

Non-Exchangeable Mean Field Markov Decision Processes with common noise : from Bellman equation to quantitative propagation of chaos

Samy Mekkaoui

GT MathsFi , CMAP 2025

Joint work with
Huyên Pham (CMAP, École Polytechnique)

22 October 2025

1 Introduction

- Context and motivations
- Problem formulation

2 Bellman fixed point on the space $\mathcal{P}_\lambda(I \times \mathcal{X})$

- Lifted MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$
- Bellman fixed point on $\mathcal{P}_\lambda(I \times \mathcal{X})$
- Equivalence with the strong formulation

3 Propagation of chaos from V_N towards V

- The N-agent problem as a MDP on the state space \mathcal{X}^N and action space A^N
- Convergence of value functions
- Approximate optimal policies

Introduction

Mean-field approach to large population stochastic control

Mean field approach to large population stochastic control

- Large number of agents N interacting dynamic agents/entities with **heterogeneous** interactions.
- Agents are **cooperative** and act following a social planner.
- When $N \rightarrow \infty$, we get an optimal control of mean-field type.
 - Symmetric agents \rightarrow McKean-Vlasov equations
 - Nonsymmetric agents \rightarrow New limiting systems
- Here, we focus on
 - **Discrete time**, and finite / continuous state space
 - **Infinite Horizon**
 - **Common noise**
 - When $N \rightarrow \infty$: **Conditional** Non exchangeable Markov Decision Process (**CNEMF-MDP**).

\rightarrow Mathematical framework of **reinforcement learning (RL)** with many interacting cooperative agents.

Framework and notations

- Universal filtered probability space $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$.
- **State** and **action** spaces: \mathcal{X} and A (compact and Polish) and $I = [0, 1]$ encoding **heterogeneity** of the agents labeled by $u \in I$.
 - $\mathcal{P}(I \times \mathcal{X})$, resp $\mathcal{P}(A)$, resp $\mathcal{P}(I \times \mathcal{X} \times A)$: set of probability measures on $I \times \mathcal{X}$, resp A , resp $I \times \mathcal{X} \times A$, with Wasserstein distance.
- Discrete time **transition dynamics**
 - **Idiosyncratic noises**: $(\epsilon_t^u)_{u \in I, t \in \mathbb{N}}$, i.i.d valued in E .
 - **Common noise**: $(\epsilon_t^0)_{t \in \mathbb{N}}$ for **all agents**, i.i.d valued in E^0 .
 - **F** measurable function from $I \times \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A) \times E \times E^0 \rightarrow \mathcal{X}$.
- **Reward** on infinite horizon.
 - Discount factor $\beta \in [0, 1)$.
 - **f** measurable bounded function from $I \times \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A) \rightarrow \mathbb{R}$.

The conditional McKean-Vlasov MDP problem

Conditional McKean-Vlasov Markov Decision Processes (**CMKV-MDP**) problem studied by Motte and Pham (see [1]):

$$V(\xi) = \inf_{\alpha \in \mathcal{A}} V^\alpha(\xi) := \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f(X_t, \alpha_t, \mathbb{P}_{(X_t, \alpha_t)}^0) \right], \quad (1)$$

where \mathcal{A} is a suitable class of control with controlled state $X^\alpha = (X_t^\alpha)_{t \in \mathbb{N}}$ dynamics given by :

$$\begin{aligned} X_{t+1}^\alpha &= F(X_t, \alpha_t, \mathbb{P}_{(X_t, \alpha_t)}^0, \epsilon_{t+1}, \epsilon_{t+1}^0), \\ X_0^\alpha &= \xi. \end{aligned} \quad (2)$$

where all the random variables are defined on an abstract filtered probability space $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$.

→ The control problem (1)-(2) can be lifted on the space of measures $\mathcal{P}(\mathcal{X})$ and show that V is law invariant, ie for 2 \mathcal{X} -valued random variables ξ and ξ' satisfying $\mathbb{P}_\xi = \mathbb{P}_{\xi'}$, we have $V(\xi) = V(\xi')$.

Introduction

Context and motivations

→ Extend the known **CMKV-MDP** theory to the case of non exchangeable interactions. Non exchangeable interactions are motivated by recent litterature on **Graphons**.

- **Graphon mean field systems** :
 - Bayraktar, Chakraborty, Ruoyu Wu (22).
 - De Crescenzo, Coppini, Pham (23).
- **Graphon mean field control** (in continuous time):
 - Cao and Laurière (25).
 - De Crescenzo, Fuhrman, Kharroubi and Pham (24).
 - Kharroubi, Mekkaoui and Pham (25).

The agents labeled by $u \in I$ interact through a weighted probability measure in the form $\frac{\int_I G(u,v) \mathbb{P}_{X_t^v}(dx) dv}{\int_I G(u,v) dv}$ where $G : I \times I \ni (u, v) \mapsto G(u, v)$ is a measurable map which measures the weight between agents u and v .

Introduction

Context and motivations

→ Extend the known **CMKV-MDP** theory to the case of non exchangeable interactions. Non exchangeable interactions are motivated by recent litterature on **Graphons**.

- **Graphon mean field systems** :
 - Bayraktar, Chakraborty, Ruoyu Wu (22).
 - De Crescenzo, Coppini, Pham (23).
- **Graphon mean field control** (in continuous time):
 - Cao and Laurière (25).
 - De Crescenzo, Fuhrman, Kharroubi and Pham (24).
 - Kharroubi, Mekkaoui and Pham (25).

The agents labeled by $u \in I$ interact through a weighted probability measure in the form $\frac{\int_I G(u,v) \mathbb{P}_{X_t}(\mathrm{d}x) \mathrm{d}v}{\int_I G(u,v) \mathrm{d}v}$ where $G : I \times I \ni (u, v) \mapsto G(u, v)$ is a measurable map which measures the weight between agents u and v .

→ We want to extend the framework of **CMKV-MDP** by introducing an adequate modelling of the heterogeneity between the agents.

Introduction

The N agent formulation in the CNEMF-MDP control problem

N -agent formulation

- **State dynamics** for the controlled systems $\mathbf{X}^N = (X^{i,N})_{i \in \llbracket 1, N \rrbracket}$

$$\begin{cases} X_0^{i,N} = x_0^i, \\ X_{t+1}^{i,N} = F_N\left(\frac{i}{N}, X_t^{i,N}, \alpha_t^{i,N}, \frac{1}{N} \sum_{j=1}^N \delta_{\left(\frac{j}{N}, X_t^{j,N}, \alpha_t^{j,N}\right)}, \epsilon_{t+1}^i, \epsilon_{t+1}^0\right), \quad t \in \mathbb{N}. \end{cases}$$

- **Value function** for the N -agent system:

$$V_N^\alpha(x_0) := \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f_N\left(\frac{i}{N}, X_t^{i,N}, \alpha_t^{i,N}, \frac{1}{N} \sum_{j=1}^N \delta_{\left(\frac{j}{N}, X_t^{j,N}, \alpha_t^{j,N}\right)}\right) \right],$$

where $\mathbf{x}_0 := (x_0^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}^N$ is the initial vector state of the agents. We then define

$$V_N(\mathbf{x}_0) := \sup_{\alpha \in \mathcal{A}} V_N^\alpha(\mathbf{x}_0).$$

where

$$\mathcal{A} := \left\{ \alpha = (\alpha_t^i)_{i \in 1, N, t \in \mathbb{N}} : \alpha^i \text{ is } \mathbb{F}^N\text{-adapted for each } i \in \llbracket 1, N \rrbracket \right\},$$

and where $\mathbb{F}^N := (\mathcal{F}_t^N)_{t \in \mathbb{N}}$ generated by $\epsilon^N = ((\epsilon_t^i)_{i \in \llbracket 1, N \rrbracket}, \epsilon_t^0)_{t \in \mathbb{N}}$ completed with a family of mutually i.i.d uniform random variables $\mathbf{U}^N = (U_t^i)_{i \in \llbracket 1, N \rrbracket, t \in \mathbb{N}}$ used for randomizing the controls $(\alpha^i)_{i \in \llbracket 1, N \rrbracket}$.

Introduction

The non exchangeable mean field limit

Strong formulation for the non exchangeable mean field limit

- **State dynamics** for the controlled systems $\mathbf{X} = (X^u)_{u \in I}$:

$$\begin{cases} X_0^u = \xi^u, \\ X_{t+1}^u = F(u, X_t^u, \alpha_t^u, \mathbb{P}_{(X_t^v, \alpha_t^v)}^0(dx, da)dv, \epsilon_{t+1}^u, \epsilon_{t+1}^0), \quad t \in \mathbb{N}, \quad u \in I. \end{cases} \quad (3)$$

- **Value function** in the strong formulation:

$$\begin{cases} V_{\text{strong}}^\alpha(\xi) &:= \int_I \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f(u, X_t^u, \alpha_t^u, \mathbb{P}_{(X_t^v, \alpha_t^v)}^0(dx, da)dv) \right] du, \\ V_{\text{strong}}(\xi) &:= \sup_{\alpha \in \mathcal{A}^{\text{strong}}} V_{\text{strong}}^\alpha(\xi), \quad \xi \in \mathcal{I}. \end{cases}$$

where $\xi = (\xi^u)_{u \in I}$ denotes the collection of random initial values and

$$\mathcal{A}^{\text{strong}} := \{ \alpha = (\alpha_t^u)_{u \in I, t \in \mathbb{N}} : \alpha_t^u = \alpha_t(u, \Gamma^u, (\epsilon_s^u)_{s \leq t}, (\epsilon_s^0)_{s \leq t}) \text{ for every } t \in \mathbb{N} \}.$$

where

→ Γ^u denotes the **initial information** available for agent u supposed to admit an extra random variable $U^u \sim \mathcal{U}([0, 1])$ independant of ξ^u and $\mathcal{G}^u = \sigma(\Gamma^u)$ -measurable.

→ \mathcal{I} denotes an **admissible class** of initial conditions ensuring the measurability of $u \mapsto \mathbb{P}_{(X_t^u, \alpha_t^u, (\epsilon_s^u)_{s \leq t})}$ for any $t \in \mathbb{N}$, hence the well posedness of the cost functional

Introduction

Weak formulation for the mean field limit

Weak formulation for the non exchangeable mean field limit

- **State dynamics** for the controlled system X

$$\begin{cases} X_0 = \xi, \\ X_{t+1} = F(\mathbf{U}, X_t, \alpha_t, \mathbb{P}^0_{(\mathbf{U}, X_t, \alpha_t)}, \epsilon_{t+1}, \epsilon_{t+1}^0), \quad t \in \mathbb{N}. \end{cases} \quad (4)$$

- **Value function** in the weak formulation:

$$\begin{cases} V_{\text{weak}}^\alpha(\xi) &:= \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f(\mathbf{U}, X_t, \alpha_t, \mathbb{P}^0_{(\mathbf{U}, X_t, \alpha_t)}) \right], \\ V_{\text{weak}}(\xi) &:= \sup_{\alpha \in \mathcal{A}^{\text{weak}}} V_{\text{weak}}^\alpha(\xi), \quad \xi \in \mathcal{I}. \end{cases}$$

where \mathbf{U} is a uniform random variable $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ encoding the **heterogeneity** and where

$$\mathcal{A}^{\text{weak}} := \{ \alpha = (\alpha_t)_{t \in \mathbb{N}} : \alpha_t = \alpha_t(\mathbf{U}, \Gamma, (\epsilon_s)_{s \leq t}, (\epsilon_s^0)_{s \leq t}) \text{ for every } t \in \mathbb{N} \}.$$

→ The weak formulation (4) should be understood as a **relaxed formulation** of (3) which avoids measurability issues but lacks of pathwise interpretation.

Introduction

Goal of this presentation

We will work under the **weak formulation** and show further a connection with the **strong formulation**.

Objectives :

- Show how the control problem (4)- (5) called **CNEMF-MDP** can be recasted as a standard mean field control problem on the space

$$\mathcal{P}_\lambda(I \times \mathcal{X}) := \{\mu \in \mathcal{P}(I \times \mathcal{X}) : \text{pr}_1 \# \mu = \lambda\} \quad (5)$$

where $\text{pr}_1 : I \times \mathcal{X} \ni (u, x) \mapsto \text{pr}_1(u, x) = u$ and $\#$ is the pushforward notation. We will then characterize the value function V_{weak} as a fixed point of a suitable **Bellman operator** on $\mathcal{P}_\lambda(I \times \mathcal{X})$.

Introduction

Goal of this presentation

We will work under the **weak formulation** and show further a connection with the **strong formulation**.

Objectives :

- Show how the control problem (4)- (5) called **CNEMF-MDP** can be recasted as a standard mean field control problem on the space

$$\mathcal{P}_\lambda(I \times \mathcal{X}) := \{\mu \in \mathcal{P}(I \times \mathcal{X}) : \text{pr}_1 \# \mu = \lambda\} \quad (5)$$

where $\text{pr}_1 : I \times \mathcal{X} \ni (u, x) \mapsto \text{pr}_1(u, x) = u$ and $\#$ is the pushforward notation. We will then characterize the value function V_{weak} as a fixed point of a suitable **Bellman operator** on $\mathcal{P}_\lambda(I \times \mathcal{X})$.

- Show a **quantitative propagation of chaos** for the convergence of the value function of the N -agent MDP V_N towards V_{weak} and V_{strong} for all $\mathbf{x} := (x^i)_{i \in \{1, N\}}$ satisfying a regularity condition to be precised later and show how to construct **approximate optimal policies** for the N -agent MDP from optimal randomized feedback control of the **CNEMF-MDP**.

Introduction

Goal of this presentation

We will work under the **weak formulation** and show further a connection with the **strong formulation**.

Objectives :

- Show how the control problem (4)- (5) called **CNEMF-MDP** can be recasted as a standard mean field control problem on the space

$$\mathcal{P}_\lambda(I \times \mathcal{X}) := \{\mu \in \mathcal{P}(I \times \mathcal{X}) : \text{pr}_1 \# \mu = \lambda\} \quad (5)$$

where $\text{pr}_1 : I \times \mathcal{X} \ni (u, x) \mapsto \text{pr}_1(u, x) = u$ and $\#$ is the pushforward notation. We will then characterize the value function V_{weak} as a fixed point of a suitable **Bellman operator** on $\mathcal{P}_\lambda(I \times \mathcal{X})$.

- Show a **quantitative propagation of chaos** for the convergence of the value function of the N -agent MDP V_N towards V_{weak} and V_{strong} for all $\mathbf{x} := (x^i)_{i \in \{1, N\}}$ satisfying a regularity condition to be precised later and show how to construct **approximate optimal policies** for the N -agent MDP from optimal randomized feedback control of the **CNEMF-MDP**.
- Propose a simple application of our non exchangeable mean field model to the case of **targeting advertising**.

1 Introduction

- Context and motivations
- Problem formulation

2 Bellman fixed point on the space $\mathcal{P}_\lambda(I \times \mathcal{X})$

- Lifted MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$
- Bellman fixed point on $\mathcal{P}_\lambda(I \times \mathcal{X})$
- Equivalence with the strong formulation

3 Propagation of chaos from V_N towards V

- The N-agent problem as a MDP on the state space \mathcal{X}^N and action space A^N
- Convergence of value functions
- Approximate optimal policies

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

Some regularity assumptions

Regularity assumptions on f and F

- Regularity on the **state transition** function F

$$\mathbb{E}[d(F(\mathbf{u}, x, a, \mu, \epsilon_1^1, e^0), F(\mathbf{u}, x', a, \mu', \epsilon_1^1, e^0))] \leq L_F(d(x, x') + \mathcal{W}(\mu, \mu')). \quad (6)$$

- Regularity on the **reward** function f

$$|f(\mathbf{u}, x, a, \mu) - f(\mathbf{u}, x', a, \mu')| \leq L_f(d(x, x') + \mathcal{W}(\mu, \mu')). \quad (7)$$

for every $\mathbf{u} \in I$, $x, x' \in \mathcal{X}$, $a \in A$, $\mu, \mu' \in \mathcal{P}(I \times \mathcal{X} \times A)$ and $e^0 \in E^0$.

- The Lipschitz assumption on F is made on expectation, and not pathwisely.
- The definition of the mean-field limit doesn't require any regularity assumption on the label \mathbf{u} .

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

Define the measurable map $\tilde{F} : I \times \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A) \times E \times E^0 \rightarrow I \times \mathcal{X}$ as

$$\tilde{F}(u, x, a, \mu, e, e^0) = (u, F(u, x, a, \mu, e, e^0)).$$

- Set $\mu_{t+1} = \mathbb{P}_{(U, X_{t+1})}^0 \in \mathcal{P}_\lambda(I \times \mathcal{X})$. Then (using the pushforward notation $\#$):

$$\mu_{t+1} = \tilde{F}(\cdot, \cdot, \cdot, \mathbb{P}_{(U, X_t, \alpha_t)}^0, \cdot, \epsilon_{t+1}^0) \# (\mathbb{P}_{(U, X_t, \alpha_t)}^0 \otimes \lambda_\epsilon) \quad \mathbb{P}\text{-a.s.}, \quad t \in \mathbb{N}. \quad (8)$$

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

Define the measurable map $\tilde{F} : I \times \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A) \times E \times E^0 \rightarrow I \times \mathcal{X}$ as

$$\tilde{F}(u, x, a, \mu, e, e^0) = (u, F(u, x, a, \mu, e, e^0)).$$

- Set $\mu_{t+1} = \mathbb{P}_{(U, X_{t+1})}^0 \in \mathcal{P}_\lambda(I \times \mathcal{X})$. Then (using the pushforward notation $\#$):

$$\mu_{t+1} = \tilde{F}(\cdot, \cdot, \cdot, \mathbb{P}_{(U, X_t, \alpha_t)}^0, \cdot, \epsilon_{t+1}^0) \# (\mathbb{P}_{(U, X_t, \alpha_t)}^0 \otimes \lambda_\epsilon) \quad \mathbb{P}\text{-a.s.}, \quad t \in \mathbb{N}. \quad (8)$$

Considering the \mathbb{F}^0 -adapted control process $\alpha_t = \mathbb{P}_{(U, X_t, \alpha_t)}^0$ (Note that this process has to satisfy $\text{pr}_{12} \# \alpha_t = \mu_t$), and from a suitable **measurable coupling** ensuring that one can find a measurable map

$$\mathbf{p} : \mathcal{P}_\lambda(I \times \mathcal{X}) \times \mathcal{P}_\lambda(I \times \mathcal{X} \times A) \rightarrow \mathcal{P}_\lambda(I \times \mathcal{X} \times A),$$

such that $\text{pr}_{12} \# \mathbf{p}(\mu, \mathbf{a}) = \mu$ and if $\text{pr}_{12} \# \mathbf{a} = \mu$, then $\mathbf{p}(\mu, \mathbf{a}) = \cdot$. It follows that (11) can be rewritten as

$$\mu_{t+1} = \hat{F}(\mu_t, \alpha_t, \epsilon_{t+1}^0), \quad \mathbb{P}\text{-a.s.}, \quad t \in \mathbb{N}, \quad (9)$$

with $\hat{F}(\mu, \mathbf{a}, e^0) := \tilde{F}(\cdot, \cdot, \cdot, \mathbf{p}(\mu, \mathbf{a}), \cdot, e^0) \# ((\mathbf{p}(\mu, \mathbf{a}) \otimes \lambda_\epsilon)$.

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

Lifting the MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$

- Similarly and with law of conditional expectations,

$$V_{\text{weak}}^\alpha(\xi) = \hat{V}^\alpha(\mu_0) = \mathbb{E}\left[\sum_{t \in \mathbb{N}} \beta^t \hat{f}(\mu_t, \alpha_t)\right], \quad \text{with } \mu_0 = \mathbb{P}_{(U, \xi)} \in \mathcal{P}_\lambda(I \times \mathcal{X}). \quad (10)$$

for some measurable function $\hat{f} : \mathcal{P}_\lambda(I \times \mathcal{X}) \times \mathcal{P}_\lambda(I \times \mathcal{X} \times A) \rightarrow \mathbb{R}$ explicitly derived from f :

$$\hat{f}(\mu, \mathbf{a}) := \int_{I \times \mathcal{X} \times A} f(u, x, a, \mathbf{p}(\mu, \mathbf{a})) \mathbf{p}(\mu, \mathbf{a})(du, dx, da).$$

Defining \mathcal{A} as the set of \mathbb{F}^0 -adapted processed valued in $\mathbf{A} = \mathcal{P}_\lambda(I \times \mathcal{X} \times A)$ and denoting $\nu \in \mathcal{A}$, we define

$$\begin{cases} \hat{V}^\nu(\mu_0) &= \mathbb{E}\left[\sum_{t \in \mathbb{N}} \beta^t \hat{f}(\mu_t, \nu_t)\right], \\ \hat{V}(\mu_0) &= \sup_{\nu \in \mathcal{A}} \hat{V}^\nu(\mu_0). \end{cases} \quad (11)$$

with dynamics $\mu_{t+1} = \hat{F}(\mu_t, \nu_t, \epsilon_{t+1}^0)$.

→ From (10), we can see that $V_{\text{weak}}^\alpha(\xi) \leq \hat{V}(\mu)$ when $\mu = \mathbb{P}_{(U, \xi)}$, and the goal is to show the equality.

Bellman operator on the lifted MDP

Definition of the Bellman operator \mathcal{T}

- Bellman operator \mathcal{T} defined on $L_m^\infty(\mathcal{P}_\lambda(I \times \mathcal{X}))$ = bounded \mathbb{R} -valued measurable maps on $\mathcal{P}_\lambda(I \times \mathcal{X})$.

$$[\mathcal{T}W](\mu) := \sup_{\mathbf{a} \in \mathbf{A}} \{ \hat{f}(\mu, \mathbf{a}) + \beta \mathbb{E}[W(\hat{F}(\mu, \mathbf{a}, \epsilon_1^0))] \}, \quad \mu \in \mathcal{P}_\lambda(I \times \mathcal{X}).$$

- operator \mathcal{T} of the lifted MDP: For $\mathcal{W} \in L_m^\infty(\mathcal{P}_\lambda(I \times \mathcal{X}))$,

$$[\mathcal{T}W](\mu) = \sup_{\mathbf{a} \in L^0(I \times \mathcal{X} \times [0,1]; \mathbf{A})} [\mathbb{T}^{\mathbf{a}}W](\mu),$$

where $\mathbb{T}^{\mathbf{a}}$ is an operator defined on $L^\infty(\mathcal{P}_\lambda(I \times \mathcal{X}))$ by

$$[\mathbb{T}^{\mathbf{a}}W](\mu) := \mathbb{E} \left[f(\xi, \mathbf{a}(\xi, \tilde{U}), \mathbb{P}_{(\xi, \mathbf{a}(\xi, \tilde{U}))}) + \beta W(\mathbb{P}_{\tilde{F}(\xi, \mathbf{a}(\xi, \tilde{U}), \mathbb{P}_{(\xi, \mathbf{a}(\xi, \tilde{U}), \epsilon_1^0, \epsilon_1^0)})}^0) \right],$$

for any $(\xi = (U, \xi), \tilde{U}) \sim \mu \otimes \mathcal{U}([0, 1])$.

Theorem

- **Law invariance.** For any ξ and ξ' \mathcal{X} -valued random variables s.t $\mathbb{P}_{(U,\xi)} = \mathbb{P}_{(U,\xi')}$, we have $V_{\text{weak}}(\xi) = V_{\text{weak}}(\xi')$. We then define $V(\mu) := V_{\text{weak}}(\xi)$, for $\mu = \mathbb{P}_{(U,\xi)} \in \mathcal{P}_\lambda(I \times \mathcal{X})$.

Theorem

- **Law invariance.** For any ξ and ξ' \mathcal{X} -valued random variables s.t. $\mathbb{P}_{(U,\xi)} = \mathbb{P}_{(U,\xi')}$, we have $V_{\text{weak}}(\xi) = V_{\text{weak}}(\xi')$. We then define $V(\mu) := V_{\text{weak}}(\xi)$, for $\mu = \mathbb{P}_{(U,\xi)} \in \mathcal{P}_\lambda(I \times \mathcal{X})$.
- **Dynamic Programming.** We have V_{weak} fixed point for the operator \mathcal{T} :

$$V_{\text{weak}}(\mu) = [\mathcal{T}V_{\text{weak}}](\mu), \quad \mu \in \mathcal{P}_\lambda(I \times \mathcal{X})$$

- **Existence of optimal randomized feedback control** a^* **for** $V_{\text{weak}}(\xi)$ in the form:

Theorem

- **Law invariance.** For any ξ and ξ' \mathcal{X} -valued random variables s.t $\mathbb{P}_{(U,\xi)} = \mathbb{P}_{(U,\xi')}$, we have $V_{\text{weak}}(\xi) = V_{\text{weak}}(\xi')$. We then define $V(\mu) := V_{\text{weak}}(\xi)$, for $\mu = \mathbb{P}_{(U,\xi)} \in \mathcal{P}_\lambda(I \times \mathcal{X})$.
- **Dynamic Programming.** We have V_{weak} fixed point for the operator \mathcal{T} :

$$V_{\text{weak}}(\mu) = [\mathcal{T}V_{\text{weak}}](\mu), \quad \mu \in \mathcal{P}_\lambda(I \times \mathcal{X})$$

- **Existence of optimal randomized feedback control a^* for $V_{\text{weak}}(\xi)$ in the form:**

$$\alpha_t^* = a^*(\mathbb{P}_{(U,X_t)}^0, U, X_t, \tilde{U}_t) \quad (12)$$

where $(\tilde{U}_t)_{t \in \mathbb{N}}$ sequence of *i.i.d* uniform random variables for some measurable function $a^*(\mu, u, x, \tilde{u})$ on $\mathcal{P}_\lambda(I \times \mathcal{X}) \times I \times \mathcal{X} \times [0, 1]$.

Theorem

- **Law invariance.** For any ξ and ξ' \mathcal{X} -valued random variables s.t. $\mathbb{P}_{(U,\xi)} = \mathbb{P}_{(U,\xi')}$, we have $V_{\text{weak}}(\xi) = V_{\text{weak}}(\xi')$. We then define $V(\mu) := V_{\text{weak}}(\xi)$, for $\mu = \mathbb{P}_{(U,\xi)} \in \mathcal{P}_\lambda(I \times \mathcal{X})$.
- **Dynamic Programming.** We have V_{weak} fixed point for the operator \mathcal{T} :

$$V_{\text{weak}}(\mu) = [\mathcal{T}V_{\text{weak}}](\mu), \quad \mu \in \mathcal{P}_\lambda(I \times \mathcal{X})$$

- **Existence of optimal randomized feedback control a^* for $V_{\text{weak}}(\xi)$ in the form:**

$$\alpha_t^* = a^*(\mathbb{P}_{(U,X_t)}^0, U, X_t, \tilde{U}_t) \quad (12)$$

where $(\tilde{U}_t)_{t \in \mathbb{N}}$ sequence of *i.i.d* uniform random variables for some measurable function $a^*(\mu, u, x, \tilde{u})$ on $\mathcal{P}_\lambda(I \times \mathcal{X}) \times I \times \mathcal{X} \times [0, 1]$.

- **Hölder property of the value function.** There exists a positive constant $\gamma \leq 1$ such that the value function function is γ -Hölder ie

$$|V(\mu) - V(\mu')| \leq K_* \mathcal{W}(\mu, \mu')^\gamma, \quad \forall (\mu, \mu') \in \mathcal{P}(I \times \mathcal{X}).$$

The strong formulation

Problem formulation

Formulation of the strong formulation

- **State dynamics** for the controlled systems $\mathbf{X} = (X^u)_{u \in I}$:

$$\begin{cases} X_0^u = \xi^u, \\ X_{t+1}^u = F(u, X_t^u, \alpha_t^u, \mathbb{P}_{(X_t^v, \alpha_t^v)}^0(dx, da)dv, \epsilon_{t+1}^u, \epsilon_{t+1}^0), \end{cases} \quad t \in \mathbb{N}, \quad u \in I.$$

- **Value function** in the **strong** formulation :

$$V_{\text{strong}}^\alpha(\xi) := \int_I \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f(u, X_t^u, \alpha_t^u, \mathbb{P}_{(X_t^v, \alpha_t^v)}^0(dx, da)dv) \right] du, \quad \xi = (\xi^u)_{u \in I}.$$

The **value function** of the conditional non exchangeable mean field control Markov decision processes CNEMF-MDP is then defined by

$$V_{\text{strong}}(\xi) := \sup_{\alpha \in \mathcal{A}} V_{\text{strong}}^\alpha(\xi), \quad \xi \in \mathcal{I}.$$

- Note that the uncountable collection of *i.i.d* random variables $(\epsilon^u)_{u \in I}$ induces some measurability issues for the formulation of the strong formulation compared to the weak formulation.

The strong formulation

Equivalence of value functions between weak and strong formulation

Proposition (Equivalence of value functions).

Let $\xi = (\xi^u)_{u \in I}$ and ξ be random variables such that $\mathbb{P}_{\xi^u} = \mathbb{P}_{\xi|U=u}$ for λ a.e $u \in I$. Then , we have

$$V_{\text{strong}}(\xi) = V_{\text{weak}}(\xi) = V(\mu), \quad \mu = \mathbb{P}_{(U,\xi)} = \mathbb{P}_{\xi^u}(\mathrm{d}x)\mathrm{d}u.$$

Proof.

The main idea of the proof follows from the fact that given an optimal randomized feedback policy a for the weak formulation a , it gives an optimal feedback control for the strong formulation by setting for the same a^* in (12).

$$\alpha_t^{\text{strong},u} = a^*(\mathbb{P}_{X_t^v}^0(\mathrm{d}x)\mathrm{d}v, u, X_t^u, \tilde{U}_t^u),$$

since $V_{\text{weak}}^{\alpha^{\text{weak}}}(\xi) = V_{\text{strong}}^{\alpha^{\text{strong}}}(\xi)$ when α^{weak} and α^{strong} are associated to the same policy a . □

→ We now denote indifferently V to denote V_{strong} or V_{weak} .

1 Introduction

- Context and motivations
- Problem formulation

2 Bellman fixed point on the space $\mathcal{P}_\lambda(I \times \mathcal{X})$

- Lifted MDP on $\mathcal{P}_\lambda(I \times \mathcal{X})$
- Bellman fixed point on $\mathcal{P}_\lambda(I \times \mathcal{X})$
- Equivalence with the strong formulation

3 Propagation of chaos from V_N towards V

- The N-agent problem as a MDP on the state space \mathcal{X}^N and action space A^N
- Convergence of value functions
- Approximate optimal policies

The N -agent problem as a MDP on state space \mathcal{X}^N and action space A^N .

Formulation of the N -agent MDP

- **State dynamics** for the N -agent controlled systems $\mathbf{X}^N = (X_i^N)_{i \in \llbracket 1, N \rrbracket}$

$$\begin{cases} X_0^i = x_0^i, \\ X_{t+1}^i = F_N(\frac{i}{N}, X_t^i, \alpha_t^i, \frac{1}{N} \sum_{j=1}^N \delta_{(\frac{j}{N}, X_t^j, \alpha_t^j)}, \epsilon_{t+1}^j, \epsilon_{t+1}^0), \quad t \in \mathbb{N}. \end{cases} \quad (13)$$

where $\mathbf{x}_0 := (x_0^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}^N$ is the initial vector state of the agents.

- **Value function** for the N agent MDP.

$$V_N^\alpha(\mathbf{x}_0) := \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f_N(\frac{i}{N}, X_t^i, \alpha_t^i, \frac{1}{N} \sum_{j=1}^N \delta_{(\frac{j}{N}, X_t^j, \alpha_t^j)}) \right], \quad (14)$$

$$V_N(\mathbf{x}_0) := \sup_{\alpha \in \mathcal{A}} V_N^\alpha(\mathbf{x}_0).$$

The N -agent problem as a MDP on the space \mathcal{X}^N .

MDP on the space \mathcal{X}^N .

- **State dynamics** (13) can be written :

$$\mathbf{X}_{t+1} = \mathbf{F}_N(\mathbf{X}_t, \boldsymbol{\alpha}_t, \boldsymbol{\epsilon}_{t+1}), \quad (15)$$

with **state transition function** $\mathbf{F}_N : \mathcal{X}^N \times A^N \times (E^N \times E^0) \rightarrow \mathcal{X}^N$ is given for $\mathbf{x} = (x^i)_{i \in \llbracket 1, N \rrbracket}$, $\mathbf{a} = (a^i)_{i \in \llbracket 1, N \rrbracket}$ and $\mathbf{e} = ((e^i)_{i \in \llbracket 1, N \rrbracket}, e^0)$ by

$$\mathbf{F}_N(\mathbf{x}, \mathbf{a}, \mathbf{e}) := \left(F_N\left(\frac{i}{N}, x^i, a^i, \frac{1}{N} \sum_{i=1}^N \delta_{(\frac{i}{N}, x^i, a^i)}, e^i, e^0\right) \right)_{i \in \llbracket 1, N \rrbracket},$$

- **Value function** (14) for the N agent MDP :

$$V_N^\alpha(\mathbf{x}_0) = \mathbb{E} \left[\sum_{t \in \mathbb{N}} \beta^t f_N(\mathbf{X}_t, \boldsymbol{\alpha}_t) \right]. \quad (16)$$

with **reward function** $f_N : \mathcal{X}^N \times A^N \rightarrow \mathbb{R}$ is given by

$$f_N(\mathbf{x}, \mathbf{a}) := \frac{1}{N} \sum_{i=1}^N f_N\left(\frac{i}{N}, x^i, a^i, \frac{1}{N} \sum_{i=1}^N \delta_{(\frac{i}{N}, x^i, a^i)}\right), \quad \mathbf{x} = (x^i)_{i \in \llbracket 1, N \rrbracket}, \quad \mathbf{a} = (a^i)_{i \in \llbracket 1, N \rrbracket}.$$

Definition of the Bellman operator on $L_m^\infty(\mathcal{X}^N)$

- **Bellman operator** for the N -agent MDP def

$$[\mathcal{T}_N W](\mathbf{x}) := \sup_{\mathbf{a} \in A^N} \mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}), \quad \mathbf{x} \in \mathcal{X}^N.$$

where

$$\mathbb{T}_N^{\mathbf{a}} W(\mathbf{x}) := f_N(\mathbf{x}, \mathbf{a}) + \beta \mathbb{E} \left[W(\mathbf{F}_N(\mathbf{x}, \mathbf{a}, \epsilon_1)) \right], \quad \mathbf{x} \in \mathcal{X}^N, \quad \mathbf{a} \in A^N.$$

Propagation of chaos of value functions

Regularity of the initial conditions

Assumption on the regularity of the initial condition

For a given $\mathbf{x} := (x^1, x^2, \dots, x^N) \in \mathcal{X}^N$, we say that \mathbf{x} is regular if the following condition holds true. There exists a constant $C > 0$ such that for any $i, j \in \{1, \dots, N\}$,

$$d(x^i, x^j) \leq C \frac{|i - j|}{N}. \quad (17)$$

The set of regular \mathbf{x} will be denoted in the following $\mathcal{X}_{\text{reg}}^N$.

- The assumption (17) is crucial in the derivation of the **propagation of chaos** result

Propagation of chaos of value functions

Regularity in the label state

Assumption on the regularity of f and F with respect to the label state

(i) Let $N \in \mathbb{N}^*$. The mapping

$$I \ni \mathbf{u} \mapsto f(\mathbf{u}, x, a, \mu) \in \mathbb{R}, \quad (18)$$

has a bounded variation on the interval $I_j = [\frac{j-1}{N}, \frac{j}{N}[$ which we denoted by $V_{\frac{j-1}{N}}^{\frac{j}{N}}(f)$ (by omitting the dependance in (x, a, μ) which satisfies

$$\max_{1 \leq j \leq N} \sup_{(x, a, \mu) \in \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A)} V_{\frac{j-1}{N}}^{\frac{j}{N}}(f) \leq \frac{C}{\sqrt{N}}. \quad (19)$$

for every $j \in \{1, \dots, N\}$ and for every $(x, a, \mu) \in \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A)$.

(ii)

$$\mathbb{E}[d(F(\mathbf{u}, x, a, \mu, \epsilon_1^1, e^0), F(\mathbf{u}', x', a', \mu', \epsilon_1^1, e^0))] \leq K_F(d((\mathbf{u}, x, a), (\mathbf{u}', x', a')) + \mathcal{W}(\mu, \mu')), \quad (20)$$

for every $\mathbf{u}, \mathbf{u}' \in I$, $x, x' \in \mathcal{X}$, $a, a' \in A$ and $\mu, \mu' \in \mathcal{P}(I \times \mathcal{X} \times A)$, $e^0 \in E^0$.

Propagation of chaos of value functions

Convergence of F^n and f^n towards F and f

(Convergence of f_N and F_N towards f and F).

There exists two positive decreasing sequences $(\epsilon_N^f)_{N \in \mathbb{N}^*}$ and $(\epsilon_N^F)_{N \in \mathbb{N}^*}$ converging to 0 as $N \rightarrow \infty$ such that

(1)

$$\max_{1 \leq j \leq N} \sup_{(x, a, \mu) \in \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A)} \left| f\left(\frac{j}{N}, x, a, \mu\right) - f_N\left(\frac{j}{N}, x, a, \mu\right) \right| \leq \epsilon_N^f \quad (21)$$

(2)

$$\max_{1 \leq j \leq N} \sup_{(x, a, \mu) \in \mathcal{X} \times A \times \mathcal{P}(I \times \mathcal{X} \times A)} \mathbb{E} \left[d\left(F\left(\frac{j}{N}, x, a, \mu, \epsilon_1^i, \epsilon_1^0\right), F_N\left(\frac{j}{N}, x, a, \mu, \epsilon_1^i, \epsilon_1^0\right)\right) \right] \leq \epsilon_N^F \quad (22)$$

Theorem : Convergence of value functions and propagation of chaos

- For V value function on $\mathcal{P}_\lambda(I \times \mathcal{X})$ of the CNEMF-MDP, we set the lifted function \tilde{V} defined on \mathcal{X}^N by

$$\tilde{V}(\mathbf{x}) := V(\mu_N^\lambda[\mathbf{u}, \mathbf{x}]), \quad \text{for } \mathbf{x} = (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}^N,$$

where $\mu_N^\lambda[\mathbf{u}, \mathbf{x}] := \left(\sum_{j=1}^N \mathbb{1}_{I_j}(u) \delta_{x^j}(\mathrm{d}x) \right) \mathrm{d}u \in \mathcal{P}_\lambda(I \times \mathcal{X})$.

- There exists some positive constant C such that for all $\mathbf{x} := (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}_{\text{reg}}^N$, we have

$$|V_N(\mathbf{x}) - V(\mu_N^\lambda[\mathbf{u}, \mathbf{x}])| \xrightarrow{N \rightarrow \infty} 0. \quad (23)$$

Moreover, **propagation of chaos** rate of convergence takes the following form

$$|V_N(\mathbf{x}) - V(\mu_N^\lambda[\mathbf{u}, \mathbf{x}])| \leq C \left(M_N^\gamma + O(N^{-\frac{\gamma}{2}}) + \epsilon_N^f + (\epsilon_N^F)^\gamma \right).$$

with $M_N := \sup_{\nu \in \mathcal{P}(I \times \mathcal{X} \times A)} \mathbb{E} \left[\mathcal{W}(\nu_N, \nu) \right], \quad (\nu_N \text{ empirical measure of } \nu).$

It extends the result from [1] with the additional errors:

- $O(N^{-\frac{\gamma}{2}})$ which represents the error due to the **label convergence**.
- ϵ_N^f and ϵ_N^F which represent the errors due to the convergence of the **state dynamics functions** and the **reward functions**.

Theorem : Approximate optimal policies

Let $a^* : \mathcal{P}_\lambda(I \times \mathcal{X}) \times I \times \mathcal{X} \times [0, 1] \rightarrow A$ be an optimal **randomized feedback policy** for the **CNEMF-MDP** satisfying the following regularity condition

$$\mathbb{E} \left[| (a(\mu, u, x, U) - a_\epsilon(\mu, u', x', U)) | \right] \leq K (|u - u'| + |x - x'|), \quad (24)$$

for a positive constant $K \geq 0$. Then, denoting

$$\pi_r^{a^*, N}(\mathbf{x}, \tilde{\mathbf{u}}) := (a^*(\mu_N^\lambda[\mathbf{u}, \mathbf{x}], \frac{i}{N}, x^i, \tilde{u}^i))_{i \in \llbracket 1, N \rrbracket}, \quad (25)$$

for $\mathbf{x} := (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}_{\text{reg}}^N$, $\tilde{\mathbf{u}} = (\tilde{u}^i)_{i \in \{1, N\}}$. and defined the **randomized feedback control** $\alpha_t^{r, N} \in \mathcal{A}$ as

$$\alpha_t^{r, N} = \pi_r^{a^*, N}(\mathbf{X}_t, \mathbf{U}_t), \quad t \in \mathbb{N}, \quad (26)$$

where $\{\mathbf{U}_t = (U_t^i)_{i \in \{1, N\}}, t \in \mathbb{N}\}$ is a family of mutually i.i.d uniform random variables on $[0, 1]$, is an $O(M_N^\gamma + N^{-\frac{\gamma}{2}} + \epsilon_N^f + (\epsilon_N^F)^\gamma)$ optimal control for the N -agent MDP.

Conclusion

Main results of our work

Conclusion of our work

- **CNEMF-MDP** lifted to optimization problem on the space $\mathcal{P}_\lambda(I \times \mathcal{X})$ with relaxed controls valued in $\mathbf{A} = \mathcal{P}_\lambda(I \times \mathcal{X} \times A)$ with marginal constraint \rightarrow Standard MFC on the Wasserstein space $\mathcal{P}_\lambda(I \times \mathcal{X})$.
 - Characterization of the **value function** as a fixed point of a Bellman operator.
 - Equivalence formulation between **weak** and **strong** formulation.
 - Existence of an optimal randomized feedback control policy a^* .

Conclusion

Main results of our work

Conclusion of our work

- **CNEMF-MDP** lifted to optimization problem on the space $\mathcal{P}_\lambda(I \times \mathcal{X})$ with relaxed controls valued in $\mathbf{A} = \mathcal{P}_\lambda(I \times \mathcal{X} \times A)$ with marginal constraint \rightarrow Standard MFC on the Wasserstein space $\mathcal{P}_\lambda(I \times \mathcal{X})$.
 - Characterization of the **value function** as a fixed point of a Bellman operator.
 - Equivalence formulation between **weak** and **strong** formulation.
 - Existence of an optimal randomized feedback control policy a^* .
- Convergence of the **value function** V_N of the N -agent MDP towards the value function V for initial state agents $\mathbf{x} := (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}_{\text{reg}}^N$ with **explicit** rate of convergence

Conclusion

Main results of our work

Conclusion of our work

- **CNEMF-MDP** lifted to optimization problem on the space $\mathcal{P}_\lambda(I \times \mathcal{X})$ with relaxed controls valued in $\mathbf{A} = \mathcal{P}_\lambda(I \times \mathcal{X} \times A)$ with marginal constraint \rightarrow Standard MFC on the Wasserstein space $\mathcal{P}_\lambda(I \times \mathcal{X})$.
 - Characterization of the **value function** as a fixed point of a Bellman operator.
 - Equivalence formulation between **weak** and **strong** formulation.
 - Existence of an optimal randomized feedback control policy a^* .
- Convergence of the **value function** V_N of the N -agent MDP towards the value function V for initial state agents $\mathbf{x} := (x^i)_{i \in \llbracket 1, N \rrbracket} \in \mathcal{X}_{\text{reg}}^N$ with **explicit** rate of convergence
- Optimal randomized feedback control for **CNEMF-MDP** \rightarrow Quantitative approximate optimal policy for the N -agent MDP.

Conclusion

Further works on non exchangeable mean field systems

Future works on non exchangeable mean field systems

- Model with controlled interactions $\mathbf{X}^N = (X^{i,N})_{i \in \llbracket 1, N \rrbracket}$

$$\begin{cases} X_{t+1}^{i,N} = F_N\left(\frac{i}{N}, X_t^{i,N}, \frac{1}{N} \sum_{j=1}^N \delta_{(\frac{j}{N}, X_t^{j,N}, \beta_t^{i,j,N})}, \epsilon_{t+1}^i, \epsilon_{t+1}^0\right), & t \in \mathbb{N}. \\ X_0^{i,N} = x_0^i, \end{cases}$$

where $\beta^{i,j,N}$ should be represented as the interaction term between agents i and j .

→ Include controlled graphon interactions $\frac{1}{N} \sum_{j=1}^N G_N(\frac{i}{N}, \frac{j}{N}, \beta_t^{i,j,N})$.

Conclusion

Further works on non exchangeable mean field systems

Future works on non exchangeable mean field systems

- Model with controlled interactions $\mathbf{X}^N = (X^{i,N})_{i \in \llbracket 1, N \rrbracket}$

$$\begin{cases} X_{t+1}^{i,N} = F_N\left(\frac{i}{N}, X_t^{i,N}, \frac{1}{N} \sum_{j=1}^N \delta_{(\frac{j}{N}, X_t^{j,N}, \beta_t^{i,j,N})}, \epsilon_{t+1}^i, \epsilon_{t+1}^0\right), & t \in \mathbb{N}. \\ X_0^{i,N} = x_0^i, \end{cases}$$

where $\beta^{i,j,N}$ should be represented as the interaction term between agents i and j .

→ Include controlled graphon interactions $\frac{1}{N} \sum_{j=1}^N G_N(\frac{i}{N}, \frac{j}{N}, \beta_t^{i,j,N})$.

- Numerical algorithms in the context of the label state formulation

(1) In continuous time :

→ In a model based setting : by DPP principle or by backward algorithms based on Pontryagin formulation.

→ In a model-free setting : By policy gradient and actor-critic algorithms.

Conclusion

Further works on non exchangeable mean field systems

Future works on non exchangeable mean field systems

- Model with controlled interactions $\mathbf{X}^N = (X^{i,N})_{i \in \llbracket 1, N \rrbracket}$

$$\begin{cases} X_{t+1}^{i,N} = F_N\left(\frac{i}{N}, X_t^{i,N}, \frac{1}{N} \sum_{j=1}^N \delta_{(\frac{j}{N}, X_t^{j,N}, \beta_t^{i,j,N})}, \epsilon_{t+1}^i, \epsilon_{t+1}^0\right), & t \in \mathbb{N}. \\ X_0^{i,N} = x_0^i, \end{cases}$$

where $\beta^{i,j,N}$ should be represented as the interaction term between agents i and j .

→ Include controlled graphon interactions $\frac{1}{N} \sum_{j=1}^N G_N(\frac{i}{N}, \frac{j}{N}, \beta_t^{i,j,N})$.

- Numerical algorithms in the context of the label state formulation
 - (1) In continuous time :
 - In a model based setting : by DPP principle or by backward algorithms based on Pontryagin formulation.
 - In a model-free setting : By policy gradient and actor-critic algorithms.
 - (2) In discrete time :
 - Develop reinforcement learning algorithms based on the DPP principle.

Conclusion

Further works on non exchangeable mean field systems

Future works on non exchangeable mean field systems

- Model with controlled interactions $\mathbf{X}^N = (X^{i,N})_{i \in \llbracket 1, N \rrbracket}$

$$\begin{cases} X_{t+1}^{i,N} = F_N\left(\frac{i}{N}, X_t^{i,N}, \frac{1}{N} \sum_{j=1}^N \delta_{(\frac{j}{N}, X_t^{j,N}, \beta_t^{i,j,N})}, \epsilon_{t+1}^i, \epsilon_{t+1}^0\right), & t \in \mathbb{N}. \\ X_0^{i,N} = x_0^i, \end{cases}$$

where $\beta^{i,j,N}$ should be represented as the interaction term between agents i and j .

→ Include controlled graphon interactions $\frac{1}{N} \sum_{j=1}^N G_N(\frac{i}{N}, \frac{j}{N}, \beta_t^{i,j,N})$.

- Numerical algorithms in the context of the label state formulation
 - In continuous time :
 - In a model based setting : by DPP principle or by backward algorithms based on Pontryagin formulation.
 - In a model-free setting : By policy gradient and actor-critic algorithms.
 - In discrete time :
 - Develop reinforcement learning algorithms based on the DPP principle.
- LQ control problem for non exchangeable mean field systems (with common noise).

Problem formulation for the targeted advertising problem

Problem formulation

In the sequel, we extend a targeted advertising model developed in [3] Section 3.5 to the non exchangeable mean field setting. To this end, let C denotes a company, I an influencer working for the advertising of C which can impact customer choices. We model the state space $\mathcal{X} = \{0, 1\}$ where $x = 1$ (resp $x = 0$) indicates that x is (is not) a customer of C . We also model the action space $A = \{0, 1\}$ where $a = 1$ (resp $a = 0$) indicates if I will display (or not) an ad.

- **Reward function** : for $u \in I, x \in \mathcal{X}, a \in A$,

$$f(u, x, a) = x - c^u a,$$

where c^u is an ad cost for agent $u \in I$. It means that if the u labeled user is a customer of C ($x = 1$), it contributes to the revenue of the company but if C had to send him an ad ($a = 1$), it costs c^u to the company.

- **State transition function**. For $(u, x, a) \in I \times \mathcal{X} \times A, e \in [0, 1]$ and $\mu \in \mathcal{P}(I \times \mathcal{X})$,

$$F(u, x, a, \mu, e) = \begin{cases} \mathbf{1}_{e > \int_I G(u, v) \mu^v(\{0\}) dv - 2\eta^u a} & \text{if } x = 0, \\ \mathbf{1}_{e < \int_I G(u, v) \mu^v(\{1\}) dv + 2\eta^u a} & \text{if } x = 1. \end{cases}$$

The parameters of the state transition function F can be interpreted as follows. $\eta^u > 0$ represents the efficiency of an ad to become or remain a customer of C for agent u and a large e indicates an intention to switch operators. The interpretation of the state transition function is then the following: If agent u is in state 0, he's more likely to become a customer of C if $\mu^u(\{0\})$ is low and he receives an ad ($a = 1$).



M. Motte and H. Pham. *Mean-field Markov decision processes with common noise and open-loop controls*, *The Annals of Applied Probability*, 32(2):1421–1458, 2022.



M. Motte and H. Pham. *Quantitative propagation of chaos for mean field Markov decision process with common noise*. *Electronic Journal of Probability*, 28:1–24, 2023.



M. Motte. *Mathematical models for large populations, behavioral economics, and targeted advertising*. PhD thesis, Université Paris Cité, 2021.



S.Mekkaoui, H.Pham *Analysis of Non-Exchangeable Mean Field Markov Decision Processes with common noise : From Bellman equation to quantitative propagation of chaos*. *Work in Progress*

THANK YOU FOR YOUR ATTENTION