Mean Field Games: Numerical Methods and Applications in Machine Learning Part 8: Learning in MFGs

## Mathieu LAURIÈRE

https://mlauriere.github.io/teaching/MFG-PKU-8.pdf

Peking University Summer School on Applied Mathematics July 26 – August 6, 2021

## 1. Introduction

2. Learning/Optimization Methods

3. Reinforcement Learning Methods

4. Unifying RL for MFC and MFG: a Two Timescale Approach

## Warning

▲ Terminology "learning":

1 2 A Terminology "learning":

• Game theory, economics, ...:

Fudenberg & Levine [FL09]<sup>1</sup>: "The theory of learning in games [...] examines how, which, and what kind of equilibrium might arise as a consequence of a long-run nonequilibrium process of learning, adaptation, and/or imitation"

<sup>&</sup>lt;sup>1</sup> Fudenberg, D., & Levine, D. K. (2009). Learning and equilibrium. Annu. Rev. Econ., 1(1), 385-420.

A Terminology "learning":

• Game theory, economics, ...:

Fudenberg & Levine [FL09]<sup>1</sup>: "The theory of learning in games [...] examines how, which, and what kind of equilibrium might arise as a consequence of a long-run nonequilibrium process of learning, adaptation, and/or imitation"

Machine learning, RL, ...:

Mitchell  $[M^+97]^2$ : "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E."

<sup>&</sup>lt;sup>1</sup>Fudenberg, D., & Levine, D. K. (2009). Learning and equilibrium. Annu. Rev. Econ., 1(1), 385-420.

<sup>&</sup>lt;sup>2</sup>Mitchell, T. M. (1997). *Machine Learning*. New York: McGraw-Hill. ISBN: 978-0-07-042807-2

## Learning/Optimization Algorithms in Games

#### Learning/optimization methods:

- Fixed point iteration
  - Banach-Picard iterations
  - idem + damping/mixing/smoothing
  - Fictitious Play (FP)
- Online Mirror Descent (OMD)
- ...

3 4

#### Learning/optimization methods:

- Fixed point iteration
  - Banach-Picard iterations
  - idem + damping/mixing/smoothing
  - Fictitious Play (FP)
- Online Mirror Descent (OMD)
- Ο...

#### in

- Games, particularly in economics, see e.g. Fudenberg & Levine [FL<sup>+</sup>98]<sup>3</sup>
- Non-atomic games. see e.g. Hadikhanloo et al. [HLMS21]<sup>4</sup>
- Mean Field Games, see e.g. Hadikhanloo [Had18]<sup>5</sup>

<sup>&</sup>lt;sup>3</sup>Fudenberg, D., & Levine, D. (1998). The Theory of Learning in Games. *The MIT Press.* 

<sup>&</sup>lt;sup>4</sup> Hadikhanloo, S., Laraki, R., Mertikopoulos, P., & Sorin, S. (2021). Learning in nonatomic games, Part I: Finite action spaces and population games. arXiv preprint arXiv:2107.01595.

<sup>&</sup>lt;sup>5</sup>Hadikhanloo, S. (2018). *Learning in Mean Field Games* (Doctoral dissertation, Université Paris sciences et lettres).

#### Generic structure: repeated game (iterations)

- Update the representative agent behavior
  - value function
  - policy (control)
- Update the population behavior

#### Generic structure: repeated game (iterations)

- Update the representative agent behavior
  - value function
  - policy (control)
- Update the population behavior

#### Where is there learning?

- $\rightarrow$  First type of "Learning": meta-algorithm / outside loop
- $\rightarrow~$  Second type of "Learning": agent's viewpoint / inner loop

## MFG Setup

#### Generic Mean Field model: for a typical infinitesimal agent

#### • Dynamics: discrete time

$$X_{n+1}^{\boldsymbol{\alpha},\boldsymbol{\mu}} = F(X_n^{\boldsymbol{\alpha},\boldsymbol{\mu}}, \boldsymbol{\alpha_n}, \boldsymbol{\mu_n}, \boldsymbol{\epsilon_{n+1}}, \boldsymbol{\epsilon_{n+1}}^0), \quad n \ge 0, \qquad X_0^{\boldsymbol{\alpha},\boldsymbol{\mu}} \sim \boldsymbol{\mu_0}$$

- $X_n^{\alpha,\mu} \in \mathcal{X} \subseteq \mathbb{R}^d$ : state,  $\alpha_n \in \mathcal{U} \subseteq \mathbb{R}^k$ : action
- $\epsilon_n \sim \nu$ : idiosyncratic noise,  $\epsilon_n^0 \sim \nu^0$ : common noise (random env.)
- $p(x'|x, a, \mu)$ : corresponding transition probability distribution
- $\mu_n \in \mathcal{P}(\mathcal{X} \times \mathcal{A})$ : a state-action distribution
- $\pi_n$ : a policy; randomized actions:  $\alpha_n \sim \pi_n(\cdot | s_n, \mu_n)$

## MFG Setup

#### Generic Mean Field model: for a typical infinitesimal agent

#### Dynamics: discrete time

$$X_{n+1}^{\boldsymbol{\alpha},\boldsymbol{\mu}} = F(X_n^{\boldsymbol{\alpha},\boldsymbol{\mu}}, \boldsymbol{\alpha}_n, \boldsymbol{\mu}_n, \boldsymbol{\epsilon}_{n+1}, \boldsymbol{\epsilon}_{n+1}^0), \quad n \ge 0, \qquad X_0^{\boldsymbol{\alpha},\boldsymbol{\mu}} \sim \boldsymbol{\mu}_0$$

- $X_n^{\alpha,\mu} \in \mathcal{X} \subseteq \mathbb{R}^d$ : state,  $\alpha_n \in \mathcal{U} \subseteq \mathbb{R}^k$ : action
- $\epsilon_n \sim \nu$ : idiosyncratic noise,  $\epsilon_n^0 \sim \nu^0$ : common noise (random env.)
- ▶  $p(x'|x, a, \mu)$ : corresponding transition probability distribution
- $\mu_n \in \mathcal{P}(\mathcal{X} \times \mathcal{A})$ : a state-action distribution
- $\pi_n$ : a policy; randomized actions:  $\alpha_n \sim \pi_n(\cdot | s_n, \mu_n)$

• Reward 
$$\underline{\Lambda}$$
:  $\mathbb{J}(\pi;\mu) = \mathbb{E}_{\epsilon,\epsilon^0} \left[ \sum_{n=0}^{\infty} \gamma^n r \left( X_n^{\alpha,\mu}, \alpha_n, \mu_n \right) \right]$ 

Two scenarios:

• Cooperative (MFC): Find  $\pi^*$  s.t.

$$\pi^*$$
 maximizes  $\pi \mapsto J^{MFC}(\pi) = \mathbb{J}(\pi; \mu^{\pi})$  where  $\mu_n^{\pi} = \mathbb{P}^0_{X_n^{\alpha, \mu^{\pi}}}$ 

Non-Cooperative (MFG): Find  $(\hat{\pi}, \hat{\mu})$  s.t.

$$\begin{cases} \hat{\pi} \text{ maximizes } \pi \mapsto J^{MFG}(\pi; \hat{\mu}) = \mathbb{J}(\pi; \hat{\mu}) \\ \hat{\mu}_n = \mathbb{P}^0_{X_n^{\hat{\alpha}, \hat{\mu}}} \end{cases}$$

## Best Response and Population Behavior Maps

We focus on MFG and write  $J = J^{MFG}$ . For simplicity let's forget the common noise.

Two important functions:

Best Response map:

$$\mathsf{BR}: \boldsymbol{\mu} \mapsto \boldsymbol{\pi} \in \operatorname{argmax} J^{MFG}(\cdot; \boldsymbol{\mu})$$

• **Population Behavior** induced when everyone using a policy:

$$PB: \pi \mapsto \mu: \mu_{n+1} = \Phi(\mu_n, \pi_n)$$

where:

$$\Phi(\mu,\pi)(x) := \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} p(x|x_0, a, \mu) \pi(a|x_0, \mu) \mu(x_0), \qquad x \in \mathcal{S}$$

represents a one-step transition of the population distribution

## Best Response and Population Behavior Maps

We focus on MFG and write  $J = J^{MFG}$ . For simplicity let's forget the common noise.

Two important functions:

Best Response map:

$$\mathsf{BR}: \boldsymbol{\mu} \mapsto \boldsymbol{\pi} \in \operatorname{argmax} J^{MFG}(\cdot; \boldsymbol{\mu})$$

• **Population Behavior** induced when everyone using a policy:

$$\mathsf{PB}: \pi \mapsto \mu: \mu_{n+1} = \Phi(\mu_n, \pi_n)$$

where:

$$\Phi(\mu,\pi)(x) := \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} p(x|x_0, a, \mu) \pi(a|x_0, \mu) \mu(x_0), \qquad x \in \mathcal{S}$$

represents a one-step transition of the population distribution

**Mean Field Nash equilibrium:**  $(\hat{\mu}, \hat{\pi})$  such that

$$\begin{cases} \hat{\mu} = PB(\pi) \\ \hat{\pi} = BR(\hat{\mu}) \end{cases}$$

 $\underline{\wedge} \ \hat{\mu}$  can be unique without  $\hat{\pi}$  being unique!

## 1. Introduction

## 2. Learning/Optimization Methods

3. Reinforcement Learning Methods

4. Unifying RL for MFC and MFG: a Two Timescale Approach

#### Generic structure: repeated game (iterations)

- Update the representative agent behavior
  - value function
  - policy (control)
- Update the population behavior

#### Where is there learning?

- $\rightarrow$  First type of "Learning": meta-algorithm / outside loop
- → Second type of "Learning": agent's viewpoint / inner loop

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

**Convergence:** holds under strict contraction property for the map:

 $\mu^{(\mathbf{k})} \mapsto \mu^{(\mathbf{k}+1)}$ 

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

**Convergence:** holds under strict contraction property for the map:

 $\mu^{(\mathbf{k})} \mapsto \mu^{(\mathbf{k}+1)}$ 

Typically ensured by assuming that

• 
$$\mu^{(k)} \mapsto \pi^{(k+1)}$$
  
•  $\pi^{(k+1)} \mapsto \mu^{(k+1)}$ 

are Lipschitz with small enough Lipschitz constants

A First assumption is hard to check! Can be relaxed with entropy regularization

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

**Convergence:** holds under strict contraction property for the map:

 $\mu^{(\mathbf{k})} \mapsto \mu^{(\mathbf{k}+1)}$ 

Typically ensured by assuming that

• 
$$\mu^{(k)} \mapsto \pi^{(k+1)}$$
  
•  $\pi^{(k+1)} \mapsto \mu^{(k+1)}$ 

are Lipschitz with small enough Lipschitz constants

A First assumption is hard to check! Can be relaxed with entropy regularization

Remark: version with damping/mixing/smoothing

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

**Convergence:** holds under strict contraction property for the map:

 $\mu^{(\mathbf{k})} \mapsto \mu^{(\mathbf{k}+1)}$ 

Typically ensured by assuming that

• 
$$\mu^{(k)} \mapsto \pi^{(k+1)}$$
  
•  $\pi^{(k+1)} \mapsto \mu^{(k+1)}$ 

are Lipschitz with small enough Lipschitz constants

A First assumption is hard to check! Can be relaxed with entropy regularization

Remark: version with damping/mixing/smoothing

See e.g., Caines et al. [HMC<sup>+</sup>06]; Guo et al. [GHXZ19]; Anahtarci et al. [AKS20]; ...

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\overline{\mu}^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$
- Update population's average behavior:  $\overline{\mu}^{(k+1)} = \frac{k}{k+1}\overline{\mu}^{(k+1)} + \frac{1}{k+1}\mu^{(k+1)}$

Convergence: holds under (Lasry-Lions) monotonicity structure for the MFG

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\overline{\mu}^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$
- Update population's average behavior:  $\overline{\mu}^{(k+1)} = \frac{k}{k+1}\overline{\mu}^{(k+1)} + \frac{1}{k+1}\mu^{(k+1)}$

**Convergence:** holds under (Lasry-Lions) **monotonicity** structure for the MFG Typically ensured by assuming that:

• p is independent of  $\mu$ 

• *r* is separable: 
$$r(x, a, \mu) = r(x, a) + \tilde{r}(x, \mu)$$

• 
$$\tilde{r}$$
 is monotone:  $\langle \tilde{r}(x,\mu) - \tilde{r}(x,\mu'), \mu - \mu' \rangle \leq 0$ 

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\overline{\mu}^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$
- Update population's average behavior:  $\overline{\mu}^{(k+1)} = \frac{k}{k+1}\overline{\mu}^{(k+1)} + \frac{1}{k+1}\mu^{(k+1)}$

**Convergence:** holds under (Lasry-Lions) **monotonicity** structure for the MFG Typically ensured by assuming that:

- p is independent of  $\mu$
- *r* is separable:  $r(x, a, \mu) = r(x, a) + \tilde{r}(x, \mu)$
- $\tilde{r}$  is monotone:  $\langle \tilde{r}(x,\mu) \tilde{r}(x,\mu'), \mu \mu' \rangle \leq 0$

Example: crowd aversion

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\overline{\mu}^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$
- Update population's average behavior:  $\overline{\mu}^{(k+1)} = \frac{k}{k+1}\overline{\mu}^{(k+1)} + \frac{1}{k+1}\mu^{(k+1)}$

**Convergence:** holds under (Lasry-Lions) **monotonicity** structure for the MFG Typically ensured by assuming that:

• p is independent of  $\mu$ 

• *r* is separable: 
$$r(x, a, \mu) = r(x, a) + \tilde{r}(x, \mu)$$

•  $\tilde{r}$  is monotone:  $\langle \tilde{r}(x,\mu) - \tilde{r}(x,\mu'), \mu - \mu' \rangle \leq 0$ 

Example: crowd aversion

Consequence:

$$0 \geq \left[J(\pi;\mu) - J(\pi;\mu')\right] - \left[J(\pi';\mu) - J(\pi';\mu')\right] =: \mathcal{M}(\pi,\mu,\pi',\mu')$$

If  $(\hat{\mu}, \hat{\pi})$  and  $(\hat{\mu}', \hat{\pi}')$  are two Nash equilibria,

$$\mathcal{M}(\hat{\pi}, \hat{\mu}, \hat{\pi}', \hat{\mu}') = \left[J(\hat{\pi}; \hat{\mu}) - J(\hat{\pi}'; \hat{\mu})\right] + \left[J(\hat{\pi}'; \hat{\mu}') - J(\hat{\pi}; \hat{\mu}')\right]$$
$$\geq \mathcal{E}(\hat{\pi}'; \hat{\mu}) + \mathcal{E}(\hat{\pi}; \hat{\mu}') \geq 0$$

where  $\mathcal{E}$  denotes the exploitability of  $\pi$  facing  $\mu$ :  $\mathcal{E}(\pi;\mu) = \sup J(\cdot;\mu) - J(\pi;\mu) \ge 0$ 

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\overline{\mu}^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$
- Update population's average behavior:  $\overline{\mu}^{(k+1)} = \frac{k}{k+1}\overline{\mu}^{(k+1)} + \frac{1}{k+1}\mu^{(k+1)}$

**Convergence:** holds under (Lasry-Lions) **monotonicity** structure for the MFG Typically ensured by assuming that:

• p is independent of  $\mu$ 

• *r* is separable: 
$$r(x, a, \mu) = r(x, a) + \tilde{r}(x, \mu)$$

•  $\tilde{r}$  is monotone:  $\langle \tilde{r}(x,\mu) - \tilde{r}(x,\mu'), \mu - \mu' \rangle \leq 0$ 

Example: crowd aversion

Consequence:

$$0 \geq \left[J(\pi;\mu) - J(\pi;\mu')\right] - \left[J(\pi';\mu) - J(\pi';\mu')\right] =: \mathcal{M}(\pi,\mu,\pi',\mu')$$

If  $(\hat{\mu}, \hat{\pi})$  and  $(\hat{\mu}', \hat{\pi}')$  are two Nash equilibria,

$$\mathcal{M}(\hat{\pi}, \hat{\mu}, \hat{\pi}', \hat{\mu}') = \left[J(\hat{\pi}; \hat{\mu}) - J(\hat{\pi}'; \hat{\mu})\right] + \left[J(\hat{\pi}'; \hat{\mu}') - J(\hat{\pi}; \hat{\mu}')\right]$$
$$\geq \mathcal{E}(\hat{\pi}'; \hat{\mu}) + \mathcal{E}(\hat{\pi}; \hat{\mu}') \geq 0$$

where  $\mathcal{E}$  denotes the **exploitability** of  $\pi$  facing  $\mu$ :  $\mathcal{E}(\pi; \mu) = \sup J(\cdot; \mu) - J(\pi; \mu) \ge 0$ See e.g., Cardaliaguet & Hadikhanloo [CH17]; Elie *et al.* [EPL<sup>+</sup>20]; Perrin *et al.* [PPL<sup>+</sup>20]; Geist *et al.* [GPL<sup>+</sup>21]; ...

Reminder:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

Reminder:

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

Or, using a Q-function defined as:

$$\begin{aligned} Q_{\pi,\mu}(x,a) \\ &= \mathbb{E}\Big[\sum_{n\geq 0} \gamma^n r(x_n, a_n, \mu)\Big], \quad x_{n+1} \sim p(\cdot | x_n, a_n, \mu), a_{n+1} \sim \pi(\cdot | x_{n+1}), x_0 = x, a_0 = a \\ &= r(x, a, \mu) + \gamma \mathbb{E}[Q_{\pi,\mu}(x', a')], \quad x' \sim p(\cdot | x, a, \mu), a' \sim \pi(\cdot | x') \end{aligned}$$

Reminder:

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

Or, using a Q-function defined as:

$$Q_{\pi,\mu}(x,a) = \mathbb{E}\Big[\sum_{n\geq 0} \gamma^n r(x_n, a_n, \mu)\Big], \quad x_{n+1} \sim p(\cdot|x_n, a_n, \mu), a_{n+1} \sim \pi(\cdot|x_{n+1}), x_0 = x, a_0 = a$$
$$= r(x, a, \mu) + \gamma \mathbb{E}[Q_{\pi,\mu}(x', a')], \quad x' \sim p(\cdot|x, a, \mu), a' \sim \pi(\cdot|x')$$

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's policy:  $\pi^{(k+1)}(x) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q^{(k+1)}(x, a), x \in \mathcal{S}$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

Reminder:

**Method:** For k = 0, 1, 2, ..., K:

- Update agent's policy:  $\pi^{(k+1)} \in BR(\mu^{(k)})$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

Or, using a Q-function defined as:

$$Q_{\pi,\mu}(x,a) = \mathbb{E}\Big[\sum_{n\geq 0} \gamma^n r(x_n, a_n, \mu)\Big], \quad x_{n+1} \sim p(\cdot|x_n, a_n, \mu), a_{n+1} \sim \pi(\cdot|x_{n+1}), x_0 = x, a_0 = a$$
$$= r(x, a, \mu) + \gamma \mathbb{E}[Q_{\pi,\mu}(x', a')], \quad x' \sim p(\cdot|x, a, \mu), a' \sim \pi(\cdot|x')$$

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's policy:  $\pi^{(k+1)}(x) = \underset{\pi \in \Pi}{\operatorname{argmax}_{\pi \in \Pi}} \langle Q^{(k+1)}(x, \cdot), \pi \rangle, x \in S$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's average Q-function:  $\overline{Q}^{(\Bbbk+1)} = \overline{Q}^{(\Bbbk)} + \eta Q^{(\Bbbk+1)}$
- Update agent's policy by mirroring:  $\pi^{(k+1)}(\cdot|x) = \Gamma(\overline{Q}^{(k+1)}(x,\cdot))$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's average Q-function:  $\overline{Q}^{(k+1)} = \overline{Q}^{(k)} + \eta Q^{(k+1)}$
- Update agent's policy by mirroring:  $\pi^{(k+1)}(\cdot|x) = \Gamma(\overline{Q}^{(k+1)}(x,\cdot))$

• Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$ 

where

$$\Gamma(y) := \nabla h^*(y) = \operatorname*{argmax}_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)].$$

with a regularizer  $h: \mathcal{P}(\mathcal{A}) \to \mathbb{R}$  and  $h^*: \mathbb{R}^{|\mathcal{A}|} \to \mathbb{R}$  its convex conjugate defined by  $h^*(y) = \max_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)]$ 

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's average Q-function:  $\overline{Q}^{(k+1)} = \overline{Q}^{(k)} + \eta Q^{(k+1)}$
- Update agent's policy by mirroring:  $\pi^{(k+1)}(\cdot|x) = \Gamma(\overline{Q}^{(k+1)}(x,\cdot))$

• Update population's behavior:  $\mu^{(\mathbf{k}+1)} = \operatorname{Pop}(\pi^{(\mathbf{k}+1)})$ 

where

$$\Gamma(y) := \nabla h^*(y) = \operatorname*{argmax}_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)].$$

with a regularizer  $h: \mathcal{P}(\mathcal{A}) \to \mathbb{R}$  and  $h^*: \mathbb{R}^{|\mathcal{A}|} \to \mathbb{R}$  its convex conjugate defined by  $h^*(y) = \max_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)]$ 

Convergence: typically under monotonicity structure

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's average Q-function:  $\overline{Q}^{(k+1)} = \overline{Q}^{(k)} + \eta Q^{(k+1)}$
- Update agent's policy by mirroring:  $\pi^{(k+1)}(\cdot|x) = \Gamma(\overline{Q}^{(k+1)}(x,\cdot))$

• Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$ 

where

$$\Gamma(y) := \nabla h^*(y) = \operatorname*{argmax}_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)].$$

with a regularizer  $h: \mathcal{P}(\mathcal{A}) \to \mathbb{R}$  and  $h^*: \mathbb{R}^{|\mathcal{A}|} \to \mathbb{R}$  its convex conjugate defined by  $h^*(y) = \max_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)]$ 

#### Convergence: typically under monotonicity structure

Note: Here, no need to compute a BR; just evaluate a Q function & argmax

- Update agent's Q-function:  $Q^{(k+1)} = Q_{\pi^{(k)},\mu^{(k)}}$
- Update agent's average Q-function:  $\overline{Q}^{(k+1)} = \overline{Q}^{(k)} + \eta Q^{(k+1)}$
- Update agent's policy by mirroring:  $\pi^{(k+1)}(\cdot|x) = \Gamma(\overline{Q}^{(k+1)}(x,\cdot))$
- Update population's behavior:  $\mu^{(k+1)} = Pop(\pi^{(k+1)})$

where

$$\Gamma(y) := \nabla h^*(y) = \operatorname*{argmax}_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)].$$

with a regularizer  $h: \mathcal{P}(\mathcal{A}) \to \mathbb{R}$  and  $h^*: \mathbb{R}^{|\mathcal{A}|} \to \mathbb{R}$  its convex conjugate defined by  $h^*(y) = \max_{p \in \mathcal{P}(\mathcal{A})} [\langle y, p \rangle - h(\pi)]$ 

#### Convergence: typically under monotonicity structure

Note: Here, no need to compute a BR; just evaluate a *Q* function & argmax See e.g., Hadikhanloo [Had18]; Pérolat *et al.* [PPE<sup>+</sup>21]; Geist *et al.* [GPL<sup>+</sup>21]; ...

## 1. Introduction

2. Learning/Optimization Methods

## 3. Reinforcement Learning Methods

- Examples of RL Algorithms
- Examples of Applications in MFGs

4. Unifying RL for MFC and MFG: a Two Timescale Approach

#### Generic structure: repeated game (iterations)

- Update the representative agent behavior
  - value function
  - policy (control)
- Update the population behavior

#### Where is there learning?

- $\rightarrow$  First type of "Learning": meta-algorithm / outside loop
- $\rightarrow~$  Second type of "Learning": agent's viewpoint / inner loop

## 1. Introduction

2. Learning/Optimization Methods

## 3. Reinforcement Learning Methods

- Examples of RL Algorithms
- Examples of Applications in MFGs

4. Unifying RL for MFC and MFG: a Two Timescale Approach



#### Source: [OpenAl Spinning Up]<sup>6</sup>

<sup>6</sup> https://spinningup.openai.com/en/latest/spinningup/rl\_intro2.html

Algorithm 1 Deep Q-learning with Experience Replay Initialize replay memory  $\mathcal{D}$  to capacity N Initialize action-value function Q with random weights for episode = 1, M do Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ for t = 1, T do With probability  $\epsilon$  select a random action  $a_t$ otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ Perform a gradient descent step on  $(y_i - Q(\phi_i, a_i; \theta))^2$  according to equation 3 end for end for

Source: Mnih et al. [MKS+13]7

<sup>&</sup>lt;sup>7</sup> Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D. & Riedmiller, M. (2013). Playing Atari with Deep Reinforcement Learning.

## DDPG

#### Algorithm 1 DDPG algorithm

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ . Initialize target network Q' and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$ Initialize replay buffer Rfor episode = 1, M do Initialize a random process N for action exploration Receive initial observation state  $s_1$ for t = 1, T do Select action  $a_t = \mu(s_t|\theta^\mu) + N_t$  according to the current policy and exploration noise Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$ Store transition  $(s_t, a_t, r_t, s_{t+1})$  in RSample a random minibatch of N transitions  $(s_i, a_i, r_i, s_{i+1})$  from RSet  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\Theta^{\prime})$ Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$ Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^{\mu}} J \approx \frac{1}{N} \sum_{i} \nabla_{a} Q(s, a | \theta^{Q}) |_{s=s_{i}, a=\mu(s_{i})} \nabla_{\theta^{\mu}} \mu(s | \theta^{\mu}) |_{s_{i}}$$

Update the target networks:

$$\begin{aligned} \theta^{Q'} &\leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'} \\ \theta^{\mu'} &\leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'} \end{aligned}$$

end for end for

#### Source: Lillicrap et al. [LHP+16]<sup>8</sup>

<sup>&</sup>lt;sup>8</sup>Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D. & Wierstra, D. (2016). Continuous control with deep reinforcement learning. ICLR 2016.

Source: Haarnoja et al. [HZAL18]9

<sup>&</sup>lt;sup>9</sup>Haarnoja, T., Zhou, A., Abbeel, P. & Levine, S. (2018). Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. ICML 2018.

What about the population behavior  $\mu$ ?

What about the population behavior  $\mu$ ?

- Empirical distribution  $\mu^N$
- Histogram (state space discretization)
- $\epsilon$ -net in  $\mathcal{P}(\mathcal{X})$
- Function approximation for the density:
  - Kernels
  - Neural nets: normalizing flows, ...
  - ▶ ...

## 1. Introduction

2. Learning/Optimization Methods

# 3. Reinforcement Learning Methods Examples of RL Algorithms

- Examples of RL Algorithms
- Examples of Applications in MFGs

4. Unifying RL for MFC and MFG: a Two Timescale Approach

## Systemic Risk

Revisiting: Systemic risk model of Carmona, Fouque, Sun [CFS15]

$$J((a_n)_n; (m_n)_n) = -\mathbb{E}\bigg[\sum_{n=0}^{N_T} \left(a_n^2 \underbrace{-qa_n(m_n - X_n)}_{\text{borrow if } X_n < m_n} + \kappa(m_n - X_n)^2\right) + c(m_{N_T} - X_{N_T})^2\bigg]$$

Subj. to:  $X_{n+1} = X_n + [K(m_n - X_n) + a_n] + \epsilon_{n+1} + \epsilon_{n+1}^0$ At equilibrium:  $m_n = \mathbb{E}[X_n]\epsilon^0], n \ge 0$ 

## Systemic Risk

Revisiting: Systemic risk model of Carmona, Fouque, Sun [CFS15]

$$J((a_{n})_{n}; (m_{n})_{n}) = -\mathbb{E}\left[\sum_{n=0}^{N_{T}} \left(a_{n}^{2} \underbrace{-qa_{n}(m_{n} - X_{n})}_{\text{borrow if } X_{n} < m_{n}} + \kappa(m_{n} - X_{n})^{2}\right) + c(m_{N_{T}} - X_{N_{T}})^{2}\right]$$
  
Subj. to:  $X_{n+1} = X_{n} + [K(m_{n} - X_{n}) + a_{n}] + \epsilon_{n+1} + \epsilon_{n+1}^{0}$ 

At equilibrium:  $m_n = \mathbb{E}[X_n | \epsilon^0], n \ge 0$ 

Perrin et al. [PPL+20]: Fictitious Play with Backward Induction or tabular Q-learning



Revisiting: Crowd aversion model of Alumulla, Ferreira, Gomes [AFG17] MFG on  $\mathbb{T}$ ,

$$f(x, m, v) = \frac{1}{2} |v|^2 + \tilde{f}(x) + \ln(m(x)),$$

with  $\tilde{f}(x) = 2\pi^2 \left[ -\sum_{i=1}^d c \sin(2\pi x_i) + \sum_{i=1}^d |c \cos(2\pi x_i)|^2 \right] - 2\sum_{i=1}^d c \sin(2\pi x_i),$ then the solution is given by  $u(x) = c \sum_{i=1}^d \sin(2\pi x_i)$  and  $m(x) = e^{2u(x)} / \int e^{2u}$  tł

Revisiting: Crowd aversion model of Alumulla, Ferreira, Gomes [AFG17] MFG on  $\mathbb{T}$ , н

$$f(x, m, v) = \frac{1}{2} |v|^2 + \tilde{f}(x) + \ln(m(x)),$$
  
with  $\tilde{f}(x) = 2\pi^2 \left[ -\sum_{i=1}^d c \sin(2\pi x_i) + \sum_{i=1}^d |c \cos(2\pi x_i)|^2 \right] - 2\sum_{i=1}^d c \sin(2\pi x_i),$   
then the solution is given by  $u(x) = c \sum_{i=1}^d \sin(2\pi x_i)$  and  $m(x) = e^{2u(x)} / \int e^{2u}$ 





## Flocking

Revisiting: Flocking aversion model of Nourian, Caines, Malhamé [NCM11] Perrin *et al.* [PLP<sup>+</sup>21]: For continuous space problems: **Deep RL** 

- Deep RL (SAC) for the policy ( $\approx$  control)
- Deep NN (normalizing flow) for the population distribution

state = (position, velocity) = 
$$(x, v) \in \mathbb{R}^{2d}$$
, 
$$\begin{cases} x_{n+1} = x_n + v_n \Delta t, \\ v_{n+1} = v_n + a_n \Delta t + \epsilon_{n+1}, \end{cases}$$
with running cost: 
$$f_{\beta}^{\text{flock}}(x, v, \mu) = \left\| \int_{\mathbb{R}^{2d}} \frac{(v - v')}{(1 + \|x - x'\|^2)^{\beta}} \, d\mu(x', v') \right\|^2,$$

where  $\beta \ge 0$ , and  $\mu$  is the position-velocity distribution.

## Flocking

Revisiting: Flocking aversion model of Nourian, Caines, Malhamé [NCM11] Perrin *et al.* [PLP<sup>+</sup>21]: For continuous space problems: **Deep RL** 

- Deep RL (SAC) for the policy ( $\approx$  control)
- Deep NN (normalizing flow) for the population distribution

state = (position, velocity) = 
$$(x, v) \in \mathbb{R}^{2d}$$
, 
$$\begin{cases} x_{n+1} = x_n + v_n \Delta t, \\ v_{n+1} = v_n + \mathbf{a}_n \Delta t + \epsilon_{n+1}, \end{cases}$$
with running cost: 
$$f_{\beta}^{\text{flock}}(x, v, \mu) = \left\| \int_{\mathbb{R}^{2d}} \frac{(v - v')}{(1 + \|x - x'\|^2)^{\beta}} \, d\mu(x', v') \right\|^2,$$

where  $\beta \ge 0$ , and  $\mu$  is the position-velocity distribution.



Video: https://www.youtube.com/watch?v=TdXysW\_FA3k

## **Building Evacuation**

A model for crowd motion during building evacuation:

$$r(x, a, \mu) = -\eta \log(\mu(x)) + 10 \times \mathbb{1}_{floor=0}$$

Pérolat et al. [PPE+21]: OMD (no RL for now)



Initial distribution

## **Building Evacuation**

A model for crowd motion during building evacuation:

$$r(x, a, \mu) = -\eta \log(\mu(x)) + 10 \times \mathbb{1}_{floor=0}$$

Pérolat et al. [PPE+21]: OMD (no RL for now)



## 1. Introduction

- 2. Learning/Optimization Methods
- 3. Reinforcement Learning Methods

## 4. Unifying RL for MFC and MFG: a Two Timescale Approach

**MFControl:** Fix a control v, compute induced distribution  $\mu^v$ , update v, ... **MFGame:** Fix a distribution  $\mu$ , compute best response  $v^{\mu}$ , update  $\mu$ , ...

**MFControl:** Fix a control v, compute induced distribution  $\mu^v$ , update  $v, \ldots$ **MFGame:** Fix a distribution  $\mu$ , compute best response  $v^{\mu}$ , update  $\mu$ , ...

**Unification:** update both  $v, \mu$  simultaneously but at different rates  $\rho^v, \rho^{\mu}$ 

- $\rho^{v} < \rho^{\mu} \Rightarrow v$  evolves slowly  $\Rightarrow$  MFControl
- $\rho^{v} > \rho^{\mu} \Rightarrow \mu$  evolves slowly  $\Rightarrow$  MFGame

**MFControl:** Fix a control v, compute induced distribution  $\mu^v$ , update  $v, \ldots$ **MFGame:** Fix a distribution  $\mu$ , compute best response  $v^{\mu}$ , update  $\mu$ , ...

**Unification:** update both  $v, \mu$  simultaneously but at different rates  $\rho^v, \rho^{\mu}$ 

- $\rho^{v} < \rho^{\mu} \Rightarrow v$  evolves slowly  $\Rightarrow$  MFControl
- $\rho^v > \rho^\mu \Rightarrow \mu$  evolves slowly  $\Rightarrow$  MFGame

**Implementation:** Finite state space  $\mathcal{X}$  and finite action space  $\mathcal{A}$ , stationary problem **Q-learning:** Given  $\mu$ , **optimal** cost-to-go when starting at *x* using action *a* 

$$Q(x,a) = f(x,\mu,a) + \sum_{x' \in \mathcal{X}} p(x'|x,\mu,a) \underbrace{\min_{a'} Q(x',a')}_{=V(x')}.$$

Note: optimal control is  $\hat{v}_Q(x) = \operatorname{argmin}_a Q(x, a)$ .

**MFControl:** Fix a control v, compute induced distribution  $\mu^v$ , update  $v, \ldots$ **MFGame:** Fix a distribution  $\mu$ , compute best response  $v^{\mu}$ , update  $\mu$ , ...

**Unification:** update both  $v, \mu$  simultaneously but at different rates  $\rho^v, \rho^\mu$ 

- $\rho^{v} < \rho^{\mu} \Rightarrow v$  evolves slowly  $\Rightarrow$  MFControl
- $\rho^v > \rho^\mu \Rightarrow \mu$  evolves slowly  $\Rightarrow$  MFGame

**Implementation:** Finite state space  $\mathcal{X}$  and finite action space  $\mathcal{A}$ , stationary problem **Q-learning:** Given  $\mu$ , **optimal** cost-to-go when starting at *x* using action *a* 

$$Q(x, a) = f(x, \mu, a) + \sum_{x' \in \mathcal{X}} p(x'|x, \mu, a) \underbrace{\min_{a'} Q(x', a')}_{=V(x')}.$$

Note: optimal control is  $\hat{v}_{Q}(x) = \operatorname{argmin}_{a} Q(x, a)$ . The scheme can be written as:  $\begin{cases} Q_{k+1} &= Q_{k} + \rho_{k}^{Q} \mathcal{T}(Q_{k}, \mu_{k}) \\ \mu_{k+1} &= \mu_{k} + \rho_{k}^{\mu} \mathcal{P}(Q_{k}, \mu_{k}), \end{cases}$ 

 $\text{ where } \begin{cases} \mathcal{T}(Q,\mu)(x,a) = f(x,a,\mu) + \gamma \sum_{x'} p(x'|x,a,\mu) \min_{a'} Q(x',a') - Q(x,a), \\ \mathcal{P}(Q,\mu)(x) = (\mu P^{Q,\mu})(x) - \mu(x), & \text{ with } P^{Q,\mu}(x,x') = p(x'|x,\hat{v}_Q(x),\mu) \end{cases}$ 

**MFControl:** Fix a control v, compute induced distribution  $\mu^v$ , update  $v, \ldots$ **MFGame:** Fix a distribution  $\mu$ , compute best response  $v^{\mu}$ , update  $\mu$ , ...

**Unification:** update both  $v, \mu$  simultaneously but at different rates  $\rho^v, \rho^\mu$ 

- $\rho^{v} < \rho^{\mu} \Rightarrow v$  evolves slowly  $\Rightarrow$  MFControl
- $\rho^v > \rho^\mu \Rightarrow \mu$  evolves slowly  $\Rightarrow$  MFGame

**Implementation:** Finite state space  $\mathcal{X}$  and finite action space  $\mathcal{A}$ , stationary problem **Q-learning:** Given  $\mu$ , **optimal** cost-to-go when starting at *x* using action *a* 

$$Q(x, a) = f(x, \mu, a) + \sum_{x' \in \mathcal{X}} p(x'|x, \mu, a) \underbrace{\min_{a'} Q(x', a')}_{=V(x')}.$$

Note: optimal control is  $\hat{v}_Q(x) = \operatorname{argmin}_a Q(x, a)$ . The scheme can be written as:  $\begin{cases} Q_{k+1} = Q_k + \rho_k^Q \mathcal{T}(Q_k, \mu_k) \\ \mu_{k+1} = \mu_k + \rho_k^\mu \mathcal{P}(Q_k, \mu_k), \end{cases}$ 

where  $\begin{cases} \mathcal{T}(Q,\mu)(x,a) = f(x,a,\mu) + \gamma \sum_{x'} p(x'|x,a,\mu) \min_{a'} Q(x',a') - Q(x,a), \\ \mathcal{P}(Q,\mu)(x) = (\mu P^{Q,\mu})(x) - \mu(x), & \text{with } P^{Q,\mu}(x,x') = p(x'|x,\hat{v}_Q(x),\mu) \end{cases}$ 

**Convergence:** based on Borkar's **two timescale** approach (includes sto. approx.) Rem.: For MFG only see e.g. [Mguni *et al.* [MJdC18], Subramanian *et al.* [SM19] Restricted environment: the agent needs to estimate the distribution



#### Restricted environment: the agent needs to estimate the distribution



#### Numerical illustration: Linear-quadratic example



## Summary

## References I

- [AFG17] Noha Almulla, Rita Ferreira, and Diogo Gomes, *Two numerical approaches to stationary mean-field games*, Dyn. Games Appl. 7 (2017), no. 4, 657–682. MR 3698446
- [AFL20] Andrea Angiuli, Jean-Pierre Fouque, and Mathieu Laurière, *Unified reinforcement q-learning for mean field game and control problems*, arXiv preprint arXiv:2006.13912 (2020).
- [AKS20] Berkay Anahtarci, Can Deha Kariksiz, and Naci Saldi, *Q-learning in regularized mean-field games*, arXiv preprint arXiv:2003.12151 (2020).
- [CFS15] René Carmona, Jean-Pierre Fouque, and Li-Hsien Sun, *Mean field games and systemic risk*, Commun. Math. Sci. **13** (2015), no. 4, 911–933. MR 3325083
- [CH17] Pierre Cardaliaguet and Saeed Hadikhanloo, Learning in mean field games: the fictitious play, ESAIM Control Optim. Calc. Var. 23 (2017), no. 2, 569–591. MR 3608094
- [EPL<sup>+</sup>20] Romuald Elie, Julien Perolat, Mathieu Laurière, Matthieu Geist, and Olivier Pietquin, On the convergence of model free learning in mean field games, Proceedings of the AAAI Conference on Artificial Intelligence, vol. 34, 2020, pp. 7143–7150.
- [FL<sup>+</sup>98] Drew Fudenberg, David K Levine, et al., *The theory of learning in games*, MIT Press Books **1** (1998).

## References II

- [FL09] Drew Fudenberg and David K Levine, *Learning and equilibrium*, Annu. Rev. Econ. 1 (2009), no. 1, 385–420.
- [GHXZ19] Xin Guo, Anran Hu, Renyuan Xu, and Junzi Zhang, *Learning mean-field games*, Advances in Neural Information Processing Systems **32** (2019), 4966–4976.
- [GPL+21] Matthieu Geist, Julien Pérolat, Mathieu Laurière, Romuald Elie, Sarah Perrin, Olivier Bachem, Rémi Munos, and Olivier Pietquin, *Concave utility reinforcement learning:* the mean-field game viewpoint, arXiv preprint arXiv:2106.03787 (2021).
- [Had18] Saeed Hadikhanloo, *Learning in mean field games*, Ph.D. thesis, PSL Research University, 2018.
- [HLMS21] Saeed Hadikhanloo, Rida Laraki, Panayotis Mertikopoulos, and Sylvain Sorin, Learning in nonatomic games, part i: Finite action spaces and population games, arXiv preprint arXiv:2107.01595 (2021).
- [HMC<sup>+</sup>06] Minyi Huang, Roland P Malhamé, Peter E Caines, et al., Large population stochastic dynamic games: closed-loop mckean-vlasov systems and the nash certainty equivalence principle, Communications in Information & Systems 6 (2006), no. 3, 221–252.
- [HZAL18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine, *Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor*, International conference on machine learning, PMLR, 2018, pp. 1861–1870.

## References III

- [LHP<sup>+</sup>16] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra, *Continuous control with deep* reinforcement learning., ICLR (Poster), 2016.
- [M<sup>+</sup>97] Tom M Mitchell et al., *Machine learning*.
- [MJdC18] David Mguni, Joel Jennings, and Enrique Munoz de Cote, *Decentralised learning in systems with many, many strategic agents*, Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- [MKS<sup>+</sup>13] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller, *Playing atari with deep reinforcement learning*, arXiv preprint arXiv:1312.5602 (2013).
- [NCM11] Mojtaba Nourian, Peter E Caines, and Roland P Malhamé, Mean field analysis of controlled cucker-smale type flocking: Linear analysis and perturbation equations, IFAC Proceedings Volumes 44 (2011), no. 1, 4471–4476.
- [PLP+21] Sarah Perrin, Mathieu Laurière, Julien Pérolat, Matthieu Geist, Romuald Élie, and Olivier Pietquin, *Mean field games flock! the reinforcement learning way*, arXiv preprint arXiv:2105.07933. Accepted to IJCA'21 (2021).
- [PPE<sup>+</sup>21] Julien Perolat, Sarah Perrin, Romuald Elie, Mathieu Laurière, Georgios Piliouras, Matthieu Geist, Karl Tuyls, and Olivier Pietquin, *Scaling up mean field games with* online mirror descent, arXiv preprint arXiv:2103.00623 (2021).

- [PPL+20] Sarah Perrin, Julien Pérolat, Mathieu Laurière, Matthieu Geist, Romuald Elie, and Olivier Pietquin, *Fictitious play for mean field games: Continuous time analysis and applications*, Advances in Neural Information Processing Systems (2020).
- [SM19] Jayakumar Subramanian and Aditya Mahajan, *Reinforcement learning in stationary mean-field games*, Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, 2019, pp. 251–259.

Unless otherwise specified, the images are from https://unsplash.com