

Mean Field Games:
Numerical Methods and
Applications in Machine Learning
Part 7: Mean Field Reinforcement Learning

Mathieu LAURIÈRE

<https://mlauriere.github.io/teaching/MFG-PKU-7.pdf>

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

RECAP

Outline

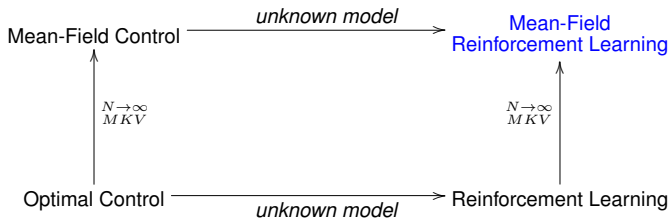
1. Introduction

2. Mean Field Reinforcement Learning

3. Model-Free Policy Gradient

4. Q-Learning

From Optimal Control to MFRL



- **Markov Decision Process (MDP):** $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where:
 - \mathcal{S} : state space, \mathcal{A} : action space,
 - $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$: transition kernel, $p(\cdot | s, a)$ gives next state's distribution
 - $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$: reward function, $\gamma \in (0, 1)$: discount factor
- **Goal:** Find (stationary, mixed) policy $\pi^* : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ maximizing:

$$R(\pi) = \mathbb{E} \left[\sum_{n \geq 0} \gamma^n r(s_n, a_n) \right], \quad \text{with } a_n \sim \pi(\cdot | s_n), s_{n+1} \sim p(\cdot | s_n, a_n)$$

- **Markov Decision Process (MDP):** $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where:
 - \mathcal{S} : state space, \mathcal{A} : action space,
 - $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$: transition kernel, $p(\cdot | s, a)$ gives next state's distribution
 - $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$: reward function, $\gamma \in (0, 1)$: discount factor
- **Goal:** Find (stationary, mixed) policy $\pi^* : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ maximizing:

$$R(\pi) = \mathbb{E} \left[\sum_{n \geq 0} \gamma^n r(s_n, a_n) \right], \quad \text{with } a_n \sim \pi(\cdot | s_n), s_{n+1} \sim p(\cdot | s_n, a_n)$$

- **Model:** p, r

- **Markov Decision Process (MDP):** $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where:
 - \mathcal{S} : state space, \mathcal{A} : action space,
 - $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$: transition kernel, $p(\cdot | s, a)$ gives next state's distribution
 - $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$: reward function, $\gamma \in (0, 1)$: discount factor
- **Goal:** Find (stationary, mixed) policy $\pi^* : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ maximizing:

$$R(\pi) = \mathbb{E} \left[\sum_{n \geq 0} \gamma^n r(s_n, a_n) \right], \quad \text{with } a_n \sim \pi(\cdot | s_n), s_{n+1} \sim p(\cdot | s_n, a_n)$$

- **Model:** p, r
- **Two settings:**
 - (1) **Known model** : **Optimal control** theory & methods
 - (2) **Sample transitions & rewards**: **Reinforcement Learning (RL)** framework

We want to **learn** the best control by performing **experiments** of the form:

Given the current state S_t ,

(1) Take an action A_t

(2) Observe reward R_{t+1} & new state S_{t+1}

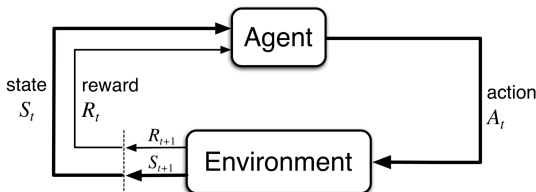
¹ Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.

We want to **learn** the best control by performing **experiments** of the form:

Given the current state S_t ,

(1) Take an action A_t

(2) Observe reward R_{t+1} & new state S_{t+1}



Source: [Sutton, Barto]¹

¹ Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.

- **Learning the policy:**

- ▶ Policy Gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \nabla J(\theta^{(k)}), \quad \pi^{(k)}(a|s) = \pi(s|a, \theta^{(k)})$$

- **Learning the policy:**

- ▶ Policy Gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \nabla J(\theta^{(k)}), \quad \pi^{(k)}(a|s) = \pi(s|a, \theta^{(k)})$$

- ▶ PPO, TRPO
- ▶ ...

- **Learning the policy:**

- ▶ Policy Gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \nabla J(\theta^{(k)}), \quad \pi^{(k)}(a|s) = \pi(s|a, \theta^{(k)})$$

- ▶ PPO, TRPO

- ▶ ...

- **Learning the value function:**

- ▶ Q-learning

$$Q^*(s, a) = r(s, a) + \gamma \max_{\pi} \mathbb{E}_{s' \sim p(\cdot|s, a), a' \sim \pi(\cdot|s')} \left[Q^*(s', a') \right]$$

Note: $V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$, $v^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$

● Learning the policy:

- ▶ Policy Gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \nabla J(\theta^{(k)}), \quad \pi^{(k)}(a|s) = \pi(s|a, \theta^{(k)})$$

- ▶ PPO, TRPO
- ▶ ...

● Learning the value function:

- ▶ Q-learning

$$Q^*(s, a) = r(s, a) + \gamma \max_{\pi} \mathbb{E}_{s' \sim p(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q^*(s', a')]$$

Note: $V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$, $v^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$

- ▶ Deep Q-neural network (DQN)
- ▶ ...

● Learning the policy:

- ▶ Policy Gradient

$$\theta^{(k+1)} = \theta^{(k)} - \eta^{(k)} \nabla J(\theta^{(k)}), \quad \pi^{(k)}(a|s) = \pi(s|a, \theta^{(k)})$$

- ▶ PPO, TRPO
- ▶ ...

● Learning the value function:

- ▶ Q-learning

$$Q^*(s, a) = r(s, a) + \gamma \max_{\pi} \mathbb{E}_{s' \sim p(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q^*(s', a')]$$

Note: $V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$, $v^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$

- ▶ Deep Q-neural network (DQN)
- ▶ ...

● Hybrid:

- ▶ Deep Deterministic Policy Gradient (DDPG)
- ▶ Soft Actor Critic (SAC)
- ▶ ...

Outline

1. Introduction

2. Mean Field Reinforcement Learning

3. Model-Free Policy Gradient

4. Q-Learning

● Dynamics: discrete time

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

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

- **Dynamics:** discrete time

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

- ▶ $X_n^{\alpha, \mu} \in \mathcal{X} \subseteq \mathbb{R}^d$: 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.)**
- ▶ $p(x'|x, \alpha, \mu)$: corresponding transition probability distribution
- ▶ $\mu_n \in \mathcal{P}(\mathcal{X} \times \mathcal{A})$: a **state-action distribution**
- ▶ π_n : a policy; randomized **actions**: $\alpha_n \sim \pi_n(\cdot | s_n, \mu_n)$

- **Cost:** $\mathbb{J}(\pi; \mu) = \mathbb{E}_{\epsilon, \epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n f(X_n^{\alpha, \mu}, \alpha_n, \mu_n) \right]$

- **Two scenarios:**

- ▶ **Cooperative (MFC):** Find π^* s.t.

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

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

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

- **Key Remark:**

$$\alpha^* \in \underset{\alpha}{\operatorname{argmin}} J^{MFC}(\alpha) = \mathbb{E}_{\epsilon, \epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n f(X_n^\alpha, \alpha_n, \mu_n^\pi) \right], \quad \mu_n^\pi = \mathbb{P}_{X_n^\alpha}^0$$

● **Key Remark:**

$$\begin{aligned}
 \alpha^* \in \underset{\alpha}{\operatorname{argmin}} J^{MFC}(\alpha) &= \mathbb{E}_{\epsilon, \epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n f(X_n^\alpha, \alpha_n, \mu_n^\pi) \right], & \mu_n^\pi &= \mathbb{P}_{X_n^\alpha}^0 \\
 &= \mathbb{E}_{\epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n \underbrace{\int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu_n^\pi) \nu_n^\pi(dx, da)}_{\text{function of } \nu_n^\pi} \right]
 \end{aligned}$$

● **Key Remark:**

$$\begin{aligned} \alpha^* \in \underset{\alpha}{\operatorname{argmin}} J^{MFC}(\alpha) &= \mathbb{E}_{\epsilon, \epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n f(X_n^\alpha, \alpha_n, \mu_n^\pi) \right], & \mu_n^\pi &= \mathbb{P}_{X_n^\alpha}^0 \\ &= \mathbb{E}_{\epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n \underbrace{\int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu_n^\pi) \nu_n^\pi(dx, da)}_{\text{function of } \nu_n^\pi} \right] \end{aligned}$$

● **Lifted problem:** population / social planner's optimization problem:

- state = population distribution μ_n^π
- value function = function of the distribution μ

● **Key Remark:**

$$\begin{aligned}\alpha^* \in \underset{\alpha}{\operatorname{argmin}} J^{MFC}(\alpha) &= \mathbb{E}_{\epsilon, \epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n f(X_n^\alpha, \alpha_n, \mu_n^\pi) \right], & \mu_n^\pi &= \mathbb{P}_{X_n^\alpha}^0 \\ &= \mathbb{E}_{\epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n \underbrace{\int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu_n^\pi) \nu_n^\pi(dx, da)}_{\text{function of } \nu_n^\pi} \right]\end{aligned}$$

● **Lifted problem:** population / social planner's optimization problem:

→ state = population distribution μ_n^π

→ value function = function of the distribution μ

● **Mean Field Markov Decision Process (MFMDP):** $(\bar{\mathcal{S}}, \bar{\mathcal{A}}, \bar{p}, \bar{r}, \gamma)$, where:

- State space: $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$
- Action space: $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$ with constraint: $pr_1(\bar{a}) = \mu$
- Transition function: $\mu' = \bar{F}(\mu, \bar{a}, \epsilon^0) \sim \bar{p}(\mu, \bar{a})$
- Reward function: $\bar{r}(\mu, \bar{a}) = - \int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu) \bar{a}(dx, da)$

● **Key Remark:**

$$\begin{aligned} \alpha^* \in \underset{\alpha}{\operatorname{argmin}} J^{MFC}(\alpha) &= \mathbb{E}_{\epsilon, \epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n f(X_n^\alpha, \alpha_n, \mu_n^\pi) \right], \quad \mu_n^\pi = \mathbb{P}_{X_n^\alpha}^0 \\ &= \mathbb{E}_{\epsilon^0} \left[\sum_{n=0}^{\infty} \gamma^n \underbrace{\int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu_n^\pi) \nu_n^\pi(dx, da)}_{\text{function of } \nu_n^\pi} \right] \end{aligned}$$

● **Lifted problem:** population / social planner's optimization problem:

→ state = population distribution μ_n^π

→ value function = function of the distribution μ

● **Mean Field Markov Decision Process (MFMDP):** $(\bar{\mathcal{S}}, \bar{\mathcal{A}}, \bar{p}, \bar{r}, \gamma)$, where:

• State space: $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$

• Action space: $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$ with constraint: $pr_1(\bar{a}) = \mu$

• Transition function: $\mu' = \bar{F}(\mu, \bar{a}, \epsilon^0) \sim \bar{p}(\mu, \bar{a})$

• Reward function: $\bar{r}(\mu, \bar{a}) = - \int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu) \bar{a}(dx, da)$

● **Goal:** $\max. \bar{J}^{\bar{\pi}}(\mu) = \mathbb{E} \left[\sum_{n=0}^{\infty} \gamma^n \bar{r}(\mu_n^{\bar{\pi}}, \bar{a}_n) \right], \bar{a}_n \sim \bar{\pi}(\cdot | \mu_n^{\bar{\pi}}), \mu_{n+1}^{\bar{\pi}} \sim \bar{p}(\cdot | \mu_n^{\bar{\pi}}, \bar{a}_n),$
 $\mu_0^{\bar{\pi}} = \mu$

● **Mean field policy:** $\bar{\pi}$ kernel $\bar{\mathcal{S}} \rightarrow \mathcal{P}(\bar{\mathcal{A}})$, randomized population-strategies \bar{a}

Theorem: DPP for MFMDP

[Carmona, L., Tan'21]²

Under suitable conditions,

$$\bar{J}^*(\mu) := \sup_{\bar{\pi}} \bar{J}^{\bar{\pi}}(\mu) = \sup_{\bar{\pi}} \left\{ \int_{\bar{\mathcal{A}}} \left[\bar{r}(\mu, \bar{a}) + \gamma \mathbb{E} \left[\bar{J}^* \left(\bar{F}(\mu, \bar{a}, \epsilon^0) \right) \right] \right] \bar{\pi}(d\bar{a} | \mu) \right\},$$

where the sup is over a subset of $\{\bar{\pi} : \bar{\mathcal{S}} \rightarrow \mathcal{P}(\bar{\mathcal{A}})\}$

Likewise for **mean field state-action value function** \bar{Q}^*

²Carmona, R., Laurière, M., & Tan, Z. (2019). Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. arXiv preprint arXiv:1910.12802. (Preliminary version. Update coming soon!)

Theorem: DPP for MFMDP

[Carmona, L., Tan'21]²

Under suitable conditions,

$$\bar{J}^*(\mu) := \sup_{\bar{\pi}} \bar{J}^{\bar{\pi}}(\mu) = \sup_{\bar{\pi}} \left\{ \int_{\bar{\mathcal{A}}} \left[\bar{r}(\mu, \bar{a}) + \gamma \mathbb{E} \left[\bar{J}^* \left(\bar{F}(\mu, \bar{a}, \epsilon^0) \right) \right] \right] \bar{\pi}(d\bar{a} | \mu) \right\},$$

where the sup is over a subset of $\{\bar{\pi} : \bar{\mathcal{S}} \rightarrow \mathcal{P}(\bar{\mathcal{A}})\}$

Likewise for **mean field state-action value function** \bar{Q}^*

Proof: based on “double lifting” [Bertsekas, Shreve'78]

²Carmona, R., Laurière, M., & Tan, Z. (2019). Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. arXiv preprint arXiv:1910.12802. (Preliminary version. Update coming soon!)

Theorem: DPP for MFMDP

[Carmona, L., Tan'21]²

Under suitable conditions,

$$\bar{J}^*(\mu) := \sup_{\bar{\pi}} \bar{J}^{\bar{\pi}}(\mu) = \sup_{\bar{\pi}} \left\{ \int_{\bar{\mathcal{A}}} \left[\bar{r}(\mu, \bar{a}) + \gamma \mathbb{E} \left[\bar{J}^* \left(\bar{F}(\mu, \bar{a}, \epsilon^0) \right) \right] \right] \bar{\pi}(d\bar{a} | \mu) \right\},$$

where the sup is over a subset of $\{\bar{\pi} : \bar{\mathcal{S}} \rightarrow \mathcal{P}(\bar{\mathcal{A}})\}$

Likewise for **mean field state-action value function** \bar{Q}^*

Proof: based on “double lifting” [Bertsekas, Shreve'78]

DPPs for MFC:

[L., Pironneau; Pham, *et al.*; Gast *et al.*; Guo *et al.*; Possamai *et al.*; ...]

²Carmona, R., Laurière, M., & Tan, Z. (2019). Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. arXiv preprint arXiv:1910.12802. (Preliminary version. Update coming soon!)

Theorem: DPP for MFMDP

[Carmona, L., Tan'21]²

Under suitable conditions,

$$\bar{J}^*(\mu) := \sup_{\bar{\pi}} \bar{J}^{\bar{\pi}}(\mu) = \sup_{\bar{\pi}} \left\{ \int_{\bar{\mathcal{A}}} \left[\bar{r}(\mu, \bar{a}) + \gamma \mathbb{E} \left[\bar{J}^* \left(\bar{F}(\mu, \bar{a}, \epsilon^0) \right) \right] \right] \bar{\pi}(d\bar{a} | \mu) \right\},$$

where the sup is over a subset of $\{\bar{\pi} : \bar{\mathcal{S}} \rightarrow \mathcal{P}(\bar{\mathcal{A}})\}$

Likewise for **mean field state-action value function** \bar{Q}^*

Proof: based on “double lifting” [Bertsekas, Shreve'78]

DPPs for MFC:

[L., Pironneau; Pham, *et al.*; Gast *et al.*; Guo *et al.*; Possamai *et al.*; ...]

Here: discrete time, infinite horizon, **common noise**, **feedback controls**, ...

→ well-suited for **RL**

→ Mean-field Q-learning algorithm

²Carmona, R., Laurière, M., & Tan, Z. (2019). Model-free mean-field reinforcement learning: mean-field MDP and mean-field Q-learning. arXiv preprint arXiv:1910.12802. (Preliminary version. Update coming soon!)

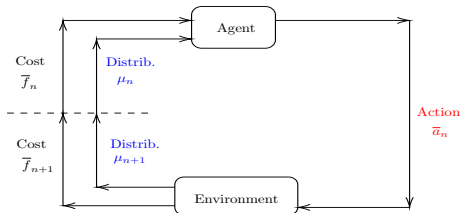
Hierarchy of settings:

- **Setting 1: known model:** computational method based on knowledge of MFMDP
 - (a) Gradient based methods
 - (b) Dynamic programming based methods

Mean Field Learning Settings

Hierarchy of settings:

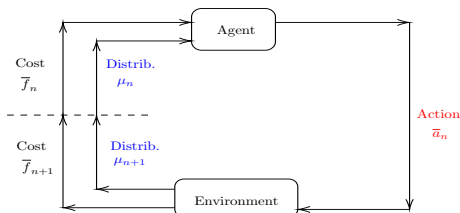
- **Setting 1: known model:** computational method based on knowledge of MFMDP
 - (a) Gradient based methods
 - (b) Dynamic programming based methods
- **Setting 2:** unknown model but **samples from MFMDP:** MF learning



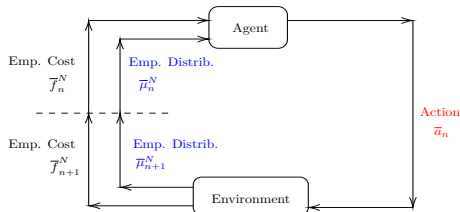
Mean Field Learning Settings

Hierarchy of settings:

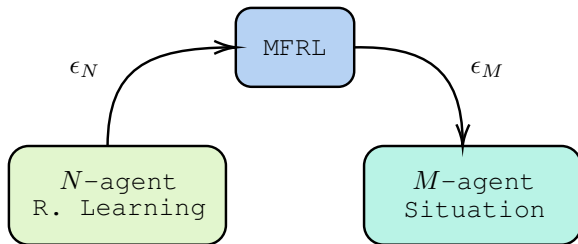
- **Setting 1: known model:** computational method based on knowledge of MFMDP
 - (a) Gradient based methods
 - (b) Dynamic programming based methods
- **Setting 2:** unknown model but **samples from MFMDP:** MF learning



- **Setting 3:** unknown model but **samples from N -agent MDP:** approx. MF learning



Mean Field Control: Finite Population Approximation



Outline

1. Introduction
2. Mean Field Reinforcement Learning
- 3. Model-Free Policy Gradient**
4. Q-Learning

Idea 1: *Make the “policy gradient” approach (see Part 5 of lecture slides) model-free*

Policy Gradient (PG) to minimize $J(\theta)$

- Control \approx **parameterized function**
- Look for the optimal parameter θ^*
- Perform **gradient descent** on the space of parameters

Idea 1: *Make the “policy gradient” approach (see Part 5 of lecture slides) model-free*

Policy Gradient (PG) to minimize $J(\theta)$

- Control \approx **parameterized function**
- Look for the optimal parameter θ^*
- Perform **gradient descent** on the space of parameters

Hierarchy of three situations, more and more complex:

(1) access to the exact **(mean field) model**:

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla J(\theta^{(k)})$$

Idea 1: *Make the “policy gradient” approach (see Part 5 of lecture slides) model-free*

Policy Gradient (PG) to minimize $J(\theta)$

- Control \approx **parameterized function**
- Look for the optimal parameter θ^*
- Perform **gradient descent** on the space of parameters

Hierarchy of three situations, more and more complex:

(1) access to the exact **(mean field) model**:

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla J(\theta^{(k)})$$

(2) access to a **mean field simulator**:

→ idem + **gradient estimation** (0^{th} -order opt.):

$$\theta^{(k+1)} = \theta^{(k)} - \eta \widetilde{\nabla} J(\theta^{(k)})$$

Idea 1: Make the “policy gradient” approach (see Part 5 of lecture slides) model-free

Policy Gradient (PG) to minimize $J(\theta)$

- Control \approx **parameterized function**
- Look for the optimal parameter θ^*
- Perform **gradient descent** on the space of parameters

Hierarchy of three situations, more and more complex:

- (1) access to the exact **(mean field) model**: $\theta^{(k+1)} = \theta^{(k)} - \eta \nabla J(\theta^{(k)})$
- (2) access to a **mean field simulator**:
→ idem + **gradient estimation** (0^{th} -order opt.): $\theta^{(k+1)} = \theta^{(k)} - \eta \widetilde{\nabla} J(\theta^{(k)})$
- (3) access to a N -agent **population simulator**:
→ idem + error on **mean \approx empirical mean (LLN)**: $\theta^{(k+1)} = \theta^{(k)} - \eta \widetilde{\nabla}^N J(\theta^{(k)})$

Idea 1: Make the “policy gradient” approach (see Part 5 of lecture slides) model-free

Policy Gradient (PG) to minimize $J(\theta)$

- Control \approx **parameterized function**
- Look for the optimal parameter θ^*
- Perform **gradient descent** on the space of parameters

Hierarchy of three situations, more and more complex:

- (1) access to the exact **(mean field) model**: $\theta^{(k+1)} = \theta^{(k)} - \eta \nabla J(\theta^{(k)})$
- (2) access to a **mean field simulator**:
→ idem + **gradient estimation** (0^{th} -order opt.): $\theta^{(k+1)} = \theta^{(k)} - \eta \widetilde{\nabla} J(\theta^{(k)})$
- (3) access to a N -agent **population simulator**:
→ idem + error on **mean \approx empirical mean (LLN)**: $\theta^{(k+1)} = \theta^{(k)} - \eta \widetilde{\nabla}^N J(\theta^{(k)})$

Theorem: For **Linear-Quadratic MFC**

[Carmona, L., Tan'19]³

In each case, convergence holds at a linear rate:

Taking $k \approx \mathcal{O}(\log(1/\epsilon))$ is sufficient to ensure $J(\theta^{(k)}) - J(\theta^*) < \epsilon$.

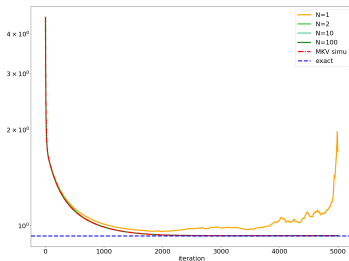
Proof: builds on [Fazel et al.'18], analysis of perturbation of Riccati equations

³Carmona, R., Laurière, M., & Tan, Z. (2019). Linear-quadratic mean-field reinforcement learning: convergence of policy gradient methods. arXiv preprint arXiv:1910.04295.

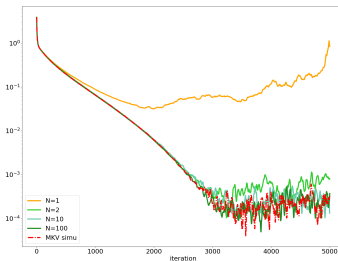
Numerical Illustration

Example: Linear dynamics, quadratic costs of the type:

$$f(x, \mu, v) = \underbrace{(\bar{\mu} - x)^2}_{\text{distance to mean position}} + \underbrace{v^2}_{\text{cost of moving}}, \quad \bar{\mu} = \underbrace{\int \mu(\xi) d\xi}_{\text{mean position}},$$



Value of the MF cost

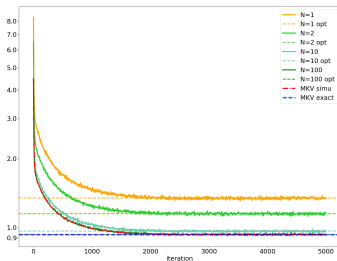


Rel. err. on MF cost

MF cost = cost in the mean field problem

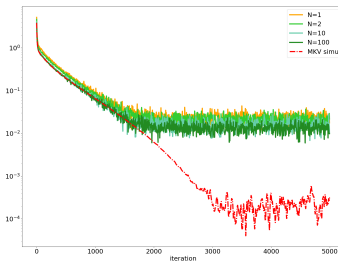
Example: Linear dynamics, quadratic costs of the type:

$$f(x, \mu, v) = \underbrace{(\bar{\mu} - x)^2}_{\text{distance to mean position}} + \underbrace{v^2}_{\text{cost of moving}}, \quad \bar{\mu} = \underbrace{\int \mu(\xi) d\xi}_{\text{mean position}},$$



Value of the social cost

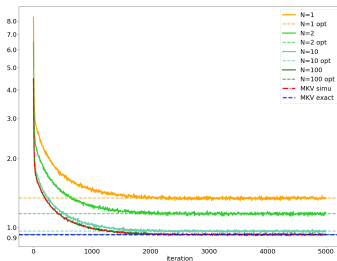
Social cost = average over the N -agents



Rel. err. on social cost

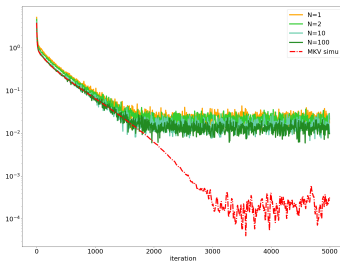
Example: Linear dynamics, quadratic costs of the type:

$$f(x, \mu, v) = \underbrace{(\bar{\mu} - x)^2}_{\text{distance to mean position}} + \underbrace{v^2}_{\text{cost of moving}}, \quad \bar{\mu} = \underbrace{\int \mu(\xi) d\xi}_{\text{mean position}},$$



Value of the social cost

Social cost = average over the N -agents



Rel. err. on social cost

Main take-away:

Trying to learn the mean-field regime solution can be efficient even for N small

Outline

1. Introduction
2. Mean Field Reinforcement Learning
3. Model-Free Policy Gradient
4. Q-Learning

Idea 2: Generalize Q-learning to Mean-Field Control

Reminder:

● **Mean Field Markov Decision Process (MFMDP):** $(\bar{\mathcal{S}}, \bar{\mathcal{A}}, \bar{p}, \bar{r}, \gamma)$, where:

- State space: $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$
- Action space: $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$ with constraint: $pr_1(\bar{a}) = \mu$
- Transition function: $\mu' = \bar{F}(\mu, \bar{a}, \epsilon^0) \sim \bar{p}(\mu, \bar{a})$
- Reward function: $\bar{r}(\mu, \bar{a}) = - \int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu) \bar{a}(dx, da)$

● **Goal:** $\max. \bar{J}^{\bar{\pi}}(\mu) = \mathbb{E} \left[\sum_{n=0}^{\infty} \gamma^n \bar{r}(\mu_n^{\bar{\pi}}, \bar{a}_n) \right]$, $\bar{a}_n \sim \bar{\pi}(\cdot | \mu_n^{\bar{\pi}})$, $\mu_{n+1}^{\bar{\pi}} \sim \bar{p}(\cdot | \mu_n^{\bar{\pi}}, \bar{a}_n)$,
 $\mu_0^{\bar{\pi}} = \mu$

Mean Field Q-Function

Idea 2: Generalize Q-learning to Mean-Field Control

Reminder:

● **Mean Field Markov Decision Process (MFMDP):** $(\bar{\mathcal{S}}, \bar{\mathcal{A}}, \bar{p}, \bar{r}, \gamma)$, where:

- State space: $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$
- Action space: $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$ with constraint: $pr_1(\bar{a}) = \mu$
- Transition function: $\mu' = \bar{F}(\mu, \bar{a}, \epsilon^0) \sim \bar{p}(\mu, \bar{a})$
- Reward function: $\bar{r}(\mu, \bar{a}) = - \int_{\mathcal{X} \times \mathcal{U}} f(x, a, \mu) \bar{a}(dx, da)$

● **Goal:** max. $\bar{J}^{\bar{\pi}}(\mu) = \mathbb{E} \left[\sum_{n=0}^{\infty} \gamma^n \bar{r}(\mu_n^{\bar{\pi}}, \bar{a}_n) \right]$, $\bar{a}_n \sim \bar{\pi}(\cdot | \mu_n^{\bar{\pi}})$, $\mu_{n+1}^{\bar{\pi}} \sim \bar{p}(\cdot | \mu_n^{\bar{\pi}}, \bar{a}_n)$,
 $\mu_0^{\bar{\pi}} = \mu$

Q-function associated to a policy π :

$$Q^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a), a' \sim \pi(\cdot | s')} \left[Q^{\pi}(s', a') \right]$$

Mean Field Q-function associated to a mean field policy $\bar{\pi}$:

$$\bar{Q}^{\bar{\pi}}(\bar{s}, \bar{a}) = \bar{r}(\bar{s}, \bar{a}) + \gamma \mathbb{E}_{\bar{s}' \sim \bar{p}(\cdot | \bar{s}, \bar{a}), \bar{a}' \sim \bar{\pi}(\cdot | \bar{s}')} \left[\bar{Q}^{\bar{\pi}}(\bar{s}', \bar{a}') \right]$$

- **Optimal MF Q-function:**

$$\bar{Q}^*(\bar{s}, \bar{a}) = \bar{r}(\bar{s}, \bar{a}) + \gamma \sup_{\bar{\pi}} \mathbb{E}_{\bar{a}' \sim \bar{\pi}(\cdot | \bar{s}), \bar{s}' \sim \bar{p}(\cdot | \bar{s}, \bar{a}')} \left[\bar{Q}^*(\bar{s}', \bar{a}') \right]$$

- **Algorithm:**

- Idealized version (synchronous):

$$\begin{aligned} \bar{Q}^{(k+1)}(\bar{s}, \bar{a}) &= \bar{r}(\bar{s}, \bar{a}) + \gamma \sup_{\bar{\pi}} \mathbb{E}_{\bar{s}' \sim \bar{p}(\cdot | \bar{s}, \bar{a}), \bar{a}' \sim \bar{\pi}(\cdot | \bar{s}')} \left[\bar{Q}^{(k)}(\bar{s}', \bar{a}') \right], \quad (\bar{s}, \bar{a}) \in \bar{\mathcal{S}} \times \bar{\mathcal{A}} \\ &= [\bar{T}^* \bar{Q}^{(k)}](\bar{s}, \bar{a}) \end{aligned}$$

- Following a trajectory (async.): $\bar{s}^{(k+1)} \sim p(\cdot | \bar{s}^{(k)}, \bar{a}^{(k)}), \bar{a}^{(k+1)} \sim \bar{\pi}^{(k+1)}(\cdot | \bar{s}^{(k)})$,

$$\begin{cases} \bar{Q}^{(k+1)}(\bar{s}, \bar{a}) = \bar{Q}^{(k)}(\bar{s}, \bar{a}), & (\bar{s}, \bar{a}) \in \bar{\mathcal{S}} \times \bar{\mathcal{A}} \\ \bar{Q}^{(k+1)}(\bar{s}^{(k+1)}, \bar{a}^{(k+1)}) \leftarrow \bar{r}(\bar{s}^{(k+1)}, \bar{a}^{(k+1)}) + \gamma \max_{\bar{a}'} \bar{Q}^{(k)}(\bar{s}^{(k+1)}, \bar{a}') \end{cases}$$

- **Implementation:** several possibilities (can be combined):

- ▶ pure (population and individual) strategies
- ▶ discretization of $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X}), \bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$
- ▶ deep Reinforcement Learning

Cyber-security example of [\[Bensoussan, Kolokoltsov\]](#) (see Part 6 of slides)

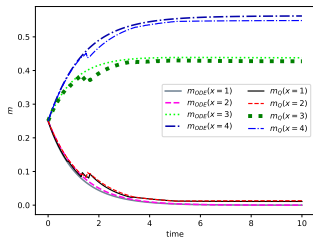
- MFC viewpoint, MF Q-learning
- pure (population and individual) strategies
- discretization of $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$, $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$

MF Q-Learning: Numerical Illustration

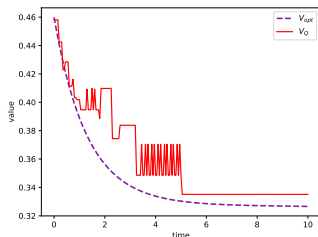
Cyber-security example of [Bensoussan, Kolokoltsov] (see Part 6 of slides)

- MFC viewpoint, MF Q-learning
- pure (population and individual) strategies
- discretization of $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$, $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$

Test 1: $m_0 = (1/4, 1/4, 1/4, 1/4)$



Evolution of m^m0 optimally controlled (m_{ODE}) or controlled using the approximate Q -function (m_Q)



V function (V_{opt}) and approximate Q -function (V_Q) along the optimal flow.

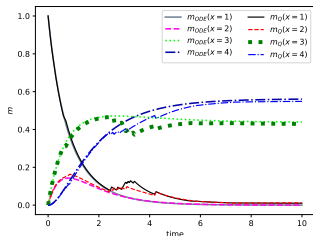
(More details in [L'21 - AMS notes])

MF Q-Learning: Numerical Illustration

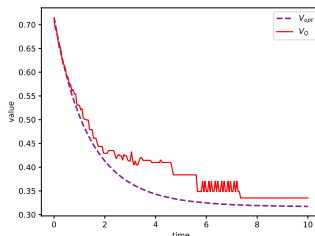
Cyber-security example of [Bensoussan, Kolokoltsov] (see Part 6 of slides)

- MFC viewpoint, MF Q-learning
- pure (population and individual) strategies
- discretization of $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$, $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$

Test 2: $m_0 = (1, 0, 0, 0)$



Evolution of m^{m_0} optimally controlled (m_{ODE}) or controlled using the approximate Q -function (m_Q)



V function (V_{opt}) and approximate Q -function (V_Q) along the optimal flow.

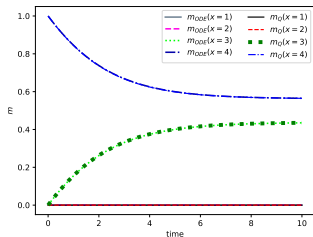
(More details in [L'21 - AMS notes])

MF Q-Learning: Numerical Illustration

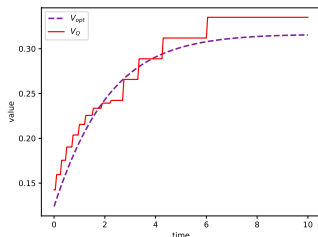
Cyber-security example of [Bensoussan, Kolokoltsov] (see Part 6 of slides)

- MFC viewpoint, MF Q-learning
- pure (population and individual) strategies
- discretization of $\bar{\mathcal{S}} = \mathcal{P}(\mathcal{X})$, $\bar{\mathcal{A}} = \mathcal{P}(\mathcal{X} \times \mathcal{U})$

Test 3: $m_0 = (0, 0, 0, 1)$



Evolution of m^{m_0} optimally controlled (m_{ODE}) or controlled using the approximate Q -function (m_Q)



V function (V_{opt}) and approximate Q -function (V_Q) along the optimal flow.

(More details in [L'21 - AMS notes])

