on the finite sets  $\check{\mathfrak{S}}$  and  $\check{\mathcal{A}}$ . However, this error vanishes as  $\varepsilon_{\mathfrak{S}} \rightarrow 0$  and  $\varepsilon_{\mathcal{A}} \rightarrow 0$ , i.e., as  $\check{\mathfrak{S}}$  and  $\check{\mathcal{A}}$  are better and better approximations of  $\mathcal{P}(S)$  and  $\mathcal{P}(A)$  respectively.

We prove this result below. The proof can be summarized in the following three steps: **(1)** For  $N_{\text{epi}}$  large enough, we have  $\check{Q}_{N_{\text{epi}}} \approx \bar{Q}^*$  on  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$ ; **(2)**  $\bar{Q}^* \approx \bar{Q}$  on  $\check{\mathfrak{S}} \times \check{\mathcal{A}}$ ; **(3)** For every  $(\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}$ ,

$\tilde{Q}^*(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) \approx \tilde{Q}^*(\mu, \tilde{a})$ . The first step relies on standard Q-learning convergence results [19], while the two other steps stem from the regularity assumptions and the approximation of  $(\mathfrak{S}, \mathcal{A})$  by  $(\check{\mathfrak{S}}, \check{\mathcal{A}})$ .

*Proof of Theorem 35.* Recall that we denote by  $\check{J}^*$  and  $\check{Q}^*$  respectively the state value function and the state-action value function of the projected MFC problem defined by (32)–(33).

We first note that, for every  $(\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}$ ,

$$\begin{aligned} & \left| \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) - \tilde{Q}^*(\mu, \tilde{a}) \right| \\ & \leq \left| \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) - \check{Q}^*(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) \right| \\ & \quad + \left| \check{Q}^*(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) - \tilde{Q}^*(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) \right| \\ & \quad + \left| \tilde{Q}^*(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})) - \tilde{Q}^*(\mu, \tilde{a}) \right|. \end{aligned}$$

We then split the proof into three steps, which consist in bounding from above each term in the right hand side.

**Step 1.** We first analyze the difference between  $\check{Q}_{N_{\text{epi}}}$  and  $\check{Q}^*$ . This comes from standard convergence results on Q-learning for finite state-action spaces. More precisely, under Assumptions **(H3)** and **(H5)**, with our choice of learning rates, and given that  $N_{\text{epi}}$  is of order (34), we can apply Theorem 4 and Corollary 34 in [19] for asynchronous Q-learning and polynomial learning rates, and we obtain that, with probability at least  $1 - \delta$ ,

$$\|\check{Q}_{N_{\text{epi}}} - \check{Q}^*\|_{\infty} = \sup_{(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \check{\mathcal{A}}} \left| \check{Q}_{N_{\text{epi}}}(\check{\mu}, \check{a}) - \check{Q}^*(\check{\mu}, \check{a}) \right| \leq \varepsilon.$$

**Step 2.** We then turn our attention to the difference between  $\check{Q}^*$  and  $\tilde{Q}^*$ . The analysis amounts to say that the projection on  $\check{\mathfrak{S}}$  realized at each step does not perturb too much the value function. Recall that for some given common noise  $\varepsilon^0$ , the operator  $\check{\Phi}^{\varepsilon^0} : \check{\mathfrak{S}} \times \check{\mathcal{A}} \rightarrow \check{\mathfrak{S}}$  is given by  $\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) = \text{Proj}_{\check{\mathfrak{S}}} \circ \bar{F}(\check{\mu}, \check{\mu} \otimes \check{a}, \varepsilon^0)$ . Likewise, we denote the transition dynamic with  $\bar{F}$  by a function  $\Phi^{\varepsilon^0} : \mathfrak{S} \times \mathcal{A} \rightarrow \mathfrak{S}$  such that:

$$\Phi^{\varepsilon^0}(\mu, \tilde{a}) = \bar{F}(\mu, \mu \otimes \tilde{a}, \varepsilon^0), \quad \forall (\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}.$$

Let us start by noting that, for every  $(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \check{\mathcal{A}}$ ,

$$\begin{aligned} & \left| \check{Q}^*(\check{\mu}, \check{a}) - \tilde{Q}^*(\check{\mu}, \check{a}) \right| \\ & \leq \gamma \mathbb{E} \left[ \left| \check{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) - \bar{J}^*(\Phi^{\varepsilon^0}(\check{\mu}, \check{a})) \right| \right] \\ & \leq \gamma \mathbb{E} \left[ \left| \check{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) - \bar{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) \right| + \left| \bar{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) - \bar{J}^*(\Phi^{\varepsilon^0}(\check{\mu}, \check{a})) \right| \right] \\ & \leq \gamma \mathbb{E} \left[ \left| \inf_{\check{a}' \in \check{\mathcal{A}}} \check{Q}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}), \check{a}') - \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}), \tilde{a}') \right| \right] + \gamma L_{\bar{J}^*} \mathbb{E} \left[ \left\| \check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) - \Phi^{\varepsilon^0}(\check{\mu}, \check{a}) \right\|_{d_{\check{\mathfrak{S}}} \right], \end{aligned}$$

where the last inequality holds by Lipschitz continuity of  $\bar{J}^*$  on  $\check{\mathfrak{S}}$ , see Assumption **(H4)**.

The second term in the last inequality can be bounded using the simplex discretization properties and Assumption **(H3)**:

$$\mathbb{E} \left[ \|\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) - \Phi^{\varepsilon^0}(\check{\mu}, \check{a})\|_{d_{\mathfrak{S}}} \right] = \mathbb{E}_{\varepsilon_1^0} \left[ \|\text{Proj}_{\mathfrak{S}} \circ \bar{F}(\check{\mu}, \check{a}, \varepsilon^0) - \bar{F}(\check{\mu}, \check{a}, \varepsilon^0)\|_{d_{\mathfrak{S}}} \right] \leq \varepsilon_{\mathfrak{S}}.$$

For the first term, let  $\check{\mu}' = \check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) \in \check{\mathfrak{S}}$  to alleviate the notation, and let us consider  $\check{a}_1^* \in \check{\mathcal{A}}$  and  $\check{a}_2^* \in \mathcal{A}$  satisfying:  $\check{Q}^*(\check{\mu}', \check{a}_1^*) = \inf_{\check{a}' \in \check{\mathcal{A}}} \check{Q}^*(\check{\mu}', \check{a}')$  and  $\tilde{Q}^*(\check{\mu}', \check{a}_2^*) = \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\check{\mu}', \tilde{a}')$ . The existence of  $\check{a}_1^*$  and  $\check{a}_2^*$  is guaranteed respectively by finiteness of  $\mathfrak{S} \times \check{\mathcal{A}}$  and by Lemma 33.

We observe that

$$\begin{aligned} & \check{Q}^*(\check{\mu}', \check{a}_1^*) - \tilde{Q}^*(\check{\mu}', \check{a}_2^*) \\ &= \left( \check{Q}^*(\check{\mu}', \check{a}_1^*) - \check{Q}^*(\check{\mu}', \text{Proj}_{\check{\mathcal{A}}}(\check{a}_2^*)) \right) + \left( \check{Q}^*(\check{\mu}', \text{Proj}_{\check{\mathcal{A}}}(\check{a}_2^*)) - \tilde{Q}^*(\check{\mu}', \text{Proj}_{\check{\mathcal{A}}}(\check{a}_2^*)) \right) \\ & \quad + \left( \tilde{Q}^*(\check{\mu}', \text{Proj}_{\check{\mathcal{A}}}(\check{a}_2^*)) - \tilde{Q}^*(\check{\mu}', \check{a}_2^*) \right) \\ &\leq 0 + \sup_{(\check{\mu}, \check{a}) \in \mathfrak{S} \times \check{\mathcal{A}}} \left| (\check{Q}^* - \tilde{Q}^*)(\check{\mu}, \check{a}) \right| + \left( \tilde{f}(\check{\mu}', \text{Proj}_{\check{\mathcal{A}}}(\check{a}_2^*)) + \gamma \mathbb{E}_{(\varepsilon^0)'} [\bar{J}^*(\bar{F}(\check{\mu}', \text{Proj}_{\check{\mathcal{A}}}(\check{a}_2^*), (\varepsilon^0)'))] \right) \\ & \quad - \left( \tilde{f}(\check{\mu}', \check{a}_2^*) + \gamma \mathbb{E}_{(\varepsilon^0)'} [\bar{J}^*(\bar{F}(\check{\mu}', \check{a}_2^*), (\varepsilon^0)'))] \right) \\ &\leq \|\check{Q}^* - \tilde{Q}^*\|_{\infty} + (L_{\tilde{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}}) \varepsilon_{\mathfrak{A}}. \end{aligned}$$

On the other hand,

$$\check{Q}^*(\check{\mu}', \check{a}_1^*) - \tilde{Q}^*(\check{\mu}', \check{a}_2^*) = - \left( \tilde{Q}^*(\check{\mu}', \check{a}_2^*) - \tilde{Q}^*(\check{\mu}', \check{a}_1^*) \right) - \left( \check{Q}^*(\check{\mu}', \check{a}_1^*) - \check{Q}^*(\check{\mu}', \check{a}_1^*) \right) \geq -\|\check{Q}^* - \tilde{Q}^*\|_{\infty}.$$

Combining the above bounds yields that for every  $(\check{\mu}, \check{a}) \in \mathfrak{S} \times \check{\mathcal{A}}$ ,

$$\left| \check{Q}^*(\check{\mu}, \check{a}) - \tilde{Q}^*(\check{\mu}, \check{a}) \right| \leq \gamma \left( \|\check{Q}^* - \tilde{Q}^*\|_{\infty} + (L_{\tilde{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}}) \varepsilon_{\mathfrak{A}} \right) + \gamma L_{\bar{J}^*} \varepsilon_{\mathfrak{S}}.$$

Consequently,

$$\|\check{Q}^* - \tilde{Q}^*\|_{\infty} \leq \frac{\gamma}{1 - \gamma} \left( (L_{\tilde{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}}) \varepsilon_{\mathfrak{A}} + L_{\bar{J}^*} \varepsilon_{\mathfrak{S}} \right).$$

**Step 3.** Last, we look at the difference between  $\check{Q}^*(\text{Proj}_{\mathfrak{S}}(\mu), \text{Proj}_{\check{\mathcal{A}}}(\tilde{a}))$  and  $\tilde{Q}^*(\mu, \tilde{a})$ . For every  $\mu \in \mathfrak{S}$  and  $\tilde{a} \in \mathcal{A}$ , letting  $\check{\mu} = \text{Proj}_{\mathfrak{S}}(\mu)$  and  $\check{a} = \text{Proj}_{\check{\mathcal{A}}}(\tilde{a})$  to alleviate the notation, we have  $\|\check{\mu} - \mu\|_{d_{\mathfrak{S}}} \leq \varepsilon_{\mathfrak{S}}$  and  $\|\check{a} - \tilde{a}\|_{d_{\mathcal{A}}} \leq \varepsilon_{\mathfrak{A}}$ . We obtain

$$\begin{aligned} \left| \check{Q}^*(\check{\mu}, \check{a}) - \tilde{Q}^*(\mu, \tilde{a}) \right| &\leq \left| \tilde{f}(\check{\mu}, \check{a}) - \tilde{f}(\mu, \tilde{a}) \right| + \gamma \mathbb{E} \left[ \left| \inf_{\check{a}' \in \check{\mathcal{A}}} \check{Q}^*(\Phi(\check{\mu}, \check{a}), \check{a}') - \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\Phi(\mu, \tilde{a}), \tilde{a}') \right| \right] \\ &\leq L_{\tilde{f}} (\|\check{\mu} - \mu\|_{d_{\mathfrak{S}}} + \|\check{a} - \tilde{a}\|_{d_{\mathcal{A}}}) + \gamma \mathbb{E} \left[ |\bar{J}^*(\bar{F}(\check{\mu}, \check{a}, \varepsilon^0) - \bar{J}^*(\bar{F}(\check{\mu}, \check{a}, \varepsilon^0)))| \right] \\ &\leq L_{\tilde{f}} (\varepsilon_{\mathfrak{S}} + \varepsilon_{\mathfrak{A}}) + \gamma L_{\bar{J}^*} \mathbb{E} \left[ \|\bar{F}(\check{\mu}, \check{a}, \varepsilon^0) - \bar{F}(\mu, \tilde{a}, \varepsilon^0)\|_{d_{\mathfrak{S}}} \right] \\ &\leq (L_{\tilde{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}}) (\varepsilon_{\mathfrak{S}} + \varepsilon_{\mathfrak{A}}), \end{aligned}$$

where we used the Lipschitz continuity of  $\tilde{f}, \bar{J}^*, \bar{F}$  and the assumption on  $\mathfrak{S}$ , see Assumptions **(H3)**, **(H4)** and the simplex discretization properties.  $\square$

5.3.2. *Q-learning with controls that are pure at level-0.* The above method is designed for the case where one looks for optimal actions that are potentially randomized at the individual level. Searching in the space  $\mathcal{P}(A)$  comes with a computational cost that is reflected in the bounds through the cardinality of the discrete simplex  $\check{\mathcal{A}}$ . In some situations it can be interesting to directly search for actions that are pure at the individual level.

In this case, instead of (29), the set of strategy functions is (for simplicity we keep the notation  $\mathcal{A}$ ):

$$\mathcal{A} := \{\tilde{a} : S \rightarrow A\} = A^S.$$

In Algorithm 1, we replace  $\check{a} \in \check{\mathcal{A}}$  by  $\tilde{a} \in \mathcal{A}$ .

**Theorem 36.** *Let  $\delta \in (0, 1)$  and  $\varepsilon > 0$ . Assume Assumptions (H3)–(H5) hold. Consider learning rates  $(\eta_n)_n$  satisfying: There exists  $\kappa \in (1/2, 1)$  such that for every  $(\check{\mu}, \tilde{a}) \in \check{\mathfrak{S}} \times \tilde{\mathcal{A}}$ ,  $\eta_n := \eta_n(\check{\mu}, \tilde{a}) = 1/(1 + C(n, \check{\mu}, \tilde{a}))^\kappa$  for each  $n \geq 0$ , where  $C(n, \check{\mu}, \tilde{a})$  is the number of times up to  $n$  that the pair  $(\check{\mu}, \tilde{a})$  has been visited in Algorithm 1. If the number of episodes  $N_{\text{epi}}$  is of order*

$$(35) \quad \Omega \left( \left( \frac{(T_{\text{cov}}(\delta))^{1+3\kappa} \check{J}_{\text{bound}}^2 \ln(|\check{\mathfrak{S}}| |A|^{|\mathcal{S}|} \check{J}_{\text{bound}} / (2\delta\beta\varepsilon))}{\beta^2 \varepsilon^2} \right)^{\frac{1}{\kappa}} + \left( \frac{(T_{\text{cov}}(\delta))}{\beta} \ln \left( \frac{\check{J}_{\text{bound}}}{\varepsilon} \right) \right)^{\frac{1}{1-\kappa}} \right),$$

then with probability  $1 - \delta$ , for all  $(\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}$ ,

$$\left| \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \tilde{a}) - \bar{Q}^*(\mu, \mu \otimes \tilde{a}) \right| \leq \varepsilon',$$

where  $\varepsilon' = \varepsilon + \left( \frac{\gamma}{1-\gamma} L_{\bar{J}^*} + L_{\bar{f}} + \gamma L_{\bar{J}^*} L_{\bar{F}} \right) \varepsilon_{\mathfrak{S}}$ .

The above result provides convergence guarantee for the Q-function. Let us now derive a consequence in terms of the optimizer. To this end, we will use the following additional assumption on the gap between the values of the best and second-best actions, which is rather standard in approximation algorithms based on tabular Q-functions [20, 8].

**Assumption (H6). Action gap:** *There exists  $K_A > 0$  such that:*

$$\tilde{Q}^*(\check{\mu}, \tilde{a}) - \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\check{\mu}, \tilde{a}') \geq K_A, \quad \check{\mu} \in \check{\mathfrak{S}}, \tilde{a} \in \mathcal{A} \setminus \arg \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\check{\mu}, \tilde{a}').$$

To recover minimizers or approximate minimizers, it will be convenient to work with the following operators. In general, they are defined on the vector space  $\mathbb{R}^{\mathfrak{m}}$ . For  $\tau > 0$  and  $x = (x_1, \dots, x_{\mathfrak{m}}) \in \mathbb{R}^{\mathfrak{m}}$ , we define  $\text{softmin}_{\tau} : \mathbb{R}^{\mathfrak{m}} \rightarrow \mathbb{R}^{\mathfrak{m}}$  by

$$\text{softmin}_{\tau}(x) = (e^{-\tau x_1}, \dots, e^{-\tau x_{\mathfrak{m}}}) / \sum_j e^{-\tau x_j}.$$

For  $x \in \mathbb{R}^{\mathfrak{m}}$ , we define  $\text{argmine} : \mathbb{R}^{\mathfrak{m}} \rightarrow [0, 1]^{\mathfrak{m}}$  by

$$\text{argmine}(x) = (\mathbf{1}_{i \in \arg \min(x)})_{i=1}^{\mathfrak{m}} / |\arg \min(x)|.$$

where  $\arg \min(x) = \{j \in \{1, \dots, \mathfrak{m}\} : x_j = \min\{x_1, \dots, x_{\mathfrak{m}}\}\}$ . In the sequel, we use these operators with the dimension  $\mathfrak{m} = |\mathcal{A}| = |A|^{|\mathcal{S}|}$ . For any function  $q : \mathcal{A} \rightarrow \mathbb{R}$ , we identify  $q$  with the vector  $(q(\tilde{a}))_{\tilde{a} \in \mathcal{A}}$ .

**Corollary 37.** *Assume the same assumptions as in Theorem 36 hold and, in addition, that Assumption (H6) holds. Let  $\check{Q}_{N_{\text{epi}}}$  be the table returned by Algorithm 1, and let  $\varepsilon'$  be as in Theorem 36. Then for every  $\check{\mu} \in \check{\mathfrak{S}}$ ,*

$$\|\text{softmin}_{\tau}(\check{Q}_{N_{\text{epi}}}(\check{\mu}, \cdot)) - \text{argmine}(\tilde{Q}^*(\check{\mu}, \cdot))\|_2 \leq \tau \varepsilon' \sqrt{|\mathcal{A}|} + 2e^{-\tau K_A} |\mathcal{A}|.$$

The proof is provided below. The argmine in the second term is here in case there are several optimal controls. The softmin regularizes the best action predicted by the estimation  $\tilde{Q}_{N_{\text{epi}}}$  of the function  $\tilde{Q}^*$ .

**Remark 38.** *Imagine we want the error bound in Corollary 37, to be smaller than some  $\delta > 0$ . It is sufficient to have: for the second term:  $\tau \geq \frac{1}{K_A} \log\left(\frac{|\mathcal{A}|}{\delta/4}\right)$ ; and for the first term:  $\varepsilon' \leq \delta/(2\tau\sqrt{|\mathcal{A}|}) = \delta K_A/(2|\mathcal{A}|^{1/2} \log(\frac{|\mathcal{A}|}{\delta/4}))$ . Then both terms in the error bound will be smaller than  $\delta/2$ . Notice that, contrary to Theorem 35, here we do not need to approximate the probability space of action  $\mathcal{P}(A)$  by  $\tilde{\mathcal{A}}$  with an  $\varepsilon_{\mathcal{A}}$ -net, hence the error bound in Theorem 36 is independent of any  $\varepsilon_{\mathcal{A}}$ . So it is possible to choose  $\tau$  and to make  $\varepsilon'$  as small as we want.*

*Proof of Corollary 37.* We use [23, Proposition 4], which states that  $\text{softmin}_\tau$  is  $\tau$ -Lipschitz and [27, Lemma 7], which states that for  $(x_i)_{i=1,\dots,\mathfrak{m}}$ ,

$$\|\text{softmin}_\tau(x) - \text{argmine}(x)\|_2 \leq 2\mathfrak{m}e^{-\tau\delta},$$

where  $\delta = \inf_{x_j > \inf(x)} x_j - \inf(x)$ , and  $\delta = \infty$  if all  $x_i$  are equal. We can apply this latter result to  $\tilde{Q}^*(\tilde{\mu}, \cdot)$  thanks to assumption **(H6)**, with  $\mathfrak{m} = |\tilde{\mathcal{A}}|$  and  $\delta = K_A$ . Combining this with Theorem 35, we have, for every  $\tilde{\mu}$ ,

$$\begin{aligned} & \left\| \text{softmin}_\tau\left(\tilde{Q}^*(\tilde{\mu}, \cdot)\right) - \text{argmine}\left(\tilde{Q}^*(\tilde{\mu}, \cdot)\right) \right\|_2 \\ & \leq \left\| \text{softmin}_\tau\left(\tilde{Q}^*(\tilde{\mu}, \cdot)\right) - \text{softmin}_\tau\left(\tilde{Q}^*(\tilde{\mu}, \cdot)\right) \right\|_2 + \left\| \text{softmin}_\tau\left(\tilde{Q}^*(\tilde{\mu}, \cdot)\right) - \text{argmine}\left(\tilde{Q}^*(\tilde{\mu}, \cdot)\right) \right\|_2 \\ & \leq \tau \left\| \tilde{Q}^*(\tilde{\mu}, \cdot) - \tilde{Q}^*(\tilde{\mu}, \cdot) \right\|_2 + 2|\mathcal{A}|e^{-\tau K_A} \\ & \leq \tau\sqrt{|\mathcal{A}|} \sup_{\tilde{a}' \in \mathcal{A}} \left| \tilde{Q}^*(\tilde{\mu}, \tilde{a}') - \tilde{Q}^*(\tilde{\mu}, \tilde{a}') \right| + 2|\mathcal{A}|e^{-\tau K_A} \\ & \leq \tau\varepsilon'\sqrt{|\mathcal{A}|} + 2e^{-\tau K_A}|\mathcal{A}|. \end{aligned}$$

□

#### 5.4. Deep reinforcement learning for MFMDP.

The above method has the advantage to be simple enough to let us carry out a detailed analysis. However, it cannot be used in practice for large state or actions spaces because of the prohibitive computational cost due to the discretization of the simplexes. An alternative is to work directly with continuous spaces, in which case the policies and value functions cannot be represented in a tabular way. Instead, we can rely on function approximation. To this end, we now propose to use methods from deep reinforcement learning which are more suitable for continuous spaces. The motivations are twofold.

First, if  $S$  and  $A$  are finite but we want to learn an optimal policy that is potentially randomized at level-0, the discretization approach proposed in § 5.3.1 has a complexity that increases with the number of points in the discretization of  $\mathcal{P}(A)$ , which itself increases exponentially quickly with the cardinality of  $A$ . For this reason, it can be interesting to tackle directly  $\tilde{A} = \mathcal{P}(A)$  as a continuous action space and to use deep RL methods for continuous action space MDPs.

Second, some MFC problems are naturally posed with a continuous state space  $S$ . In this case, under mild conditions, the optimal policy is in fact non-randomized not only at the level-1 but even at the

level-0. However, the state of the mean field MDP is an element of the infinite dimensional space  $\mathcal{P}(S)$ . From a numerical viewpoint, we need two ingredients: (1) a finite-dimensional approximation of the mean field and (2) a parameterized approximation of the value function or the policy taking this finite-dimensional representation of the mean field state as an input. For the second point, we will again use deep neural networks. For the first point, for the sake of definiteness, we choose to simply replace  $\mathcal{P}(S)$  by  $\mathcal{P}(\tilde{S})$  where  $\tilde{S}$  is a discretization of  $S$  with a finite number of points. We assume that, given  $\tilde{\mu} \in \mathcal{P}(\tilde{S})$  and  $\bar{a} : S \rightarrow A$ , one can get from the environment a sample of the next state and the associated cost  $\tilde{f}(\tilde{\mu}, \bar{a})$ . The problem thus boils down to an MDP with finite dimensional (but potentially continuous) state and action spaces. Such MDPs can be solved with a variety of deep RL algorithms. In the sequel, we provide numerical illustrations based on the Deep Deterministic Policy Gradient (DDPG) proposed in [32]. It relies on two neural networks, one for the Q-function (the critic) and one for the policy (the actor). The heart of the algorithm consists in updating alternatively the critic by minimizing an empirical square error and the actor by making one step of gradient descent. To improve exploration, a Gaussian noise  $\epsilon_{n+1}^a$  is added to the action prescribed by the actor. Furthermore, for more stability, target networks are also added. The algorithm is summarized in our setting in Algorithm 2 in the Appendix E.

## 6. Numerical Examples

### 6.1. Example 1: Cyber security model.

For a first testbed, we start with a finite state problem. We revisit the cyber security example introduced in [30], but here from the point of view of a central planner (such as a large company or a state) trying to protect its computers against the attacks of a hacker. The situation can thus be phrased as a MFC problem.

In this model, the population consists of a large group of computers which can be either defended (D) or undefended (U), and either infected (I) or susceptible (S) of infection. Hence the set  $S$  has four elements corresponding to the four possible combinations: DI, DS, UI, US. The action set is  $A = \{0, 1\}$ , where 0 is interpreted as the fact that the central planner is satisfied with the current level of protection (D or U) of the computer under consideration, whereas 1 means that she wants to change this level of protection. In the latter case, the update occurs at a (fixed) rate  $\lambda > 0$ . If the controls are pure at level-0, at each of the four states, the central planner only chooses one action per state and applies it to all the computers at that state. If the controls are mixed at level-0, then for each state, she chooses a distribution over actions and then each computer in this state picks independently an action according to the chosen distribution. When infected, each computer may recover at rate  $q_{rec}^D$  or  $q_{rec}^U$  depending on whether it is defended or not. On the other hand, a computer may be infected either directly by a hacker, at rate  $v_H q_{inf}^D$  (resp.  $v_H q_{inf}^U$ ) if it is defended (resp. undefended), or by undefended infected computers, at rate  $\beta_{UU}\mu(\{UI\})$  (resp.  $\beta_{UD}\mu(\{UI\})$ ) if it is undefended (resp. defended), or by defended infected computers, at rate  $\beta_{DU}\mu(\{DI\})$  (resp.  $\beta_{DD}\mu(\{DI\})$ ) if it is undefended (resp. defended). Here  $v_H$  can be interpreted as the attack intensity parameter.

In short, the transition matrix is given by:

$$(36) \quad P^{\mu,a} = \begin{pmatrix} \dots & P_{DS \rightarrow DI}^{\mu,a} & \lambda a & 0 \\ q_{rec}^D & \dots & 0 & \lambda a \\ \lambda a & 0 & \dots & P_{US \rightarrow UI}^{\mu,a} \\ 0 & \lambda a & q_{rec}^U & \dots \end{pmatrix}$$

where

$$\begin{aligned} P_{DS \rightarrow DI}^{\mu, a} &= v_H q_{inf}^D + \beta_{DD} \mu(\{DI\}) + \beta_{UD} \mu(\{UI\}), \\ P_{US \rightarrow UI}^{\mu, a} &= v_H q_{inf}^U + \beta_{UU} \mu(\{UI\}) + \beta_{DU} \mu(\{DI\}), \end{aligned}$$

and all the instances of  $\dots$  should be replaced by the negative of the sum of the entries of the row in which  $\dots$  appears on the diagonal. At each time step, the central planner pays a protection cost  $k_D > 0$  for each defended computer, and a penalty  $k_I > 0$  for each infected computer. The instantaneous cost in the MFMDP is thus defined as:

$$\bar{f}(\mu, \bar{a}) = k_D \mu(\{DI, DS\}) + k_I \mu(\{DI, UI\}), \quad (\mu, \bar{a}) \in \bar{S} \times \bar{A}.$$

The optimal control and optimal flow of distributions can be characterized by a forward-backward ODE system which can be obtained in way similar to what is done in the MFG setting e.g. in [12, § 7.2.3]. We will use this solution as a benchmark.

**Tabular Q-learning.** For the sake of illustration, we present results obtained by tabular Q-learning with simplex discretization as described in § 5.3. The state space for the population distribution is  $\bar{S}$ , which is identified with the simplex  $\check{\mathfrak{S}} = \{(\mu^{(i)})_{i=1, \dots, 4} \in [0, 1]^4 : \sum_i \mu^{(i)} = 1\}$ . To follow the original setting considered in [30], we consider pure controls, both at the common and idiosyncratic levels (level-0 and level-1). So we identify  $\bar{A}$  with the set of functions  $A^S$ , which is finite and of cardinality  $2^4 = 16$ .

We replace  $\check{\mathfrak{S}}$  by the finite set:

$$\check{\mathfrak{S}} = \left\{ (\mu^{(i)})_{i=1, \dots, 4} \in [0, 1/N_m, \dots, 1 - 1/N_m, 1]^4 : \sum_i \mu^{(i)} = 1 \right\},$$

where  $[0, 1/N_m, \dots, 1 - 1/N_m, 1]$  is a uniform grid over  $[0, 1]$  with  $N_m + 1 \geq 2$  points. We then aim at computing the Q-function for the projected MDP with finite state space  $\check{\mathfrak{S}}$  and action space  $A^S$ , that we still denote by  $\tilde{Q}^*$  although we do not consider mixed actions at the level-0. We note that, in the absence of common noise, the MFMDP is completely deterministic hence it would be enough to query once each state-action pair from the environment in order to learn the level-1 reward function and transition function, and hence to be able to compute perfectly the Q-function. However, for the sake of illustration, we stick to applying Algorithm 1, replacing both  $\mathcal{A}$  and  $\check{\mathcal{A}}$  by  $A^S$ .

In order to be able to compare with the benchmark solution obtained by the ODE method, we consider that the time steps are of size not 1 but  $\Delta t = 0.01$ . Although the problem is set on an infinite horizon, we truncate the training episodes and the plots at the horizon  $T = 10$ .

After  $N_{\text{epi}}$  episodes of Q-learning, we obtain an approximation  $\check{Q}_{N_{\text{epi}}}$  of the Q-function, from which we can recover an approximation  $\bar{a}_{N_{\text{epi}}}$  of the optimal control by taking the argmax, namely:  $\bar{a}_{N_{\text{epi}}}(\mu, \cdot) = \arg \max_{\check{a} \in A^S} \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu), \check{a})$ . We compare the flow of distributions induced by this control  $\bar{a}_{N_{\text{epi}}}$  with the optimally controlled flow computed by the ODE method. This method also allows us to compute for each  $t$  the value  $\bar{J}^*(\mu_t^*)$  along the optimal flow  $(\mu_t^*)_{t \in [0, T]}$ , in line with Lemma 34, we compare it with  $\max_{A^S} \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu_t^*), \cdot) = \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\check{\mathfrak{S}}}(\mu_t^*), \bar{a}_{N_{\text{epi}}}(\mu_t^*))$ . Figures 1–3 show the results for three initial conditions  $\mu_0$ . We see that the learnt value function approximately matches the  $\bar{J}^*$  value function, and the induced flows of distributions approximately match the benchmark ones. For these simulations, we used  $N_m = 30$ ,  $\gamma = 0.5$ , and the following parameters:

$$\begin{cases} \beta_{UU} = 0.3, \beta_{UD} = 0.4, \beta_{DU} = 0.3, \beta_{DD} = 0.4, \\ q_{rec}^D = 0.5, q_{rec}^U = 0.4, q_{inf}^D = 0.4, q_{inf}^U = 0.3, \\ v_H = 0.6, \lambda = 0.8, k_D = 0.3, k_I = 0.5. \end{cases}$$

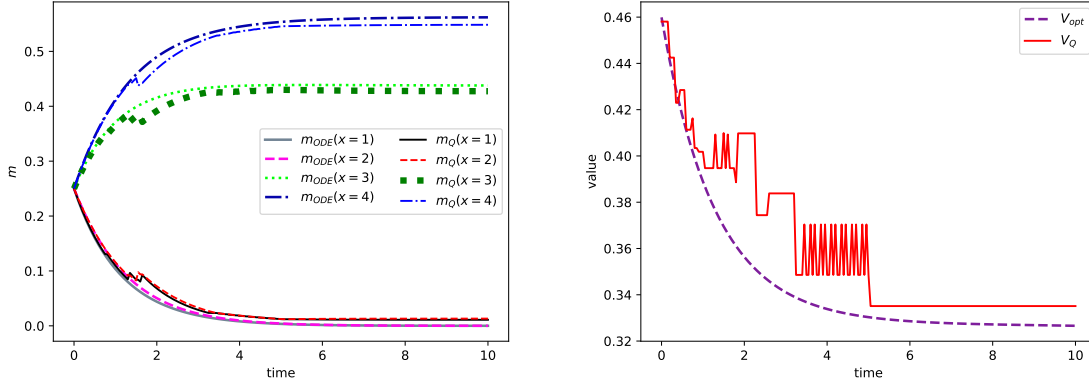


FIGURE 1. Example 1: Cyber security model. Test case 1:  $m_0 = (1/4, 1/4, 1/4, 1/4)$ . Left: Evolution of the distribution when using the benchmark optimal control ( $m_{ODE}$ ) or the control recovered from the learnt Q-function ( $m_Q$ ). Right: state value function using the benchmark solution ( $V_{opt}$ ) or the learnt Q-function ( $V_Q$ ) along the optimal mean field flow. The benchmark solution is obtained using the ODE method.

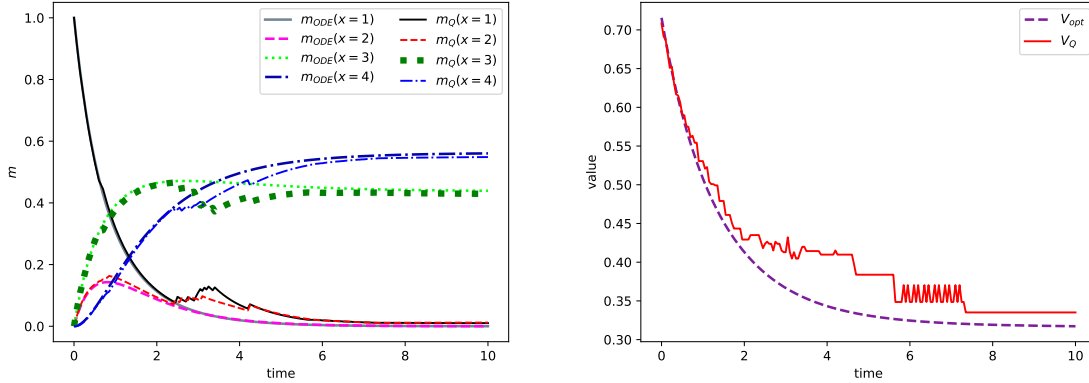


FIGURE 2. Example 1: Cyber security model. Test case 2:  $m_0 = (1, 0, 0, 0)$ . Left: Evolution of the distribution when using the benchmark optimal control ( $m_{ODE}$ ) or the control recovered from the learnt Q-function ( $m_Q$ ). Right: state value function using the benchmark solution ( $V_{opt}$ ) or the learnt Q-function ( $V_Q$ ) along the optimal mean field flow. The benchmark solution is obtained using the ODE method.

**Deep Deterministic Policy Gradient.** The solution can also be learnt by using the deep RL algorithm given in Algorithm 2 instead of mean field Q-learning given in Algorithm 1. In the present case, this approach has the advantage of avoiding the discretization of  $\mathcal{P}(\mathcal{S})$  since we instead directly deal with the distribution as a vector in dimension 4. Since, with this method, it is possible to allow the control to take continuous values, we replace  $A = \{0, 1\}$  by  $A = [0, 1]$ .

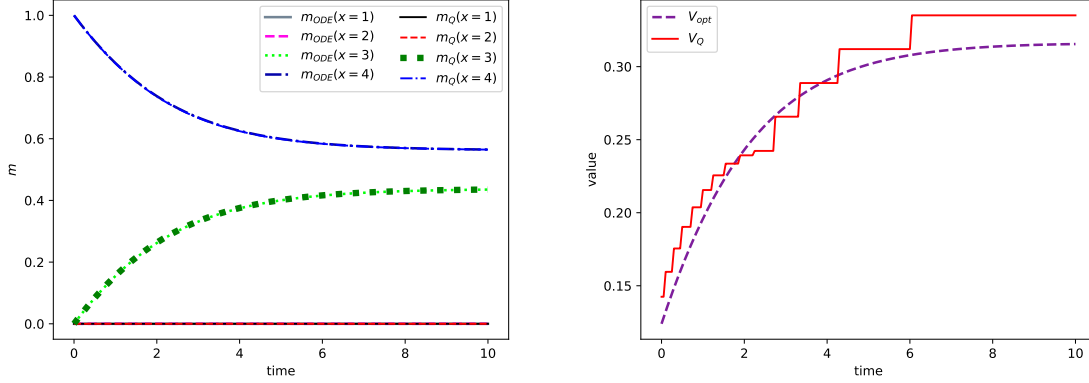


FIGURE 3. Example 1: Cyber security model. Test case 3:  $m_0 = (0, 0, 0, 1)$ . Left: Evolution of the distribution when using the benchmark optimal control ( $m_{ODE}$ ) or the control recovered from the learnt Q-function ( $m_Q$ ). Right: state value function using the benchmark solution ( $V_{opt}$ ) or the learnt Q-function ( $V_Q$ ) along the optimal mean field flow. The benchmark solution is obtained using the ODE method.

Furthermore, to make things more interesting, we consider that the attack intensity parameter  $v_H$  is stochastic. Since its value affects the evolution of the whole population, we model this using a common noise. We replace the state distribution by a conditional state distribution, conditioned on the realization of  $\epsilon^0$  up to the current time. To wit, let  $(\epsilon_n^0)_{n \geq 1}$  be a sequence of i.i.d. random variables with Gaussian distribution. Let  $v_{H,n+1} = v_{H,n} + \epsilon_{n+1}^0$ ,  $n \geq 0$ ,  $v_{H,0}$  given. The evolution from time  $n$  to time  $n+1$  of the state distribution is by the transition matrix defined in (36) but with the constant  $P_{DS \rightarrow DI}^{\mu,a}$ ,  $P_{US \rightarrow UI}^{\mu,a}$  replaced by the following stochastic coefficients that evolve in time due to the fact that  $v_H$  is replaced by a stochastic process:

$$\begin{aligned} P_{DS \rightarrow DI,n}^{\mu,a} &= v_{H,n} q_{inf}^D + \beta_{DD} \mu(\{DI\}) + \beta_{UD} \mu(\{UI\}), \\ P_{US \rightarrow UI,n}^{\mu,a} &= v_{H,n} q_{inf}^U + \beta_{UU} \mu(\{UI\}) + \beta_{DU} \mu(\{DI\}), \end{aligned}$$

Using the DDPG method described above, we train the neural networks by picking at each episode a random initial distribution  $\mu$  and a random sequence of common noises  $\epsilon^0$ . Fig. 4 displays the evolution of the population when using the learnt control starting from five initial distributions of the testing set and one initial distribution of the training set. The testing set of initial distributions is:  $\{(0.25, 0.25, 0.25, 0.25), (1, 0, 0, 0), (0, 0, 0, 1), (0.3, 0.1, 0.3, 0.1), (0.5, 0.2, 0.2, 0.1)\}$ . Consistently with the case without common noise, we see that the distribution always evolves towards a configuration in which there is no defended agents, and the proportion of undefended infected and undefended susceptible are roughly 0.43 and 0.57 respectively. Due to the common noise, the distribution is not perfectly stable; it oscillates around these values. Figure 5 shows from two perspectives the evolution of the mean field state dynamics when applying the learnt optimal control. The initial distributions are on a uniform grid of the simplex and time steps are distinguished by colors. We see that for any initial distribution, as time increases, the mean field states concentrate around the aforementioned point, which lies on an edge of the simplex.

## 6.2. Example 2: Discrete distribution planning.

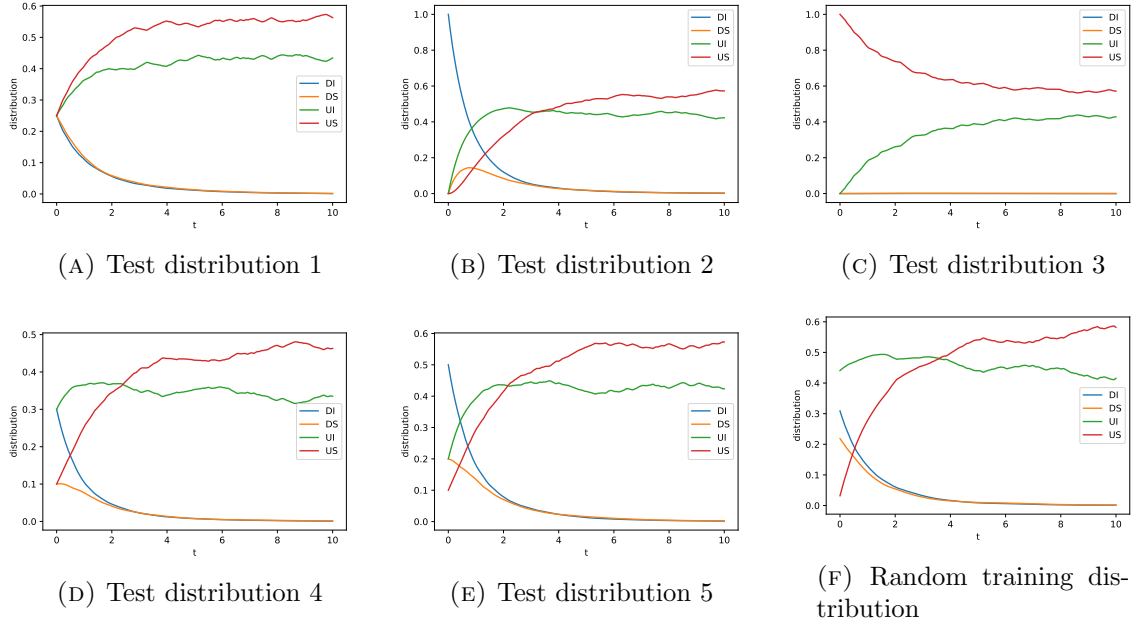


FIGURE 4. Cyber security example: Evolution of the distribution in the presence of common noise when applying the control learnt by DDPG on a testing set of five initial distributions and one random initial distribution of the training set.

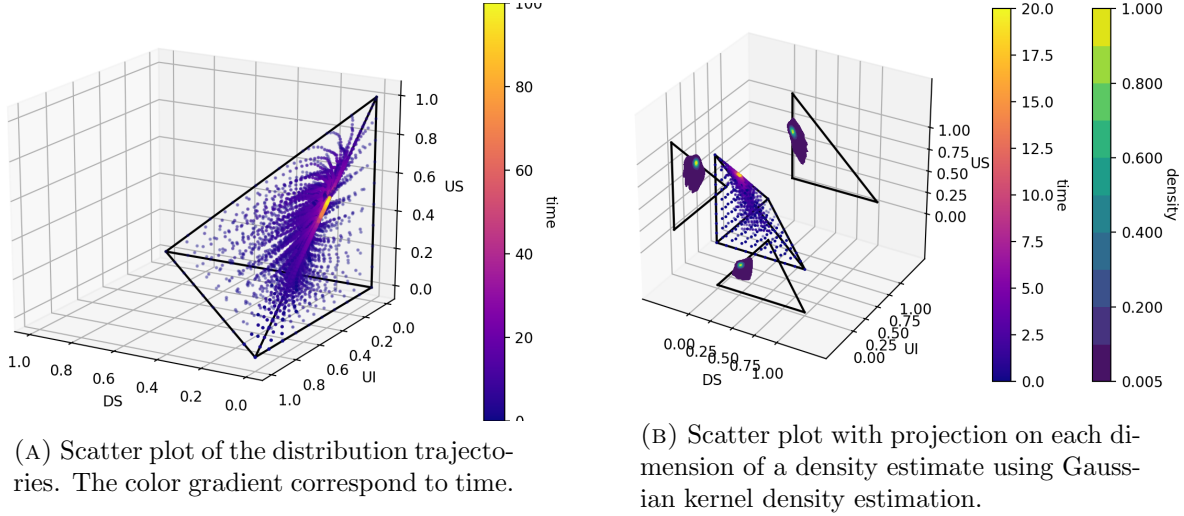


FIGURE 5. Cyber security example: Trajectories of the distribution's last three coordinates (DS, UI, US) in the simplex, when subject to common noise and controlled by the control learnt with DDPG. The initial distributions are from a uniform discretization of the simplex with mesh size 0.1 in each dimension.

We now consider an MFC problem in which the goal is to match a target distribution. We take a model with  $N_{states} = 10$  states and 3 actions (left, stay, right):  $X = \{1, \dots, 10\}$  and  $A = \{L, S, R\}$ . When at the leftmost state (resp. rightmost), the agents cannot go left (resp. right). The incurred to a representative agent is 1 if they move (i.e., they use action left or right) plus the  $L^2$  distance between the population distribution and a target distribution. Here we chose:  $(0, 0, 0.05, 0.1, 0.2, 0.3, 0.2, 0.1, 0.05, 0, 0)$  for the target distribution. There is no idiosyncratic noise. A key point is that, in this setting, except for some specific pairs of initial and target distributions, it is not possible for the population to match the target distribution unless the agents are allowed to randomize their actions at the individual level. So we use  $\mathcal{P}(A)^X$  for the level-1 action space. Hence the action space is naturally continuous and this justifies, here again, the use of the DDPG method.

We present results without and with common noise in the dynamics. In the second case, the common noise is an i.i.d. additive  $\epsilon_n^0$  at each time step  $n$ , with  $\epsilon_n^0 = -1$  with probability 0.05,  $\epsilon_n^0 = 1$  with probability 0.05 and  $\epsilon_n^0 = 0$  with probability 0.9. In both cases, the initial distributions for training are picked randomly as follows. First we pick  $x_{min}$  and  $x_{max}$  uniformly at random between 1 and  $N_{states}$ . This determines a sub-set (with periodicity)  $\{x_{min}, \dots, x_{max}\}$ . For each point in the sub-interval, a value is picked independently and uniformly at random in  $[0, 1]$ . Then the discrete distribution is normalized to have total mass 1. For numerical reasons, we stop after a finite number of time steps. Here we took 100 time steps.

Figures 6 and 7 present the results obtained without and with common noise, respectively. In each case, the results are for five different testing distributions: four fixed test distributions, as well as one random distribution. The left column displays the state distribution: initial distribution in green, target distribution in red, last distribution in purple, and the average over the last few steps before terminal time in blue. We see that the last distribution is very close to the target one so the learnt control is successful. The two columns in the middle display the control distribution at time 0 and at terminal time. For each state, the probability of picking each action is represented by a vertical bar, with one color per possible action. We see that, at initial time, the most likely choice is to move to the right (resp. left) for states on the left (resp. right) of the domain. At terminal time, the most likely choice is to stay at the current location, because the target distribution has been reached. The right column displays the trajectory of the common noise that has affected the distribution in the present run (constant equal to 0 in Figure 6). On Figure 7, we see that even when the common noise is rather strong, the learnt control manages to move the initial distribution close to the target distribution. At the bottom of each figure is displayed the evolution of the training reward (the negative of the MFC cost) along the training iterations (also called episodes) of DDPG.

### 6.3. Example 3: Swarm motion.

We then turn our attention to a model in continuous state and action spaces. More precisely, we consider a model of swarm motion with aversion to crowded regions introduced in [6] (in the context of mean field games). Although many variants are possible, we define the model in the following way in order to have an analytical solution that can be used to assess the convergence of our proposed method. We take the interval  $[0, 1]$  with periodic boundary condition, i.e. the unit torus  $\mathbb{T}$ , as the state space  $S$ . The action space is  $A = \mathbb{R}$ . The dynamics of a typical agent is driven by (1) with  $F(x, a, \mu, e, e^0) = a + e + e^0$ . In other words, the central planner chooses the velocity of each agent. The instantaneous reward of a typical agent at location  $x$  and using action  $a$  while the population's state is  $\mu$ , is defined as:  $f(x, a, \mu) = -\frac{1}{2}|a|^2 + \varphi(x) - \ln(\mu(x))$ . Here, the first term penalizes a large velocity (it can be interpreted as a kind of cost proportional to the kinetic energy of the agent),  $\varphi$  encodes spatial preferences (by giving a lower cost for certain positions in space), and the last term

models crowd aversion (it penalizes the fact of being at a location where the density of agents is high). We choose:

$$\varphi(x) = -2\pi^2 [-\sin(2\pi x) + |\cos(2\pi x)|^2] + 2\sin(2\pi x),$$

We consider that there is no common noise ( $\varepsilon_n^0 = 0$  for all  $n$ ), and the idiosyncratic noises  $\varepsilon_n$  have a Gaussian distribution. We obtain a model which, in continuous time, admits an explicit ergodic solution that we can use as a benchmark. Indeed, in this case the optimal ergodic control is given by  $\tilde{a}(x) = 2\pi \cos(2\pi x)$  and the ergodic distribution of the corresponding MKV dynamics has density  $\mu(x) = e^{2\sin(2\pi x)} / \int e^{2\sin(2\pi x')} dx'$ .

The action space being continuous, here again the use of DDPG is justified. To implement this approach, we however need a finite dimensional representation of the distribution in order to pass it to the policy network and the value function network. We replace  $\mathcal{P}(\mathbb{T})$  by a finite dimensional simplex  $\mathcal{P}(\{0, 1/N_p, \dots, 1 - 1/N_p, 1\})$  corresponding to a uniform discretization of  $\mathbb{T}$  with  $N_p + 1$  points. The environment (whose inner working is not known to the learning agent) needs to compute the evolution of the distribution. This evolution can be directly simulated with a deterministic method based, for example, on a finite difference scheme as in [2]. However, in practice, it is likely that the environment would not work in this way but would rather correspond to moving forward a large population of agents (e.g., robots). This induces extra approximations. To illustrate that our method can be applied in such situations, here we chose to implement the environment using a probabilistic approach based on Monte Carlo simulations for a large number of particles on  $S$ . Then, to prepare the input for the Q-function, we project their positions on  $\{0, 1/N_p, \dots, 1 - 1/N_p, 1\}$  and approximate the mean field distribution by a histogram. We recall that the DDPG method uses this environment as a black-box and, for a given action  $\tilde{a} \in \mathbb{R}^{N_p}$ , can only access the resulting new distribution and the associated reward. The actor and critic networks have been implemented using a feedforward fully connected architecture with 2 hidden layers of width at most 300 neurons. We used random initial states at each episode, and the noise used on the action is a Gaussian noise with mean 0 and variance 0.02. We used Adam optimizer with initial learning rate 0.0001 and minibatches of size 16.

Figure 8 presents results obtained using this method after 160 episodes. The system has been trained on initial distributions which are Gaussian with random mean and random variance. As illustrated in the figures (left column), the system has learnt how to drive this type of initial distributions towards the analytical stationary distribution and then how to use an approximation of the stationary optimal control (right column) in order to keep the system in the stationary regime. The middle column displays the learnt control for the initial distribution. It is not expected to match the optimal ergodic control, which should be applied when the distribution has reached the ergodic regime and here it is provided only for the sake of comparison.

## REFERENCES

- [1] Achdou, Y., Camilli, F., and Capuzzo-Dolcetta, I. (2012). Mean field games: numerical methods for the planning problem. *SIAM J. Control Optim.*, 50(1):77–109.
- [2] Achdou, Y. and Capuzzo-Dolcetta, I. (2010). Mean field games: numerical methods. *SIAM J. Numer. Anal.*, 48(3):1136–1162.
- [3] Achdou, Y. and Laurière, M. (2016). Mean Field Type Control with Congestion (II): An augmented Lagrangian method. *Appl. Math. Optim.*, 74(3):535–578.
- [4] Agram, N., Bakdi, A., and Oksendal, B. (2020). Deep learning and stochastic mean-field control for a neural network model. *Available at SSRN 3639022*.

- [5] Al-Aradi, A., Correia, A., Naiff, D. d. F., Jardim, G., and Saporito, Y. (2019). Applications of the deep galerkin method to solving partial integro-differential and hamilton-jacobi-bellman equations. *arXiv preprint arXiv:1912.01455*.
- [6] Almulla, N., Ferreira, R., and Gomes, D. (2017). Two numerical approaches to stationary mean-field games. *Dyn. Games Appl.*, 7(4):657–682.
- [7] Anahtarci, B., Kariksiz, C. D., and Saldi, N. (2020). Q-learning in regularized mean-field games. *arXiv preprint arXiv:2003.12151*.
- [8] Bellemare, M. G., Ostrovski, G., Guez, A., Thomas, P. S., and Munos, R. (2016). Increasing the action gap: New operators for reinforcement learning. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- [9] Bensoussan, A., Frehse, J., and Yam, S. C. P. (2013). *Mean field games and mean field type control theory*. Springer Briefs in Mathematics. Springer, New York.
- [10] Bertsekas, D. P. and Shreve, S. (2004). *Stochastic optimal control: the discrete-time case*.
- [11] Cardaliaguet, P., Delarue, F., Lasry, J.-M., and Lions, P.-L. (2019). *The master equation and the convergence problem in mean field games*, volume 201 of *Annals of Mathematics Studies*. Princeton University Press, Princeton, NJ.
- [12] Carmona, R. and Delarue, F. (2018a). *Probabilistic theory of mean field games with applications. I*, volume 83 of *Probability Theory and Stochastic Modelling*. Springer, Cham. Mean field FBSDEs, control, and games.
- [13] Carmona, R. and Delarue, F. (2018b). *Probabilistic theory of mean field games with applications. II*, volume 84 of *Probability Theory and Stochastic Modelling*. Springer, Cham. Mean field games with common noise and master equations.
- [14] Carmona, R., Hamidouche, K., Laurière, M., and Tan, Z. (2020). Policy optimization for linear-quadratic zero-sum mean-field type games. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 1038–1043. IEEE.
- [15] Carmona, R. and Laurière, M. (2019). Convergence analysis of machine learning algorithms for the numerical solution of mean field control and games: Ii—the finite horizon case. *arXiv preprint arXiv:1908.01613*. To appear in *Annals of Probability*.
- [16] Carmona, R. and Laurière, M. (2021). Convergence analysis of machine learning algorithms for the numerical solution of mean field control and games i: The ergodic case. *SIAM Journal on Numerical Analysis*, 59(3):1455–1485.
- [17] Chassagneux, J.-F., Crisan, D., and Delarue, F. (2014). A probabilistic approach to classical solutions of the master equation for large population equilibria. *arXiv:1411.3009*.
- [18] Elie, R., Perolat, J., Laurière, M., Geist, M., and Pietquin, O. (2020). On the convergence of model free learning in mean field games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7143–7150.
- [19] Even-Dar, E. and Mansour, Y. (2003). Learning rates for Q-learning. *J. Mach. Learn. Res.*, 5:1–25.
- [20] Farahmand, A.-m. (2011). Action-gap phenomenon in reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 172–180.
- [21] Fouque, J.-P. and Zhang, Z. (2020). Deep learning methods for mean field control problems with delay. *Frontiers in Applied Mathematics and Statistics*, 6:11.
- [22] Fu, Z., Yang, Z., Chen, Y., and Wang, Z. (2019). Actor-critic provably finds nash equilibria of linear-quadratic mean-field games. In *International Conference on Learning Representations*.
- [23] Gao, B. and Pavel, L. (2017). On the properties of the softmax function with application in game theory and reinforcement learning. *arXiv preprint arXiv:1704.00805*.
- [24] Germain, M., Mikael, J., and Warin, X. (2019). Numerical resolution of mckean-vlasov fbsdes using neural networks. *arXiv preprint arXiv:1909.12678*.
- [25] Gu, H., Guo, X., Wei, X., and Xu, R. (2019). Dynamic programming principles for mean-field controls with learning. *arXiv preprint arXiv:1911.07314*.
- [26] Gu, H., Guo, X., Wei, X., and Xu, R. (2020). Mean-field controls with q-learning for cooperative marl: Convergence and complexity analysis. *arXiv preprint arXiv:2002.04131*.

- [27] Guo, X., Hu, A., Xu, R., and Zhang, J. (2019). Learning mean-field games. *Advances in Neural Information Processing Systems*, 32:4966–4976.
- [28] Huang, M., Malhamé, R. P., and Caines, P. E. (2006). Large population stochastic dynamic games: closed-loop McKean-Vlasov systems and the Nash certainty equivalence principle. *Commun. Inf. Syst.*, 6(3):221–251.
- [29] Kallenberg, O. (2017). *Random measures, theory and applications*. Springer.
- [30] Kolokoltsov, V. N. and Bensoussan, A. (2016). Mean-field-game model for botnet defense in cyber-security. *Appl. Math. Optim.*, 74(3):669–692.
- [31] Lasry, J.-M. and Lions, P.-L. (2007). Mean field games. *Jpn. J. Math.*, 2(1):229–260.
- [32] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2016). Continuous control with deep reinforcement learning. In *Proceedings of the International Conference on Learning Representations (ICLR 2016)*.
- [33] Motte, M. and Pham, H. (2019). Mean-field markov decision processes with common noise and open-loop controls. *arXiv preprint arXiv:1912.07883*.
- [34] Perrin, S., Pérolat, J., Laurière, M., Geist, M., Elie, R., and Pietquin, O. (2020). Fictitious play for mean field games: Continuous time analysis and applications. *Advances in Neural Information Processing Systems*.
- [35] Ruthotto, L., Osher, S. J., Li, W., Nurbekyan, L., and Fung, S. W. (2020). A machine learning framework for solving high-dimensional mean field game and mean field control problems. *Proceedings of the National Academy of Sciences*, 117(17):9183–9193.
- [36] Subramanian, J. and Mahajan, A. (2019). Reinforcement learning in stationary mean-field games. In *Proceedings. 18th International Conference on Autonomous Agents and Multiagent Systems*.

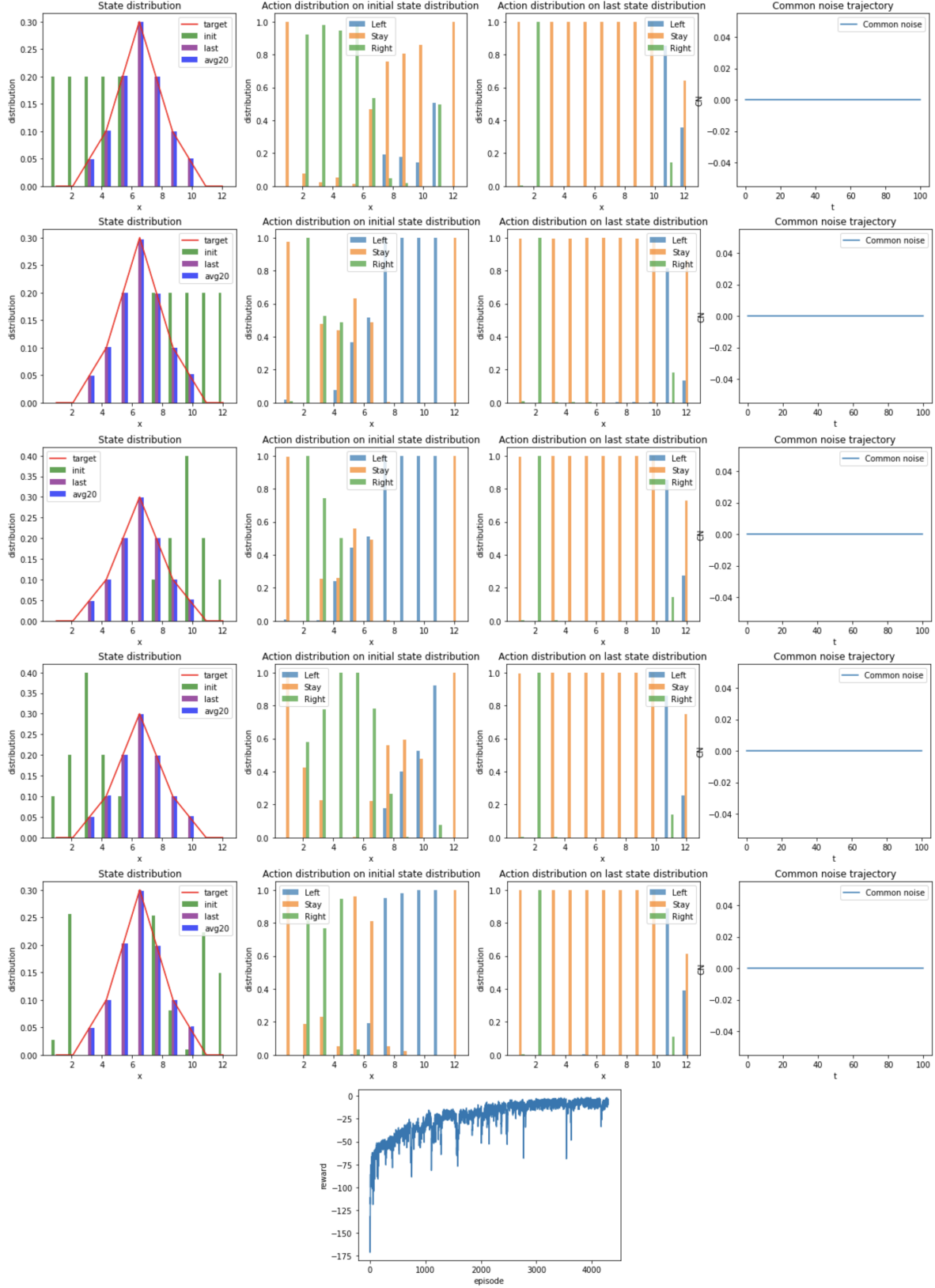


FIGURE 6. Example 2: Discrete distribution planning. Case without common noise; 5 testing distributions. Column 1: state distribution; columns 2 and 3: action distribution at initial and terminal time; column 4: common noise trajectory (identically 0 here). Bottom: evolution of the reward during training.

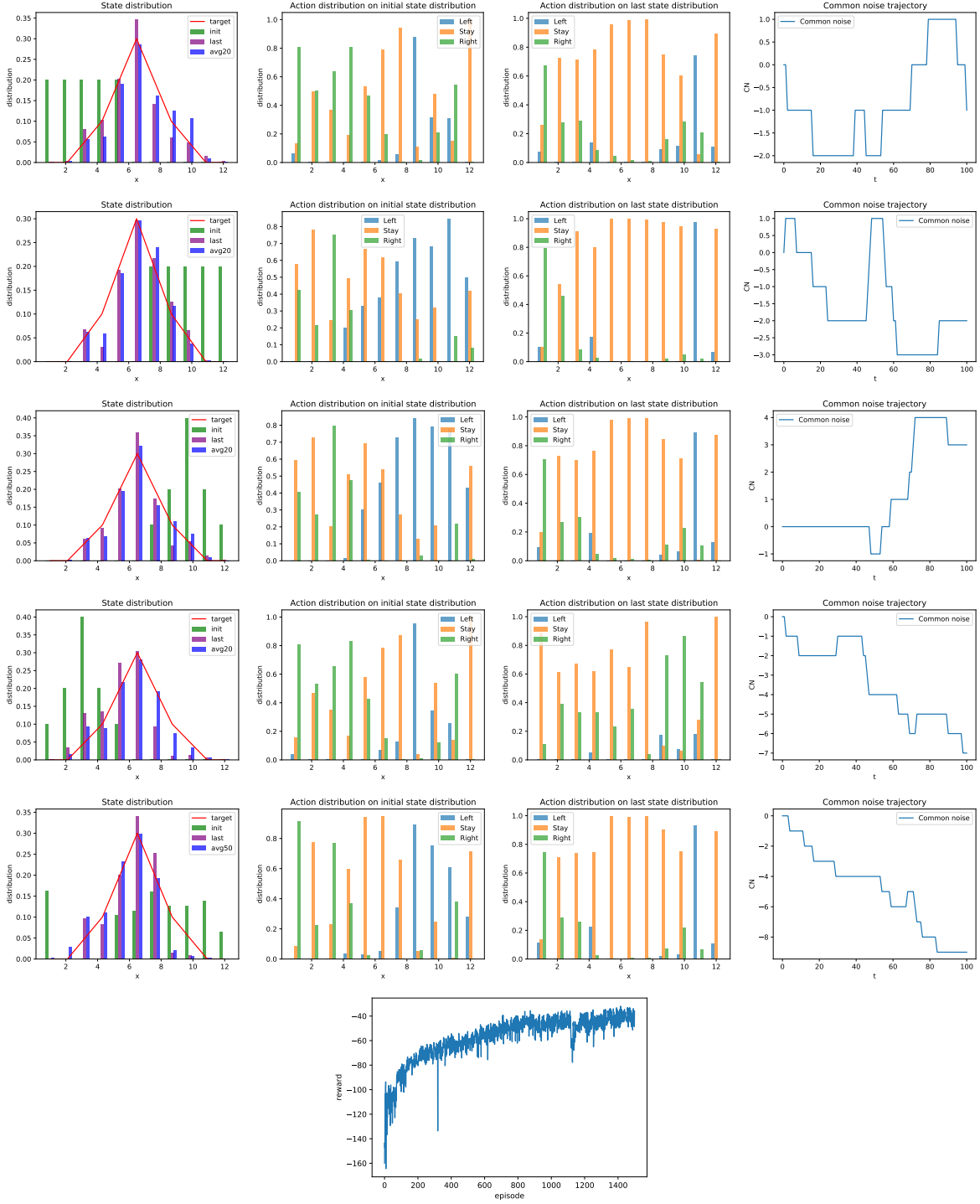


FIGURE 7. Example 2: Discrete distribution planning. Case with common noise; 5 testing distributions. Column 1: state distribution; columns 2 and 3: action distribution at initial and terminal time; column 4: common noise trajectory (identically 0 here). Bottom: evolution of the reward during training.

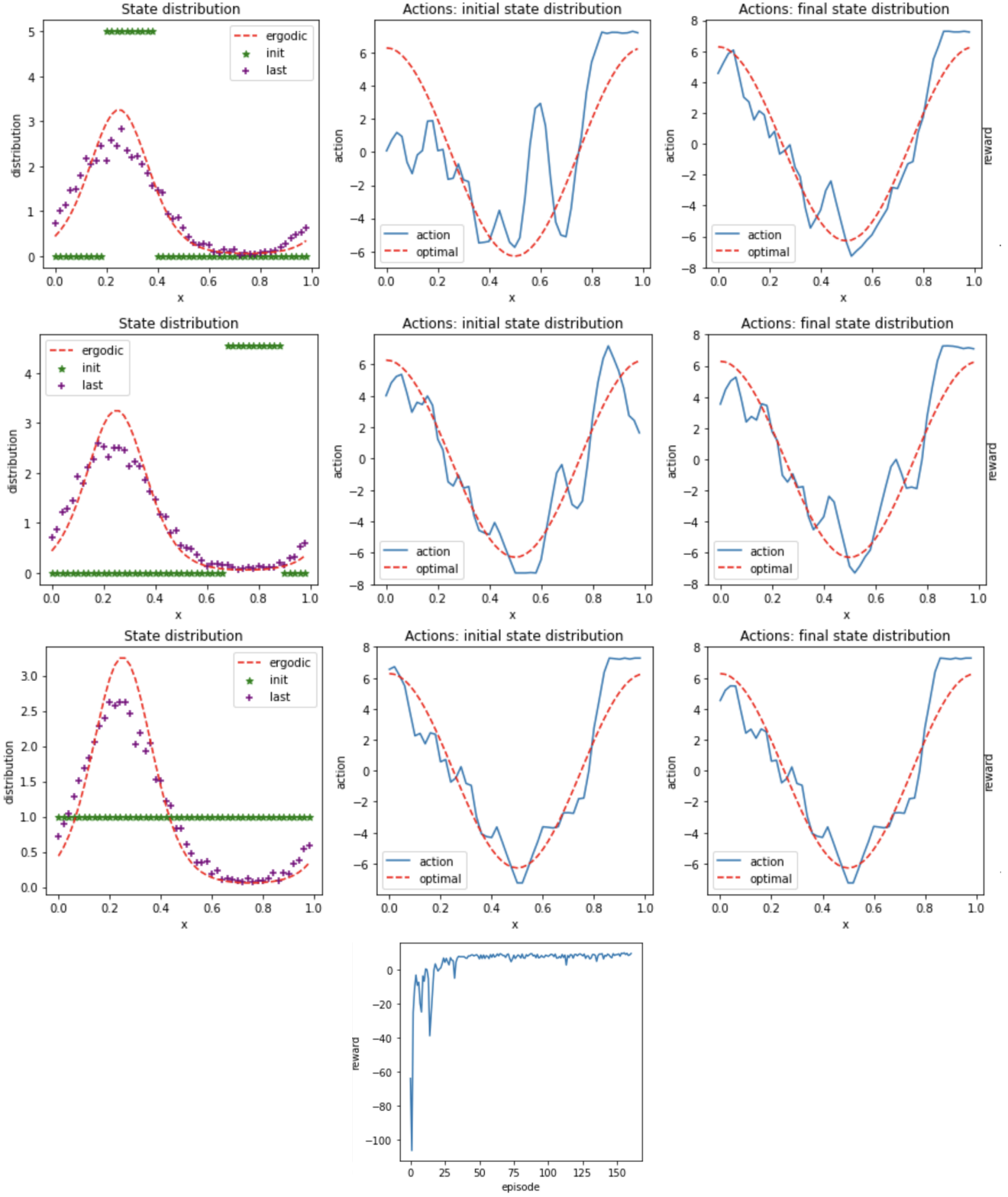


FIGURE 8. Swarm motion: Left: Distribution induced, middle: control learnt for distribution at time 0, right: control learnt for distribution at terminal time. For three different initial distributions. Red: ergodic analytical solution; green: initial distribution; purple: distribution at terminal time. Bottom: rewards.

## APPENDIX A. AUXILIARY RESULTS FOR SECTION 2

**Lemma 39.** *Given a level-0 control process  $\mathbf{a}$ , there exists a level-0 action process  $\alpha$  which is a realization of  $\mathbf{a}$ . Moreover, if another level-0 action process  $\alpha'$  is also a realization of  $\mathbf{a}$ , then for every  $n \geq 0$  and for every bounded Borel measurable function  $h : A \rightarrow \mathbb{R}$ , we have*

$$(37) \quad \mathbb{E}[h(\alpha'_n) | \mathcal{G}_n^c] = \int_A h(\alpha) \mathbf{a}_n(d\alpha) = \mathbb{E}[h(\alpha_n) | \mathcal{G}_n^c], \quad \mathbb{P} - a.s..$$

*Proof.* By definition, for each  $n \geq 0$ ,  $\mathbf{a}_n$  is of the form:

$$\mathbf{a}_n(d\alpha) = \kappa_n^{\mathbf{a}}(\mathcal{U}, \underline{\vartheta}_{n-1}, \underline{\vartheta}_n^0, \underline{\varepsilon}_n, \underline{\varepsilon}_n^0)(d\alpha), \quad \mathbb{P} - a.s.,$$

for some measurable function  $\kappa_n^{\mathbf{a}}$  on  $\Upsilon \times \Theta^{n-1} \times (\Theta^0)^n \times E^n \times (E^0)^n$  with values in  $\mathcal{P}(A)$ . Let us denote by  $\xi_n$  the random element  $(\mathcal{U}, \underline{\vartheta}_{n-1}, \underline{\vartheta}_n^0, \underline{\varepsilon}_n, \underline{\varepsilon}_n^0)$ , by  $\rho_A$  the Blackwell-Dubins function (see Lemma 2) of the space  $A$ . Let us set  $U_n = h(\vartheta_n)$  and:

$$\alpha_n(\omega) := \rho_A(\kappa_n^{\mathbf{a}}(\xi_n(\omega)), U_n(\omega)).$$

Then,  $\alpha_n$  is  $\mathcal{G}_n^{\mathbf{a}}$ -measurable, and because the  $\sigma$ -field  $\mathcal{G}_n^c = \sigma\{\xi_n\}$  is independent of  $U_n$ , we have  $\mathcal{L}(\alpha_n | \mathcal{G}_n^c) = \kappa_n^{\mathbf{a}}(\xi_n) = \mathbf{a}_n$ ,  $\mathbb{P}$ -almost surely. Equality (37) directly follows the definition of a conditional distribution.  $\square$

**Lemma 40.** *For each  $n \geq 0$ , the law of the random measure  $\mathbb{P}_{(X_n, \alpha_n)}^0$  depends only upon the open-loop policy  $\pi$  as long as  $\alpha$  is a realization of the control process generated by  $\pi$ .*

*Proof.* Let  $\alpha$  be an action process which is a realization of the control process generated by  $\pi$ . For each integer  $n \geq 0$ , we compute  $\mathbb{E}[\Phi(\mathbb{P}_{(X_n, \alpha_n)}^0)]$  for a family of bounded measurable functions  $\Phi$  on  $\mathcal{P}(S \times A)$  which generate the Borel  $\sigma$ -field of  $\mathcal{P}(S \times A)$ . For the sake of definiteness we work with functions  $\Phi$  of the form:

$$\Phi(\mu) = \prod_{j=1}^m \int_{S \times A} \varphi_j(x, \alpha) \mu(dx, d\alpha)$$

for a finite set  $\varphi_1, \dots, \varphi_m$  of bounded continuous functions on  $S \times A$ . We have:

$$(38) \quad \begin{aligned} \mathbb{E}[\Phi(\mathbb{P}_{(X_n, \alpha_n)}^0)] &= \mathbb{E}\left[\prod_{j=1}^m \int_{S \times A} \varphi_j(x, \alpha) \mathbb{P}_{(X_n, \alpha_n)}^0(dx, d\alpha)\right] \\ &= \mathbb{E}\left[\prod_{j=1}^m \mathbb{E}[\varphi_j(X_n, \alpha_n) | \sigma\{\underline{\vartheta}_n^0, \underline{\varepsilon}_n^0\}]\right]. \end{aligned}$$

Now for each  $j \in \{1, \dots, m\}$  we have

$$(39) \quad \begin{aligned} &\mathbb{E}[\varphi_j(X_n, \alpha_n) | \sigma\{\underline{\vartheta}_n^0, \underline{\varepsilon}_n^0\}] \\ &= \int \cdots \int \varphi_j(X_n(u, \underline{\vartheta}_{n-1}, \underline{\varepsilon}_n, \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0), \alpha_n(u, \underline{\vartheta}_{n-1}, \underline{\varepsilon}_n, \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0, \theta_n)) \\ &\quad \mathbb{P}_{\mathcal{U}}(du) \mathbb{P}_{\underline{\vartheta}_{n-1}}(d\underline{\vartheta}_{n-1}) \nu^n(d\underline{\varepsilon}_n) \mathbb{P}_{\vartheta_n}(d\theta_n) \Big|_{\underline{\vartheta}_n^0 = \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0 = \underline{\varepsilon}_n^0} \\ &= \int \cdots \int \left( \int_A \varphi_j(X_n(u, \underline{\vartheta}_{n-1}, \underline{\varepsilon}_n, \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0), \alpha) \pi_n(d\alpha | u, \underline{\vartheta}_{n-1}, \underline{\varepsilon}_n, \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0) \right) \\ &\quad \mathbb{P}_{\mathcal{U}}(du) \mathbb{P}_{\underline{\vartheta}_{n-1}}(d\underline{\vartheta}_{n-1}) \nu^n(d\underline{\varepsilon}_n) \Big|_{\underline{\vartheta}_n^0 = \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0 = \underline{\varepsilon}_n^0} \end{aligned}$$

where we made explicit the dependence of  $X_n$  on  $\xi_n = (\mathcal{U}, \underline{\vartheta}_{n-1}, \underline{\varepsilon}_n, \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0)$  and  $\alpha_n$  on  $(\xi_n, \vartheta_n)$ . This shows that the left hand side of (39), and hence the left hand side of (38) only depend upon the action process  $\alpha = (\alpha_n)_{n \geq 0}$  through the conditional distribution  $\pi_n(d\alpha | u, \underline{\vartheta}_{n-1}, \underline{\varepsilon}_n, \underline{\vartheta}_n^0, \underline{\varepsilon}_n^0)$ . From this we conclude that if two action processes are realizations of control processes generated by the same the open-loop policy  $\pi$ , the corresponding random measures  $\mathbb{P}_{(X_n, \alpha_n)}^0$  have the same distribution.  $\square$

**Remark 41.** For every  $(\mu, \bar{a}, e^0) \in \bar{\Gamma} \times E^0$ ,  $\bar{F}(\mu, \bar{a}, e^0)$  is defined as a probability measure on  $S$  such that for every bounded and Borel measurable function  $\phi : S \rightarrow \mathbb{R}$ ,

$$(40) \quad \int_S \bar{F}(\mu, \bar{a}, e^0)(dx') \phi(x') = \int_{S \times A \times E} \bar{a}(dx, d\alpha) \nu(de) \phi(F(x, \alpha, \bar{a}, e, e^0)).$$

It is straightforward to check that  $\bar{F}$  is Borel measurable. See for example [10, Proposition 7.29] for a proof.

**Lemma 42.** Assume **(H1)**. It holds:

- $\bar{F}$  is Borel measurable and for every  $e^0 \in E^0$ ,  $\bar{F}(\cdot, \cdot, e^0)$  is continuous in its remaining variables.
- $\bar{f}$  is bounded and lower semi-continuous.

*Proof.* The measurability of  $\bar{F}$  was argued earlier (see after Definition 12), so we only argue the continuity for  $e^0 \in E^0$  fixed. Let  $((\mu_n, \bar{a}_n))_{n \geq 0}$  be a sequence in  $\bar{\Gamma}$  which converges weakly toward  $(\mu, \bar{a})$ . Since  $\mu_n = \text{pr}_1(\bar{a}_n)$  for each  $n \geq 0$ , we necessarily have  $\mu = \text{pr}_1(\bar{a})$  so that  $(\mu, \bar{a}) \in \bar{\Gamma}$ . We pick a continuous bounded function  $\phi : S \mapsto \mathbb{R}$  and we show that:

$$(41) \quad \lim_{n \rightarrow \infty} \int_S \phi(x') \bar{F}(\mu_n, \bar{a}_n, e^0)(dx') = \int_S \phi(x') \bar{F}(\mu, \bar{a}, e^0)(dx').$$

Using Skorohod's characterization of weak convergence of probability measures, we have the existence of random variables  $(Y_n, \beta_n)$  converging  $\mathbb{P}$ -almost surely toward some  $(Y, \beta)$  and such that  $\mathbb{P}_{(Y_n, \beta_n)} = \bar{a}_n$  for each  $n \geq 0$  and  $\mathbb{P}_{(Y, \beta)} = \bar{a}$ . Consequently, the integral in the left hand side of (41) can be rewritten as:

$$\begin{aligned} \int_S \phi(x') \bar{F}(\mu_n, \bar{a}_n, e^0)(dx') &= \int_{S \times A \times E} \phi(F(x, \alpha, \bar{a}_n, e, e^0)) \bar{a}_n(dx, d\alpha) \nu(de) \\ &= \int_E \left( \mathbb{E}[\phi(F(X_n, \alpha_n, \bar{a}_n, e, e^0))] \right) \nu(de), \end{aligned}$$

which converges toward  $\int_E \left( \mathbb{E}[\phi(F(X_n, \alpha_n, \bar{a}_n, e, e^0))] \right) \nu(de)$  by Lebesgue's dominated convergence theorem because  $F(\cdot, \cdot, \cdot, e, e^0)$  is continuous and  $\phi$  is bounded continuous.

The lower semi-continuity of the one stage cost function  $\bar{f}$  follows from [10, Proposition 7.31 (a)] which only requires the lower semi-continuity of the original one-stage cost function  $f$  instead of the full continuity assumption posited in Assumption **(H1)**.  $\square$

**Lemma 43.** Let  $\bar{\pi} \in \bar{\Pi}$ . For every  $\mu \in \bar{S}$ , let  $(\mu, \bar{a})$  and  $(\mu', \bar{a}')$  be two pairs of state and action processes generated by  $(\bar{\pi}, \mu)$ . Then  $\mathbb{E}[\sum_{n \geq 0} \gamma^n \bar{f}(\mu_n, \bar{a}_n)] = \mathbb{E}[\sum_{n \geq 0} \gamma^n \bar{f}(\mu'_n, \bar{a}'_n)]$ .

*Proof.* For any fixed initial  $\mu \in \bar{S}$ , by definition of the pair of state and action processes generated by  $(\bar{\pi}, \mu)$ , we show by induction that, for all  $n \geq 0$ ,  $\mathcal{L}(\mu_n) = \mathcal{L}(\mu'_n)$  and  $\mathcal{L}((\mu_n, \bar{a}_n)) = \mathcal{L}((\mu'_n, \bar{a}'_n))$ . For  $n = 0$ ,  $\mu_0 = \mu'_0 = \mu \in \bar{S}$ . Assume that for some  $n \geq 0$ , we have  $\mathcal{L}(\mu_n) = \mathcal{L}(\mu'_n)$ , then for every bounded and Borel measurable function  $\phi : \bar{S} \times \bar{A} \rightarrow \mathbb{R}$

$$\begin{aligned} \mathbb{E}[\phi(\mu_n, \bar{a}_n)] &= \mathbb{E}[\mathbb{E}[\phi(\mu_n, \bar{a}_n) | \mu_n]] \\ &= \mathbb{E} \left[ \int_{\bar{A}} \phi(\mu_n, \bar{a}) \cdot \mathcal{L}(\bar{a}_n | \mu_n)(d\bar{a}) \right] \\ (42) \quad &= \mathbb{E} \left[ \int_{\bar{A}} \phi(\mu_n, \bar{a}) \cdot \bar{\pi}_n(\mu_n)(d\bar{a}) \right] \\ &= \mathbb{E} \left[ \int_{\bar{A}} \phi(\mu'_n, \bar{a}) \cdot \bar{\pi}_n(\mu'_n)(d\bar{a}) \right] = \mathbb{E}[\mathbb{E}[\phi(\mu'_n, \bar{a}'_n) | \mu'_n]] = \mathbb{E}[\phi(\mu'_n, \bar{a}'_n)]. \end{aligned}$$

So  $\mathcal{L}((\mu_n, \bar{a}_n)) = \mathcal{L}((\mu'_n, \bar{a}'_n))$ . Since  $\varepsilon_{n+1}^0$  is independent of  $(\mu_n, \bar{a}_n)$  and  $(\mu'_n, \bar{a}'_n)$ ,  $\mathcal{L}((\mu_n, \bar{a}_n, \varepsilon_{n+1}^0)) = \mathcal{L}((\mu'_n, \bar{a}'_n, \varepsilon_{n+1}^0))$ , which implies that the law of  $\mu_{n+1} = \bar{F}(\mu_n, \bar{a}_n, \varepsilon_{n+1}^0)$  is equal to the law of  $\bar{F}(\mu'_n, \bar{a}'_n, \varepsilon_{n+1}^0) = \mu'_{n+1}$ . Hence the conclusion.  $\square$

## APPENDIX B. PROOFS FOR SECTION 4.1

*Proof of Lemma 23.* Let us denote by  $\zeta_{n+1}$  the right hand side of (18), and let  $\phi : S \rightarrow \mathbb{R}$ ,  $h_n : (\Theta^0 \times E^0)^n \rightarrow \mathbb{R}$  and  $\psi_{n+1} : E^0 \rightarrow \mathbb{R}$  be arbitrary bounded Borel measurable functions. We have:

$$\begin{aligned}
& \mathbb{E} \left[ \psi_{n+1}(\varepsilon_{n+1}^0) h_n(\vartheta_n^0, \varepsilon_n^0) \int_S \phi(x') \zeta_{n+1}(dx') \right] \\
&= \mathbb{E} \left[ \psi_{n+1}(\varepsilon_{n+1}^0) h_n(\vartheta_n^0, \varepsilon_n^0) \int_S \phi(x') \left( \bar{F}(\mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0, \varepsilon_{n+1}^0) \right) (dx') \right] \\
&= \mathbb{E} \left[ \psi_{n+1}(\varepsilon_{n+1}^0) h_n(\vartheta_n^0, \varepsilon_n^0) \int_{S \times A \times E} \mathbb{P}_{(X_n, \alpha_n)}^0(dx, d\alpha) \nu(de) \phi \left( F(x, \alpha, \mathbb{P}_{(X_n, \alpha_n)}^0, e, \varepsilon_{n+1}^0) \right) \right] \\
&= \int_{E \times E^0} \nu(de) \nu^0(de^0) \psi_{n+1}(e^0) \mathbb{E} \left[ h_n(\vartheta_n^0, \varepsilon_n^0) \int_{S \times A} \mathbb{P}_{(X_n, \alpha_n)}^0(dx, d\alpha) \phi \left( F(x, \alpha, \mathbb{P}_{(X_n, \alpha_n)}^0, e, e^0) \right) \right] \\
&= \mathbb{E} \left[ \psi_{n+1}(\varepsilon_{n+1}^0) \mathbb{E} \left[ h_n(\vartheta_n^0, \varepsilon_n^0) \phi \left( F(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}^0, \varepsilon_{n+1}, \varepsilon_{n+1}^0) \right) \middle| \mathcal{F}_n^0, \varepsilon_{n+1}, \varepsilon_{n+1}^0 \right] \right] \\
&= \mathbb{E} \left[ \psi_{n+1}(\varepsilon_{n+1}^0) h_n(\vartheta_n^0, \varepsilon_n^0) \phi(X_{n+1}) \right],
\end{aligned}$$

where the first equality is by definition of  $\zeta_{n+1}$ , the second equality is by the definition of  $\bar{F}$  in terms of the system function  $F$  of the original MFC, the third equality is by the fact that  $\varepsilon_{n+1}^0$  is independent of all the other random quantities, the fourth equality is by definition of the conditional probability  $\mathbb{P}_{(X_n, \alpha_n)}^0$ , and the last equality is by the tower property of conditional expectation, the fact that  $X_{n+1} = F(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}^0, \varepsilon_{n+1}, \varepsilon_{n+1}^0)$  and the fact that  $(\varepsilon_{n+1}, \varepsilon_{n+1}^0)$  is independent of  $(X_n, \alpha_n)$  and  $\mathbb{P}_{(X_n, \alpha_n)}^0$  is measurable with respect to  $\mathcal{F}_n^0 = \sigma\{\vartheta_n^0, \varepsilon_n^0\}$ . This shows that  $\zeta_{n+1} = \mathbb{P}_{X_{n+1}}^0$ .  $\square$

*Proof of Lemma 24.* We prove equation (19) by induction. Equality (19) holds for  $n = 0$  by the assumptions on  $(\zeta, \bar{\eta})$ . For the sake of an argument by induction, let us assume that (19) holds for some  $n \geq 0$ . We first show that  $\mathcal{L}(\zeta_{n+1}) = \mathcal{L}(\mathbb{P}_{X_{n+1}}^0)$ . Since  $\varepsilon_{n+1}^0$  is independent of  $\mathcal{F}_n^0$ , for every bounded Borel measurable function  $\psi : \bar{S} \rightarrow \mathbb{R}$ , it holds:

$$\begin{aligned}
\mathbb{E}[\psi(\zeta_{n+1})] &= \mathbb{E} \left[ \mathbb{E} \left[ \psi \left( \bar{F}(\zeta_n, \bar{\eta}_n, \varepsilon_{n+1}^0) \right) \middle| \zeta_n, \bar{\eta}_n \right] \right] \\
&= \mathbb{E} \left[ \int_{\bar{S}} \psi(\mu) P(\zeta_n, \bar{\eta}_n)(d\mu) \right] \\
&= \mathbb{E} \left[ \int_{\bar{S}} \psi(\mu) P \left( \mathbb{P}_{X_n}^0, \mathbb{P}_{(X_n, \alpha_n)}^0 \right) (d\mu) \right] \\
&= \mathbb{E} \left[ \psi(\mathbb{P}_{X_{n+1}}^0) \right],
\end{aligned}$$

where the second equality is by definition of  $P$  in (9), the third equality is by the induction hypothesis, and the last equality is due to Lemma 23. So  $\mathcal{L}(\zeta_{n+1}) = \mathcal{L}(\mathbb{P}_{X_{n+1}}^0)$ . We then consider the joint law of  $(\zeta_{n+1}, \bar{\eta}_{n+1})$ . By the assumption that  $\mathcal{L}(\bar{\eta}_n | \zeta_n) = \kappa_n(\zeta_n)$ , we have that  $(\zeta_{n+1}, \bar{\eta}_{n+1})$  and  $(\mathbb{P}_{X_{n+1}}^0, \mathbb{P}_{(X_{n+1}, \alpha_{n+1})}^0)$  share the same regular version  $\kappa_n$  of the conditional probability. We conclude that (19) holds for  $n + 1$  instead of  $n$ .  $\square$

## APPENDIX C. PROOFS FOR SECTION 4.2

*Proof of Lemma 28.* We prove this statement by showing that there exist  $\pi \in \Pi^{OL}$  and an open-loop action process  $\alpha$  generated by  $\pi$  such that:  $J^\pi = J^\alpha = J^{\tilde{\pi}}$ . Let  $\mu_0 \in \mathcal{P}(S)$  and let  $(\mathbf{X}, \alpha)$  be a pair of state and action processes generated by  $(\tilde{\pi}, \mu_0)$ . Let  $\mathbf{a}$  be the  $\mathcal{P}(A)$ -valued process given by:

$$\mathbf{a}_n = \tilde{\pi}_n(X_n, \mathbb{P}_{X_n}^0, \vartheta_n^0), \quad n \geq 0.$$

We recall that  $\Xi_n, \xi_n$  and  $\mathcal{U}$  are defined in § 2.2.3. Since  $\mathbf{a}$  is adapted to  $\mathbb{G}^c$ , it is an admissible level-0 control process, and for every  $n \geq 0$ , there exists a Borel measurable function  $\pi_n : \Xi_n \rightarrow A$  satisfying

$$\mathbf{a}_n = \pi_n(\mathcal{U}, (\vartheta_k, \vartheta_k^0, \varepsilon_{k+1}, \vartheta_{k+1}^0)_{k=0, \dots, n-1}, \vartheta_n^0) = \pi_n(\xi_n), \quad \mathbb{P} - a.s..$$

Moreover, since  $\alpha$  is generated by  $\tilde{\pi}$ , it is adapted to  $\mathbb{G}^a$  and satisfies

$$\mathcal{L}(\alpha_n | \mathcal{G}_n^c) = \tilde{\pi}_n(X_n, \mathbb{P}_{X_n}^0, \vartheta_n^0) = \pi_n(\xi_n), \quad \mathbb{P} - a.s. \quad n \geq 0.$$

So  $\alpha$  can be viewed as an open-loop action process generated by  $\pi$ . Meanwhile, the state process  $\mathbf{X}$  constructed by equation (1) is also a state process associated with  $(\alpha, \mu_0)$  (see Definition 5). Therefore, by definition of the value function associated to an open-loop policy  $\pi$ , we have:

$$J^{\tilde{\pi}}(\mu_0) = J^\alpha(\mu_0) = \mathbb{E} \left[ \sum_{n \geq 0} \gamma^n f(X_n, \alpha_n, \mathbb{P}_{(X_n, \alpha_n)}^0) \right].$$

□

## APPENDIX D. DISINTEGRATION OF KERNELS

Given a probability measure  $P$  on a measurable space  $(C, \mathcal{C})$  and a kernel  $K$  from  $(C, \mathcal{C})$  to  $(D, \mathcal{D})$ , the composition of measure  $P$  and kernel  $K$ , denoted by  $P \hat{\otimes} K$ , is defined as a measure on the product space  $(C \times D, \mathcal{C} \otimes \mathcal{D})$  such that for every non-negative measurable function  $f : C \times D \rightarrow \mathbb{R}_+$ ,

$$(P \hat{\otimes} K)f = \int_C P(dx) \int_D f(x, y) K(x, dy).$$

Similarly, for a probability kernel  $\mu : G \rightarrow \mathcal{P}(C)$  and a probability kernel  $K : G \times C \rightarrow \mathcal{P}(D)$ , the composition of kernels  $\mu$  and  $K$ , denoted by  $\mu \hat{\otimes} K$ , is defined as a kernel from  $(G, \mathcal{G})$  to  $(C \times D, \mathcal{C} \otimes \mathcal{D})$  such that for every  $s \in G$  and for every non-negative measurable function  $f : C \times D \rightarrow \mathbb{R}_+$ ,

$$(\mu \hat{\otimes} K)(s)f = \int_C \mu(s, dx) \int_D f(x, y) K(s, x, dy).$$

## APPENDIX E. DETAILS ON Q-LEARNING, SECTION 5

## E.1. Proof of Theorem 36.

*Proof of Theorem 36.* Recall that we denote by  $\check{J}^*$  and  $\check{Q}^*$  respectively the state value function and the state-action value function of the projected MFC problem defined by (32)–(33).

We first note that, for every  $(\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}$ ,

$$\begin{aligned} \left| \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a}) - \check{Q}^*(\mu, \tilde{a}) \right| &\leq \left| \check{Q}_{N_{\text{epi}}}(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a}) - \check{Q}^*(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a}) \right| \\ &\quad + \left| \check{Q}^*(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a}) - \tilde{Q}^*(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a}) \right| \\ &\quad + \left| \tilde{Q}^*(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a}) - \tilde{Q}^*(\mu, \tilde{a}) \right|. \end{aligned}$$

We then split the proof into three steps, which consist in bounding from above each term in the right hand side.

**Step 1.** We first analyze the difference between  $\check{Q}_{N_{\text{epi}}}$  and  $\check{Q}^*$ . This comes from standard convergence results on Q-learning for finite state-action spaces. More precisely, under Assumptions **(H3)** and **(H5)**, with our choice

of learning rates, and given that  $N_{\text{epi}}$  is of order (34), we can apply Theorem 4 and Corollary 34 in [19] for asynchronous Q-learning and polynomial learning rates, and we obtain that, with probability at least  $1 - \delta$ ,

$$\|\check{Q}_{N_{\text{epi}}} - \check{Q}^*\|_\infty = \sup_{(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \mathcal{A}} \left| \check{Q}_{N_{\text{epi}}}(\check{\mu}, \check{a}) - \check{Q}^*(\check{\mu}, \check{a}) \right| \leq \varepsilon.$$

**Step 2.** We then turn our attention to the difference between  $\check{Q}^*$  and  $\tilde{Q}^*$ . The analysis amounts to say that the projection on  $\check{\mathfrak{S}}$  realized at each step does not perturb too much the value function. Recall that for some given common noise  $\varepsilon^0$ , the operator  $\check{\Phi}^{\varepsilon^0} : \check{\mathfrak{S}} \times \mathcal{A} \rightarrow \check{\mathfrak{S}}$  is given by  $\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) = \text{Proj}_{\check{\mathfrak{S}}} \circ \bar{F}(\check{\mu}, \check{\mu} \otimes \check{a}, \varepsilon^0)$ . Likewise, we denote the transition dynamic with  $\bar{F}$  by a function  $\Phi^{\varepsilon^0} : \mathfrak{S} \times \mathcal{A} \rightarrow \mathfrak{S}$  such that:

$$\Phi^{\varepsilon^0}(\mu, \tilde{a}) = \bar{F}(\mu, \mu \otimes \tilde{a}, \varepsilon^0), \quad \forall (\mu, \tilde{a}) \in \mathfrak{S} \times \mathcal{A}.$$

Let us start by noting that, for every  $(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \mathcal{A}$ ,

$$\begin{aligned} & \left| \check{Q}^*(\check{\mu}, \check{a}) - \tilde{Q}^*(\check{\mu}, \check{a}) \right| \\ & \leq \gamma \mathbb{E} \left[ \left| \check{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) - \bar{J}^*(\Phi^{\varepsilon^0}(\check{\mu}, \check{a})) \right| \right] \\ & \leq \gamma \mathbb{E} \left[ \left| \check{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) - \bar{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) \right| + \left| \bar{J}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a})) - \bar{J}^*(\Phi^{\varepsilon^0}(\check{\mu}, \check{a})) \right| \right] \\ & \leq \gamma \mathbb{E} \left[ \left| \inf_{\check{a}' \in \mathcal{A}} \check{Q}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}), \check{a}') - \inf_{\check{a}' \in \mathcal{A}} \tilde{Q}^*(\check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}), \check{a}') \right| + \gamma L_{\bar{J}^*} \mathbb{E} \left[ \left\| \check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) - \Phi^{\varepsilon^0}(\check{\mu}, \check{a}) \right\|_{d_{\check{\mathfrak{S}}}} \right] \right], \end{aligned}$$

where the last inequality holds by Lipschitz continuity of  $\bar{J}^*$  on  $\check{\mathfrak{S}}$ , see Assumption **(H4)**.

The second term in the last inequality can be bounded using the simplex discretization properties and Assumption **(H3)**:

$$\mathbb{E} \left[ \left\| \check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) - \Phi^{\varepsilon^0}(\check{\mu}, \check{a}) \right\|_{d_{\check{\mathfrak{S}}}} \right] = \mathbb{E}_{\varepsilon_1^0} \left[ \left\| \text{Proj}_{\check{\mathfrak{S}}} \circ \bar{F}(\check{\mu}, \check{a}, \varepsilon_1^0) - \bar{F}(\check{\mu}, \check{a}, \varepsilon_1^0) \right\|_{d_{\check{\mathfrak{S}}}} \right] \leq \varepsilon_{\check{\mathfrak{S}}}.$$

For the first term, let  $\check{\mu}' = \check{\Phi}^{\varepsilon^0}(\check{\mu}, \check{a}) \in \check{\mathfrak{S}}$  to alleviate the notation, and let us consider  $\check{a}_1^* \in \mathcal{A}$  and  $\check{a}_2^* \in \mathcal{A}$  satisfying:

$$\check{Q}^*(\check{\mu}', \check{a}_1^*) = \inf_{\check{a}' \in \mathcal{A}} \check{Q}^*(\check{\mu}', \check{a}') \quad \text{and} \quad \tilde{Q}^*(\check{\mu}', \check{a}_2^*) = \inf_{\check{a}' \in \mathcal{A}} \tilde{Q}^*(\check{\mu}', \check{a}').$$

The existence of  $\check{a}_1^*$  and  $\check{a}_2^*$  is guaranteed respectively by finiteness of  $\check{\mathfrak{S}} \times \mathcal{A}$ .

We observe that

$$\begin{aligned} & \check{Q}^*(\check{\mu}', \check{a}_1^*) - \tilde{Q}^*(\check{\mu}', \check{a}_2^*) \\ & = \left( \check{Q}^*(\check{\mu}', \check{a}_1^*) - \check{Q}^*(\check{\mu}', \check{a}_2^*) \right) + \left( \check{Q}^*(\check{\mu}', \check{a}_2^*) - \tilde{Q}^*(\check{\mu}', \check{a}_2^*) \right) \\ & \leq 0 + \sup_{(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \mathcal{A}} \left| (\check{Q}^* - \tilde{Q}^*)(\check{\mu}, \check{a}) \right| \\ & \leq \|\check{Q}^* - \tilde{Q}^*\|_\infty. \end{aligned}$$

On the other hand,

$$\check{Q}^*(\check{\mu}', \check{a}_1^*) - \tilde{Q}^*(\check{\mu}', \check{a}_2^*) = - \left( \tilde{Q}^*(\check{\mu}', \check{a}_2^*) - \check{Q}^*(\check{\mu}', \check{a}_1^*) \right) - \left( \check{Q}^*(\check{\mu}', \check{a}_1^*) - \check{Q}^*(\check{\mu}', \check{a}_1^*) \right) \geq -\|\check{Q}^* - \tilde{Q}^*\|_\infty.$$

Combining the above bounds yields that for every  $(\check{\mu}, \check{a}) \in \check{\mathfrak{S}} \times \mathcal{A}$ ,

$$\left| \check{Q}^*(\check{\mu}, \check{a}) - \tilde{Q}^*(\check{\mu}, \check{a}) \right| \leq \gamma \left( \|\check{Q}^* - \tilde{Q}^*\|_\infty \right) + \gamma L_{\bar{J}^*} \varepsilon_{\check{\mathfrak{S}}}.$$

Consequently,

$$\|\tilde{Q}^* - \tilde{Q}^*\|_\infty \leq \frac{\gamma}{1-\gamma} L_{\tilde{J}^*} \varepsilon_{\mathfrak{S}}.$$

**Step 3.** Last, we look at the difference between  $\tilde{Q}^*(\text{Proj}_{\mathfrak{S}}(\mu), \tilde{a})$  and  $\tilde{Q}^*(\mu, \tilde{a})$ . For every  $\mu \in \mathfrak{S}$  and  $\tilde{a} \in \mathcal{A}$ , letting  $\tilde{\mu} = \text{Proj}_{\mathfrak{S}}(\mu)$  to alleviate the notation, we have  $\|\tilde{\mu} - \mu\|_{d_{\mathfrak{S}}} \leq \varepsilon_{\mathfrak{S}}$ . We obtain

$$\begin{aligned} & \left| \tilde{Q}^*(\tilde{\mu}, \tilde{a}) - \tilde{Q}^*(\mu, \tilde{a}) \right| \\ & \leq \left| \tilde{f}(\tilde{\mu}, \tilde{a}) - \tilde{f}(\mu, \tilde{a}) \right| + \gamma \mathbb{E} \left[ \left| \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\Phi(\tilde{\mu}, \tilde{a}), \tilde{a}') - \inf_{\tilde{a}' \in \mathcal{A}} \tilde{Q}^*(\Phi(\mu, \tilde{a}), \tilde{a}') \right| \right] \\ & \leq L_{\tilde{f}} \|\tilde{\mu} - \mu\|_{d_{\mathfrak{S}}} + \gamma \mathbb{E} \left[ \left| \tilde{J}^*(\tilde{F}(\tilde{\mu}, \tilde{a}, \varepsilon^0) - \tilde{J}^*(\tilde{F}(\mu, \tilde{a}, \varepsilon^0)) \right| \right] \\ & \leq L_{\tilde{f}} \varepsilon_{\mathfrak{S}} + \gamma L_{\tilde{J}^*} \mathbb{E} \left[ \|\tilde{F}(\tilde{\mu}, \tilde{a}, \varepsilon^0) - \tilde{F}(\mu, \tilde{a}, \varepsilon^0)\|_{d_{\mathfrak{S}}} \right] \\ & \leq (L_{\tilde{f}} + \gamma L_{\tilde{J}^*} L_{\tilde{F}}) \varepsilon_{\mathfrak{S}}, \end{aligned}$$

where we used the Lipschitz continuity of  $\tilde{f}, \tilde{J}^*, \tilde{F}$  and the assumption on  $\mathfrak{S}$ , see Assumptions **(H3)**, **(H4)** and the simplex discretization properties.  $\square$

**E.2. DDPG algorithm.** In Algorithm 2, we describe the DDPG method for MFC with our notation.

---

**Algorithm 2:** DDPG for MFC

---

**Data:** A number of episodes  $N_{\text{epi}}$ ; a length  $T$  for each episode; a minibatch size  $N_{\text{batch}}$ ; a learning rate  $\tau$ .

**Result:** A strategy function for central planner represented by the target network  $\pi'_{\omega'}$ .

---

```

1 begin
2   Randomly initialize parameters  $\theta$  and  $\omega$  for critic network  $Q_\theta$  and actor network  $\pi_\omega$ 
3   Initialize  $\theta' \leftarrow \theta$  and  $\omega' \leftarrow \omega$  for target networks  $Q'_{\theta'}$  and  $\pi'_{\omega'}$ 
4   for  $k = 0, 1, \dots, N_{\text{epi}} - 1$  do
5     Initial distribution  $\tilde{\mu}_0$ 
6     Initialize replay buffer  $R_{\text{buffer}}$ 
7     for  $n = 0, 1, \dots, T - 1$  do
8       Select an action  $\bar{a}_n = \pi_\omega(\tilde{\mu}_n) + \epsilon_{n+1}^a \in \mathbb{R}^{N_p}$ , where  $\epsilon_{n+1}^a$  is the exploration noise
9       Execute  $\bar{a}_n$ , observe cost  $c_n = \tilde{f}_n(\tilde{\mu}_n, \bar{a}_n)$  and  $\tilde{\mu}_{n+1}$ 
10      Store transition  $(\tilde{\mu}_n, \bar{a}_n, c_n, \tilde{\mu}_{n+1})$  in  $R_{\text{buffer}}$ 
11      Sample a random minibatch of  $N_{\text{batch}}$  transitions  $(\tilde{\mu}_i, \bar{a}_i, c_i, \tilde{\mu}_{i+1})$  from  $R_{\text{buffer}}$ 
12      Set  $y_i = c_i + \gamma Q'_{\theta'}(\tilde{\mu}_{i+1}, \pi'_{\omega'}(\tilde{\mu}_{i+1}))$ , for  $i = 1 \dots, N_{\text{batch}}$ 
13      Update the critic by minimizing the loss:  $L(\theta) = \frac{1}{N_{\text{batch}}} \sum_i (y_i - Q_\theta(\tilde{\mu}_i, \bar{a}_i))^2$ 
14      Update the actor policy using the sampled policy gradient  $\nabla_\omega J$ :
          
$$\nabla_\omega J(\omega) \approx \frac{1}{N_{\text{batch}}} \sum_i \nabla_{\bar{a}} Q_\theta(\tilde{\mu}_i, \pi_\omega(\tilde{\mu}_i)) \nabla_\omega \pi_\omega(\tilde{\mu}_i)$$

          Update target networks:  $\theta' \leftarrow \tau\theta + (1-\tau)\theta'$  and  $\omega' \leftarrow \tau\omega + (1-\tau)\omega'$ 
15 return  $\pi'_{\omega'}$ 

```

---