Mean-field Markov Decision Process with common noise and open-loop controls

Huyên PHAM

Université Paris Cité, LPSM

Based on joint work with Médéric MOTTE, Université Paris Cité, LPSM



9th Colloquium on BSDE and mean field systems Annecy, june 27-july 1, 2022

Mean-field approach to large population stochastic control

- Large number of N interacting dynamic agents/entities
- Agents are cooperative, and act for collective welfare following a center a decision/social planner. Other interpretation:
 - Influencer controls the state of many individuals in social networks
- ▶ When $N \rightarrow \infty$: optimal control of McKean-Vlasov equation (mean-field control problem). Many papers mostly on continuous-time, and finite horizon
- ▶ Here, we focus on
 - Discrete time and possibly discrete space (graphs)
 - Infinite horizon
 - Common noise
 - When N → ∞: Conditional McKean-Vlasov Markov Decision Process (CMKV-MDP)

 \rightarrow Mathematical framework of reinforcement learning (RL) with many interacting cooperative agents (R. Carmona, M. Laurière and Z. Tan 19, Gu, Guo, Wei and Xu 19)

Motivation from targeted advertising application

- A company C: Internet retailer, candidate to election
- \bullet A social network ${\bf SN}$
 - N connected users of **SN**: state = customer/voter or not of **C**
 - Users data: cookies (track record of visited web pages)

Targeted advertising:

- An influencer (Criteo, etc) I working for C
- I displays personalized online banner ads to users according to their cookies and public data (no direct access to individual states) → Open-loop control and common noise
- Objective: optimize the targeted ad strategy, e.g., in order to attract the largest possible clients/voters given ads costs.

Outline



2 Lifted MDP on $\mathcal{P}(\mathcal{X})$ with relaxed control

3 Convergence of CMKV-MDP

4 Application to targeted advertising

Framework and notations

- A universal probability space $(\Omega, \mathcal{F}, \mathbb{P})$
- State and action spaces: X and A (compact Polish)
 - $\mathcal{P}(\mathcal{X})$, resp. $\mathcal{P}(A)$, resp. $\mathcal{P}(\mathcal{X} \times A)$: set of probability measures on \mathcal{X} , resp. A, resp. $\mathcal{X} \times A$, with Wasserstein distance
- Discrete time transition dynamics
 - Idiosyncratic noises: $(\varepsilon_t^i)_{t \in \mathbb{N}^*}$, for agent $i \in \mathbb{N}^*$, i.i.d. valued in E
 - Common noise: $(\varepsilon_t^0)_{t\in\mathbb{N}^*}$ for all agents, i.i.d. valued in E^0
 - *F*: meas. function from $\mathcal{X} \times A \times \mathcal{P}(\mathcal{X} \times A) \times E \times E^0$ into \mathcal{X}
- Reward on infinite horizon:
 - discount factor $\beta \in [0, 1)$
 - *f*: meas. bounded function from $\mathcal{X} \times A \times \mathcal{P}(\mathcal{X} \times A)$ into \mathbb{R}

Assumptions on transition dynamics and reward

 $\begin{aligned} \left(\mathbf{HF}_{lip}\right) \text{ There exists } & K_F \text{ s.t. for all } a \in A, \ e^0 \in E^0, \ x, x' \in \mathcal{X}, \ \nu, \nu' \in \mathcal{P}(\mathcal{X} \times A), \\ & \mathbb{E}\Big[d_{\mathcal{X}}\Big(F(x, a, \nu, \varepsilon_1^1, e^0), F(x', a, \nu', \varepsilon_1^1, e^0)\Big)\Big] &\leq K_F\Big(d_{\mathcal{X}}(x, x') + \mathcal{W}(\nu, \nu')\Big). \end{aligned} \\ \\ \left(\mathbf{Hf}_{lip}\right) \text{ There exists } & K_f \text{ s.t. for all } a \in A, \ x, x' \in \mathcal{X}, \ \nu, \nu' \in \mathcal{P}(\mathcal{X} \times A), \\ & \left|f(x, a, \nu) - f(x', a, \nu')\right| \leq K_F\Big(d_{\mathcal{X}}(x, x') + \mathcal{W}(\nu, \nu')\Big). \end{aligned}$

Remark: Lipschitz assumption in (HF_{lip}) is made on expectation, not pathwisely.

Information and decentralized open-loop controls in the finite population model

- A subalgebra \mathcal{G} of \mathcal{F} rich enough (used for randomization)
- ξ^i initial state in \mathcal{X} of agent i = 1, ..., N; independent of \mathcal{G}
- Decentralized open-loop control: a sequence $\alpha = (\alpha^1, \dots, \alpha^N)$ of processes valued in A^N , and adapted w.r.t.

$$\mathcal{F}_t^{N} = \sigma\{\xi^i, (\varepsilon_s^i)_{s\leq t}, (\varepsilon_s^0)_{s\leq t}, i=1,\ldots,N\} \vee \mathcal{G}, \quad t \in \mathbb{N}.$$

Remark: No symmetry assumption in control α^{i} , i = 1, ..., N.

Mean-field control problem in the N-population model

• Mean-field controlled dynamics: State process $X^{i,N,\alpha}$ of agent *i* governed by

$$\begin{cases} X_0^{i,N,\alpha} &= \xi^i \\ X_{t+1}^{i,N,\alpha} &= F(X_t^{i,N,\alpha},\alpha_t^i,\frac{1}{N}\sum_{j=1}^N \delta_{(X_t^{j,N,\alpha},\alpha_t^j)},\varepsilon_{t+1}^i,\varepsilon_{t+1}^0). \end{cases}$$
(1)

• Gain functional for each agent i = 1, ..., N:

$$V^{i,N,\alpha} = \mathbb{E}\Big[\sum_{t=0}^{\infty}\beta^t f\left(X_t^{i,N,\alpha}, \alpha_t^i, \frac{1}{N}\sum_{j=1}^N \delta_{(x_t^{j,N,\alpha}, \alpha_t^j)}\right)\Big]$$

• **Optimal gain** for the center of decision (social planner/influencer):

$$V^{N}(\xi^{1},\ldots,\xi^{N}) = \sup_{\alpha} \frac{1}{N} \sum_{i=1}^{N} V^{i,N,\alpha}.$$
 (2)

Remark: Problem (1)-(2) is a standard MDP with state space \mathcal{X}^N , action space A^N , for which it is known that optimization over (randomized/relaxed) open-loop control is equivalent to optimization over (randomized/relaxed) feedback control. But hardly tractable when N is large!

Mean-field control problem in the ∞-population model

Formally, we expect (by "propagation of chaos") a problem formulation with:

• Controlled McKean-Vlasov equation: state X^{α} of representative agent governed by

$$\begin{cases} X_0^{\alpha} = \xi \\ X_{t+1}^{\alpha} = F(X_t^{\alpha}, \alpha_t, \mathbb{P}^0_{(X_t^{\alpha}, \alpha_t)}, \varepsilon_{t+1}, \varepsilon_{t+1}^0). \end{cases}$$
(3)

where the control process α is valued in A, and adapted w.r.t.

$$\mathcal{F}_t = \sigma\{\xi, (\varepsilon_s)_{s \leq t}, (\varepsilon_s^0)_{s \leq t}\} \vee \mathcal{G}, \quad t \in \mathbb{N}.$$

• Gain functional of representative agent

$$V^{\alpha} = \mathbb{E}\Big[\sum_{t=0}^{\infty}\beta^{t}f(X_{t}^{\alpha},\alpha_{t},\mathbb{P}^{0}_{(X_{t}^{\alpha},\alpha_{t})})\Big]$$

• Optimal gain for the center of decision (social planner/influencer):

$$V(\xi) = \sup_{\alpha} V^{\alpha}.$$
 (4)

Remark: Problem (3)-(4) is a priori a nonstandard MDP due to $\mathbb{P}^{0}_{(X^{\alpha}_{t},\alpha_{t})}$, and called CMKV-MDP.

Questions addressed in this talk

- (I) Can we obtain a tractable resolution of CMKV-MDP?
 - Dynamic programming, Bellman equation?
 - Examples of explicit resolution
- (II) Convergence of *N*-agent MDP to CMKV-MDP?
- (1) V^N towards V, and at which rate?
- (2) How to get an approximate optimal control for the *N*-agent problem from an optimal control for the CMKV-MDP? At which accuracy?

Questions addressed in this talk

- (I) Can we obtain a tractable resolution of CMKV-MDP?
 - Dynamic programming, Bellman equation?
 - Examples of explicit resolution
- (II) Convergence of N-agent MDP to CMKV-MDP?
- (1) V^N towards V, and at which rate?
- (2) How to get an approximate optimal control for the *N*-agent problem from an optimal control for the CMKV-MDP? At which accuracy?

► Randomization of controls plays a crucial role! This is a noticeable difference with continuous time framework.

Related literature:

(I) In discrete time:

• Bauerle (21) closed-loop policies, Carmona, Laurière, Tan (22)

(II) mostly for continuous time MKV diffusion and for (1)

- Lacker (18), Djete (20): tightness arguments (no rate of convergence)
- Cecchin (21): Finite state: rate of CV $N^{-1/2}$
- Germain, P., Warin (21): BSDE methods under existence of a smooth solution. Rate of convergence N^{-1}
- Gangbo, Mayorga, Swiech (20), Cardaliaguet, Daudin, Jackson, Souganidis (22): viscosity solutions method. Rate of CV: N^{-γ} for some γ ∈ (0,1].

Outline



2 Lifted MDP on $\mathcal{P}(\mathcal{X})$ with relaxed control

3 Convergence of CMKV-MDP

4 Application to targeted advertising

Reformulation on $\mathcal{P}(\mathcal{X})$ with relaxed/randomized control

• $X = X^{\alpha} \rightsquigarrow \mathsf{CMKV}$ dynamics with open-loop control α :

$$X_{t+1} = F(X_t, \alpha_t, \mathbb{P}^0_{(X_t, \alpha_t)}, \varepsilon_{t+1}, \varepsilon^0_{t+1}).$$

• Set $\mu_t = \mathbb{P}^0_{X_t}$. Then (with the pushforward measure notation *):

$$\mu_{t+1} = F(\cdot, \cdot, \mathbb{P}^{0}_{(X_{t}, \alpha_{t})}, \cdot, \varepsilon^{0}_{t+1}) \star \left(\mathbb{P}^{0}_{(X_{t}, \alpha_{t})} \otimes \lambda_{\varepsilon}\right)$$

Reformulation on $\mathcal{P}(\mathcal{X})$ with relaxed/randomized control

• $X = X^{\alpha} \rightsquigarrow CMKV$ dynamics with open-loop control α :

$$X_{t+1} = F(X_t, \alpha_t, \mathbb{P}^0_{(X_t, \alpha_t)}, \varepsilon_{t+1}, \varepsilon^0_{t+1}).$$

• Set $\mu_t = \mathbb{P}^0_{X_t}$. Then (with the pushforward measure notation *):

$$\mu_{t+1} = F(\cdot, \cdot, \mathbb{P}^{0}_{(X_{t}, \alpha_{t})}, \cdot, \varepsilon^{0}_{t+1}) \star \left(\mathbb{P}^{0}_{(X_{t}, \alpha_{t})} \otimes \lambda_{\varepsilon}\right)$$

Bayes formula: $\mathbb{P}^{0}_{(X_{t},\alpha_{t})} = \mu_{t} \cdot \hat{\alpha}_{t}$, where $\hat{\alpha}_{t}$ is the probability kernel on $\mathcal{X} \times A$:

$$\hat{\alpha}_t : x \in \mathcal{X} \quad \longmapsto \quad \mathcal{L}^0(\alpha_t | X_t = x) \in \mathcal{P}(A)$$

 \rightarrow Controlled stochastic Fokker-Planck equation on $\mathcal{P}(\mathcal{X})$:

$$\mu_{t+1} = \hat{F}(\mu_t, \hat{\alpha}_t, \varepsilon_{t+1}^0), \quad t \in \mathbb{N},$$

with relaxed ($\mathcal{P}(A)$ -valued) feedback control $\hat{\alpha}$ valued in $\hat{A} \equiv L^0(\mathcal{X}; \mathcal{P}(A))$, and $\hat{F}(\mu, \hat{a}, e^0) = F(\cdot, \cdot, \mu \cdot \hat{a}, \cdot, e^0) \star ((\mu \cdot \hat{a}) \otimes \lambda_{\varepsilon})$.

• Similarly and with law of conditional expectations, we have

$$V^{\alpha} = \mathbb{E}\Big[\sum_{t=0}^{\infty}\beta^{t}\hat{f}(\mu_{t},\hat{\alpha}_{t})\Big]$$

for some function $\hat{f} : \mathcal{P}(\mathcal{X}) \times \hat{A} \to \mathbb{R}$ explicitly derived from f.

Bellman operator of the lifted MDP

• Bellman operator \mathcal{T} of the lifted MDP: for $W \in L^{\infty}_m(\mathcal{P}(\mathcal{X}))$,

$$\mathcal{T}[W](\mu) := \sup_{\hat{a} \in \hat{A}} \mathcal{T}^{\hat{a}}[W](\mu) := \sup_{\hat{a} \in \hat{A}} \left\{ \hat{f}(\mu, \hat{a}) + \beta \mathbb{E} \left[W(\hat{F}(\mu, \hat{a}, \varepsilon_1^0)) \right] \right\}.$$

▶ \mathcal{T} is well-defined and contractive on $L^{\infty}_m(\mathcal{P}(\mathcal{X})) \rightarrow$ unique fixed point, denoted V^* : $V^* = \mathcal{T}[V^*].$

Characterization by Bellman equation on $\mathcal{P}(\mathcal{X})$

Theorem

(i) Law invariance. For any ξ , $\tilde{\xi}$ s.t. $\mathbb{P}_{\xi} = \mathbb{P}_{\tilde{\xi}}$, we have $V(\xi) = V(\tilde{\xi})$. We then define $V(\mu) := V(\xi)$, for $\mu = \mathbb{P}_{\xi} \in \mathcal{P}(\mathcal{X})$.

(ii) **Dynamic Programming (DP)**. $V = V^*$, hence satisfies the Bellman fixed point equation:

$$V(\mu) = \mathcal{T}[V](\mu) = \sup_{\hat{s} \in \hat{A}} \mathcal{T}^{\hat{s}}[V](\mu), \quad \mu \in \mathcal{P}(\mathcal{X}).$$

(iii) For all $\epsilon > 0$, there exists an ϵ -optimal randomized feedback control for $V(\xi)$ in the form:

$$\alpha_t^{\varepsilon} := \mathsf{a}_{\epsilon}(\mathbb{P}^0_{X_t}, X_t, U_t), \quad t \in \mathbb{N}.$$

where $(U_t)_t$ sequence of i.i.d. \mathcal{G} -measurable ~ $\mathcal{U}([0,1])$, for some measurable function $a_{\epsilon}(\mu, x, u)$ on $\mathcal{P}(\mathcal{X}) \times \mathcal{X} \times [0,1]$ constructed from the argmax in $\mathcal{T}^{\hat{a}}$.

Outline



2 Lifted MDP on $\mathcal{P}(\mathcal{X})$ with relaxed control

3 Convergence of CMKV-MDP

4 Application to targeted advertising

Convergence of the N-agent MDP

- Usually based on propagation of chaos on state process $X^{i,N}$ towards X pathwisely or in law (symmetry arguments is crucial), and then deduce convergence of V^N towards V
- Here, we do not have in general propagation of chaos on $X^{i,N,\alpha}$ controlled by $\alpha = (\alpha^1, \ldots, \alpha^N)$, which is not assumed to be symmetric.

Convergence of the N-agent MDP

- Usually based on propagation of chaos on state process $X^{i,N}$ towards X pathwisely or in law (symmetry arguments is crucial), and then deduce convergence of V^N towards V
- Here, we do not have in general propagation of chaos on $X^{i,N,\alpha}$ controlled by $\alpha = (\alpha^1, \ldots, \alpha^N)$, which is not assumed to be symmetric.
- \bullet Instead, we prove a propagation of chaos on the Bellman operator of the N-agent MDP defined by

$$\mathcal{T}_{N}^{a}[W](x) := f(x, a) + \beta \mathbb{E} \Big[W(F(x, a, (\varepsilon_{1}^{i})_{i \in [1, N]}, \varepsilon_{1}^{0})) \Big]$$

where for $\boldsymbol{x} = (x^i)_{i \in [\![1,N]\!]} \in \mathcal{X}^N$, $\boldsymbol{a} = (a^i)_{i \in [\![1,N]\!]} \in \mathcal{A}^N$,

$$\begin{aligned} \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{a}) &:= \quad \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{f}\left(\boldsymbol{x}^{i}, \boldsymbol{a}^{i}, \frac{1}{N} \sum_{j=1}^{N} \delta_{(\boldsymbol{x}^{j}, \boldsymbol{a}^{j})}\right) \\ \boldsymbol{F}(\boldsymbol{x}, \boldsymbol{a}, (\boldsymbol{e}^{i})_{i \in [\![1,N]\!]}, \boldsymbol{e}^{0}) &:= \quad \left(\boldsymbol{F}(\boldsymbol{x}^{i}, \boldsymbol{a}^{i}, \frac{1}{N} \sum_{j=1}^{N} \delta_{(\boldsymbol{x}^{j}, \boldsymbol{a}^{j})}, \boldsymbol{e}^{i}, \boldsymbol{e}^{0})\right)_{i \in [\![1,N]\!]} \in \mathcal{X}^{N}. \end{aligned}$$

Propagation of chaos of the Bellman operator

For V value function on $\mathcal{P}(\mathcal{X})$ of the CMKV-MDP, we set \check{V} defined on \mathcal{X}^N by

$$\begin{split} \check{V}(\pmb{x}) &\coloneqq V(\mu_{\scriptscriptstyle N}[\pmb{x}]), \quad \text{for } \pmb{x} = (x^i)_{i \in [\![1,N]\!]} \in \mathcal{X}^N, \\ \text{where } \mu_{\scriptscriptstyle N}[\pmb{x}] = \frac{1}{N} \sum_{i=1}^N \delta_{x^i}. \end{split}$$

Proposition

There exists some positive constant C s.t. for all $\mathbf{x} = \mathcal{X}^N$, $N \in \mathbb{N}^*$,

$$\left|\sup_{\boldsymbol{a}\in\mathcal{A}^{N}}\mathcal{T}_{N}^{\boldsymbol{a}}[\check{V}](\boldsymbol{x}) - \sup_{\hat{a}\in\hat{A}}\mathcal{T}^{\hat{a}}[V](\mu_{N}[\boldsymbol{x}])\right| \leq CM_{N}^{\gamma},$$

where $\gamma = \min \left[1, \frac{|\ln \beta|}{(\ln 2K_F)_+}\right]$, and

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

Remark. From Fournier-Guillin 15: $M_N \xrightarrow[N \to \infty]{} 0$, and for $\mathcal{X} \times A \subset \mathbb{R}^d$

•
$$M_N = \mathcal{O}(N^{-\frac{1}{2}})$$
 for $d = 1$
• $M_N = \mathcal{O}(N^{-\frac{1}{2}}\log(1+N))$ for $d = 2$
• $M_N = \mathcal{O}(N^{-\frac{1}{d}})$ for $d \ge 3$

Convergence rate of N-agent MDP

Theorem

1. Value function. There exists some positive constant C s.t. for all $x \in \mathcal{X}^N$,

$$|V^N(\mathbf{x}) - V(\mu_N[\mathbf{x}])| \leq CM_N^{\gamma}.$$

2. Optimal control. Let $a_{\varepsilon} : \mathcal{P}(\mathcal{X}) \times \mathcal{X} \times [0,1] \to A$, be an ε -optimal randomized feedback policy for CMKV-MDP. Then, it defines an $(\varepsilon + \mathcal{O}(M_N^{\gamma}))$ -optimal randomized feedback control $\alpha^{\varepsilon} = (\alpha^{\varepsilon,i})_{i \in [\![1,N]\!]}$ for the *N*-agent MDP with

$$\alpha_t^{\varepsilon,i} = \mathsf{a}_{\varepsilon}(\mu_{\mathsf{N}}[\mathbf{X}_t], \mathbf{X}_t^i, \mathbf{U}_t^i), \quad t \in \mathbb{N}, \ i = 1, \dots, \mathsf{N},$$

where $\boldsymbol{X} = (X^i)_{i \in [\![1,N]\!]}$ is the state of the *N*-agent controlled by $\boldsymbol{\alpha}^{\varepsilon}$, and U_t^i , i = 1, ..., N, $t \in \mathbb{N}$, are i.i.d. $\sim \mathcal{U}([0,1])$.

Interpretation: Policies to agents in the population are applied randomly by the central planner in a non-symmetric assignment.

Outline

Problem formulation

2 Lifted MDP on $\mathcal{P}(\mathcal{X})$ with relaxed control





Back to targeted advertising example: spaces and noise

- State space $\mathcal{X} = \{0, 1\}$:
 - x = 1 (resp. 0): customer (not customer) of company **C**
- Action space *A* = {0,1}:

• a = 1 (resp. 0): I displays (or not) an ad

- For each **SN** user *i*:
 - ε_t^i : uniform r.v. representing e.g. time spent at day t on a forum about product sold by **C**
- For simplicity here, no common noise

Targeted advertising example: dynamics and reward

• State transition function:

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

- Large e: eager to change of operator
- $\mu(\{0\})$: proportion of **SN** users that are not customers of **C**
- $\eta > 0$: efficiency of ad for incentive to become or remain a customer of **C**
- Reward function: for $x \in \mathcal{X} = \{0, 1\}, a \in A = \{0, 1\}, a$

$$f(x,a) = x - ca,$$

• c > 0: ad cost

Lifted to deterministic control problem on [0,1]

- State variable: $p_t \equiv$ proportion of **SN** users that are customers of **C**
- (Relaxed) control variable: $q_t \equiv$ probability of displaying an ad to **SN** users (sending an ad to q_t proportion of the **SN** users)

$$p_{t+1} = p_t + q_t \min(\eta, 1 - p_t), \quad t \in \mathbb{N}.$$

• Value function on [0,1]:

 \rightarrow

$$V(p) = \sup_{q} \sum_{t=0}^{\infty} \beta^{t}(p_{t} - cq_{t}), \quad p \in [0,1].$$

Disjunction depending on the ratio cost/efficiency of ad



Three cases according to the position of $\frac{c}{n}$ relative to β and $\frac{\beta}{1-\beta}$.

 \rightarrow We next focus on the case: $\beta < \frac{c}{n} < \frac{\beta}{1-\beta}$.

Optimal control



Case $\beta < \frac{c}{\eta} < \frac{\beta}{1-\beta}$. Optimal policy (in red).

















Summary of main results

- CMKV-MDP lifted to optimization problem on space of laws with relaxed controls:
 - Dynamic Programming Bellman fixed point equation characterizing the value function
 - ε -optimal randomized feedback policy a_{ε}
- Examples of explicit resolution of the lifted MDP
- *N*-agent MDP $\xrightarrow[N \to \infty]{}$ CMKV-MDP with explicit rate of convergence.
- (Approximate) optimal randomized feedback control of CMKV-MDP \rightarrow Quantitative approximation of optimal control for the *N*-agent MDP

Some remarks

- Open loop vs feedback controls vs randomized feedback controls
 - In standard MDP, it is well-known that: sup. over open-loop/feedback/randomized feedback controls give same value
 - Here with mean-field dependence, we have:
 - Sup. over open-loop control = Sup. over randomized feedback control
 - > Sup. over feedback control.
 - This differs from continuous time McKean-Vlasov control where randomization does not yield greater gain



- M. Motte, H.P.: MEAN-FIELD MARKOV DECISION PROCESS WITH COMMON NOISE AND OPEN-LOOP CONTROLS, to appear in *Annals of Applied Probability*.
- M. Motte: MATHEMATICAL MODELS FOR LARGE POPULATIONS, BEHAVORIAL ECONOMICS, AND TARGETED ADVERTISING, PhD thesis (2021).

Thank you for your attention