

# Learning Equilibria with Bandit Feedback

Maryam Kamgarpour  
École Polytechnique Fédérale de Lausanne, Switzerland

**Workshop on Learning in Games, Toulouse, France**

02.07.2024

The EPFL logo consists of the letters "EPFL" in white, bold, sans-serif font, centered within a solid red rectangular background.The sycamore lab logo features the text "sycamore lab" in a black, lowercase, sans-serif font. A green circular icon with a stylized leaf is positioned between the words "sycamore" and "lab". Below the main text, the words "SYSTEMS CONTROL AND MULTIAGENT OPTIMIZATION RESEARCH" are written in a smaller, green, uppercase, sans-serif font.

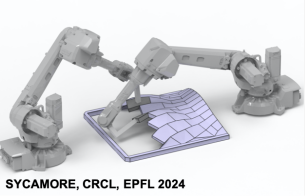
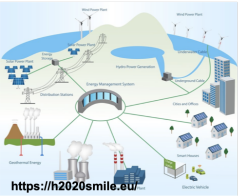
SYSTEMS CONTROL AND MULTIAGENT OPTIMIZATION RESEARCH

# Background - Control systems

From ...



to ...



## Problem of interest - Learning in games

Player  $i$  does not know  $J^i$  but can query it



How do players learn to optimize their decisions?

## Introduction

## Learning Nash equilibria

- Normal form games

- Markov games

## No-regret learning

- Normal form games

- Markov games

## Conclusions



# Outline

Introduction

Learning Nash equilibria

Normal form games

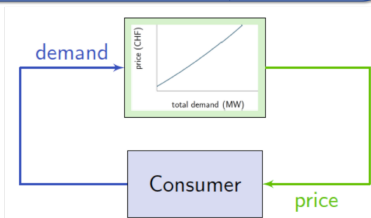
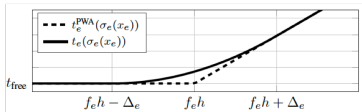
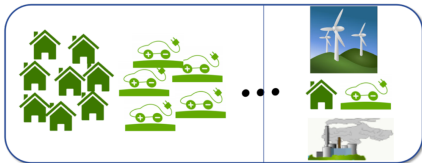
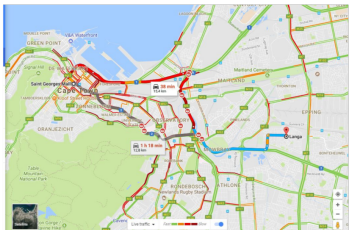
Markov games

No-regret learning

Conclusions

# Convex games

- ▶  $J^i(a^i, a^{-i})$ : convex in  $a^i$ , continuously differentiable
- ▶  $a^i \in A^i \subset \mathbb{R}^d$ : convex and compact
- ▶ Examples
  - ▶ mixed strategy extension of a finite action game
  - ▶ traffic networks, electricity market



## Nash equilibrium as a desirable solution outcome

$\mathbf{a}^* = (a^{*1}, a^{*2}, \dots, a^{*N})$  is a Nash equilibrium if for every player  $i$

$$J^i(a^{*i}, a^{*-i}) = \min_{a^i} J^i(a^i, a^{*-i})$$

► characterized by the pseudo-gradient:  $\mathbf{M} : \mathbb{R}^{Nd} \rightarrow \mathbb{R}^{Nd}$

$$\mathbf{M}(\mathbf{a}) = [\nabla_i J^i(a^i, a^{-i})]_{i=1}^N$$

$\mathbf{a}^*$  Nash equilibrium  $\iff \mathbf{M}(\mathbf{a}^*)^T (\mathbf{a} - \mathbf{a}^*) \geq 0, \forall \mathbf{a} \in \mathbf{A}$

## Learning in convex games

Player  $i$  does not know  $J^i$  but can query it



Independent payoff-based approach:

$$\theta_{t+1}^i = \text{Proj}_{A^i}(\theta_t^i - \eta_t \nabla_{\theta^i} \widehat{J^i}(\theta_t^i, \theta_t^{-i}))$$

# Learning in convex games

Player  $i$  does not know  $J^i$  but can query it



Independent payoff-based approach:

$$\theta_{t+1}^i = \text{Proj}_{A^i}(\theta_t^i - \eta_t \nabla_{\theta^i} \widehat{J^i}(\theta_t^i, \theta_t^{-i}))$$

Challenges compared to the single agent setting:

1. How can agent  $i$  estimate  $\nabla_{\theta^i} J^i(\theta)$  without knowing  $\theta$ ?
2. Under which conditions do we have convergence?

# Independent estimation of local gradients

- ▶ Finite difference:  $\widehat{\nabla_{\theta^i} J^i(\boldsymbol{\theta})} \approx \frac{J^i(\theta^i, \boldsymbol{\theta}^{-i}) - J^i(\theta^i + \delta, \boldsymbol{\theta}^{-i})}{\delta}$ 
  - ▶ requires others to stay with their action  $\implies$  coordination

- ▶ approach: randomize query  $\delta^i \sim \mathcal{N}(0, \sigma^2)$



- ▶ Construct  $\widehat{\nabla_{\theta^i} J^i(\boldsymbol{\theta})}$  with one function evaluation
  - ▶ bias:  $O(\sigma)$ , variance  $O(\frac{1}{\sigma^2})$  [Nesterov, Spokoiny 2019]

Alternatively, uniform distribution sampling [Flaxman et al. 2004]

## The game pseudo-gradient

Consider known gradients, unconstrained. Learning dynamics:

$$\begin{bmatrix} \theta_{t+1}^1 \\ \vdots \\ \theta_{t+1}^N \end{bmatrix} = \begin{bmatrix} \theta_t^1 \\ \vdots \\ \theta_t^N \end{bmatrix} - \eta_t \underbrace{\begin{bmatrix} \nabla_{\theta^1} J^1(\boldsymbol{\theta}_t) \\ \vdots \\ \nabla_{\theta^N} J^N(\boldsymbol{\theta}_t) \end{bmatrix}}_{\neq \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})}$$

- ▶ ex:  $J^1(\boldsymbol{\theta}) = \theta^1 \theta^2 = -J^2(\boldsymbol{\theta})$ ,  $\begin{bmatrix} \nabla_{\theta^1} J^1(\boldsymbol{\theta}) \\ \nabla_{\theta^2} J^2(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} \theta^1 \\ \theta^2 \end{bmatrix}$
- ▶ single agent analysis don't generally work

## Sufficient conditions for convergence

$\mathbf{a}^*$  is strongly variationally stable:  $\exists \nu > 0$ :

$$\mathbf{M}(\mathbf{a})^T(\mathbf{a} - \mathbf{a}^*) > \nu \|\mathbf{a} - \mathbf{a}^*\|^2, \forall \mathbf{a} \in \mathbf{A}$$

- ▶ example  $J^i(\mathbf{a}) = a^1 a^2 a^3 + (a^i)^2, a^i \in [-1, 2], i \in \{1, 2, 3\}$



## Sufficient conditions for convergence

$\mathbf{a}^*$  is strongly variationally stable:  $\exists \nu > 0$ :

$$\mathbf{M}(\mathbf{a})^T(\mathbf{a} - \mathbf{a}^*) > \nu \|\mathbf{a} - \mathbf{a}^*\|^2, \forall \mathbf{a} \in \mathbf{A}$$

► example  $J^i(\mathbf{a}) = a^1 a^2 a^3 + (a^i)^2, a^i \in [-1, 2], i \in \{1, 2, 3\}$

Algorithm:  $\boldsymbol{\theta}_{t+1} = \text{Proj}_{\mathbf{A}}(\boldsymbol{\theta}_t - \eta_t \hat{\mathbf{M}}(\boldsymbol{\theta}_t))$

### Theorem

Assume  $\mathbf{M}$  Lipschitz and  $\mathbf{a}^*$  strongly VS. For  $\sum_t \eta_t = \infty$ ,  
 $\sum_t \frac{\eta_t^2}{\sigma_t^2} < \infty$ ,  $\boldsymbol{\theta}_t$  converges almost surely to  $\mathbf{a}^*$

# Payoff-based learning leverages pseudo-gradient properties

## Recent progress

- ▶ Mere monotonicity of  $M(\mathbf{a}) \supseteq$  zero-sum matrix games:  
extra-gradient, optimistic gradient descent-ascent, Tikhonov regularization, ...
- ▶ Local variational stability  $\implies$  local convergence
- ▶ Convergence rates

[Tatarenko, MK, IEEE TAC 2019, IEEE TCNS 2024, ECC 2024]

[Bravo et al., 2018], [Mertikopoulos et al. 2018], [Gao, Pavel, 2022], ...

**Challenge: many games including Markov games do not satisfy above conditions**

# Markov games

- ▶ Dynamics:  $s_{h+1} \sim P(\cdot | s_h, a_h^1, \dots, a_h^N)$
- ▶ Policy  $\pi^i : S \rightarrow \Delta(A^i)$
- ▶  $V_s^i(\pi^i, \pi^{-i}) = \mathbb{E}_{P, \pi} \sum_{h=0}^{\infty} \gamma^t R^i(s_h, \pi^1(s_h), \dots, \pi^N(s_h))$
- ▶ Nash equilibrium:

$$V_s^i(\pi^{*i}, \pi^{*-i}) \geq V_s^i(\pi^i, \pi^{*-i}), \forall \pi^i, \forall i$$

note change of notation: from costs to rewards and value function for players

# Multiagent reinforcement learning approach

Given  $s_{h+1} \sim P(\cdot | s_h, a_h^1, \dots, a_h^N)$

- ▶ Parametrize a policy  $a_t^i \sim \pi_{\theta^i}(\cdot | s_t)$ ,  $\theta^i \in \mathbb{R}^d$
- ▶ Find equilibrium  $\theta^* = (\theta^1, \dots, \theta^N)$  by interacting with the system



# Policy gradient class of algorithms

Single agent RL:  $V(\theta) = \mathbb{E}_{P,\pi} \sum_{h=0}^{\infty} \gamma^t R(s_h, \pi(s_h))$

$$\theta_{t+1} = \theta_t - \eta_t \nabla_{\theta} V(\theta_t)$$

- ▶ convergence under *gradient dominance* condition [Agarwal et al., 2021],[Hu et al. 2023], [Bhandari et al. 2024], ...

Multiagent RL:  $V^i(\theta^i, \theta^{-i}) = \mathbb{E}_{P,\pi} \sum_{h=0}^{\infty} \gamma^t R^i(s_h, \pi^1(s_h), \dots, \pi^N(s_h))$

$$\theta_{t+1}^i = \theta_t^i - \eta_t \nabla_{\theta^i} V^i(\theta_t^i, \theta_t^{-i})$$

- ▶ generally non-convergent

# Challenging even in linear quadratic setting

single agent

$$J(\theta) = \mathbb{E}_{s_0} \left[ \sum_{h=0}^{\infty} s_h^T Q s_h + a_h^T R a_h \right]$$

$$s_{h+1} = A s_h + B a_h$$

$$a_h = \theta^T s_h, s_0 \sim \mathcal{D}$$

multiagent

$$J^i(\theta) = \mathbb{E}_{s_0} \left[ \sum_{h=0}^{\infty} s_h^T Q^i s_h + (a^i)^T R^i a^i \right]$$

$$s_{h+1} = A s_h + \sum_{i=1}^N B a_h^i$$

$$a_h^i = (\theta^i)^T s_h, x_0 \sim \mathcal{D}$$

---

## Global Convergence of Policy Gradient Methods for the Linear Quadratic Regulator

---

Maryam Fazel<sup>\*1</sup> Rong Ge<sup>\*2</sup> Sham M. Kakade<sup>\*1</sup> Mehran Mesbahi<sup>\*1</sup>

### Abstract

Direct policy gradient methods for reinforcement learning and continuous control problems are a popular approach for a variety of reasons: 1) they

2016) and Atari game playing (Mnih et al., 2015). Deep reinforcement learning (DeepRL) is becoming increasingly popular for tackling such challenging sequential decision making problems.

multiagent

---

## Policy-Gradient Algorithms Have No Guarantees of Convergence in Linear Quadratic Games

Eric Mazumdar  
University of California, Berkeley  
Berkeley, CA  
mazumdar@berkeley.edu

Michael I. Jordan  
University of California, Berkeley  
Berkeley, CA  
jordan@cs.berkeley.edu

Lillian J. Ratliff  
University of Washington  
Seattle, WA  
ratliff@uw.edu

S. Shankar Sastry  
University of California, Berkeley  
Berkeley, CA  
sastry@coe.berkeley.edu

### ABSTRACT

We show by counterexample that policy-gradient algorithms have no guarantees of even local convergence to Nash equilibria in continuous action and state space multi-agent settings. To do so, we analyze gradient-play in  $N$ -player general-sum linear quadratic games, a classic game setting which is recently emerging as a bench-

of multi-agent reinforcement learning have made use of policy optimization algorithms such as multi-agent actor-critic [13, 17, 30], multi-agent proximal policy optimization [2], and even simple multi-agent policy-gradients [15] in problems where the various agents have high-dimensional continuous state and action spaces like StarCraft II [32].

# Multiagent policy gradient convergence condition

Results on subclasses of Markov games or depend on equilibria

- ▶ Zero-sum [Daskalakis et al. 2020], [Wei et al. 2021], [Cen et al. 2021], [K. Zhang et al. 2023], ...
- ▶ Potential [Leonardos et al. 2022], [R. Zhang et al. 2021], [Ding et al. 2022]
- ▶ Variationally stable equilibrium [Giannou et al. 2022]  $\implies$  local convergence

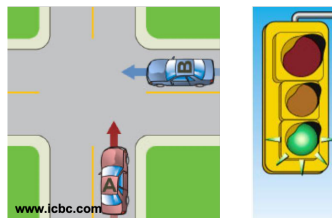
Our focus: presented as posters here

- ▶ Linear quadratic setting: conditions to be a potential game, characterizing number of equilibria
- ▶ Zero-sum Markov games: relaxing past assumptions while strengthening convergence result

## Relaxing the equilibrium notion

A probability distribution  $\mathcal{P}^*$  on  $\mathbf{A}$  is an equilibrium

$$\forall i \quad \mathbb{E}_{\boldsymbol{\theta} \sim \mathcal{P}^*} [J^i(\boldsymbol{\theta})] \leq \mathbb{E}_{\boldsymbol{\theta} \sim \mathcal{P}^*} [J^i(\tilde{\theta}^i, \boldsymbol{\theta}^{-i})], \quad \forall \tilde{\theta}^i$$



- ▶ Focus: learning algorithms that scale with number of agents



# Outline

Introduction

Learning Nash equilibria

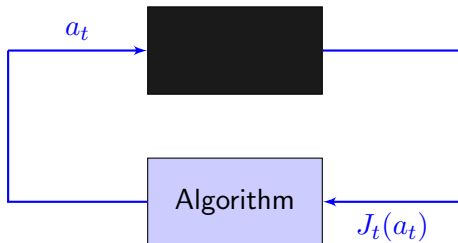
No-regret learning

- Normal form games

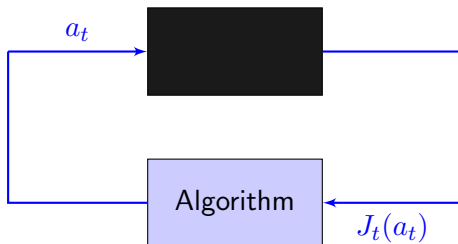
- Markov games

Conclusions

## Game as an adversarial bandit problem



## Game as an adversarial bandit problem

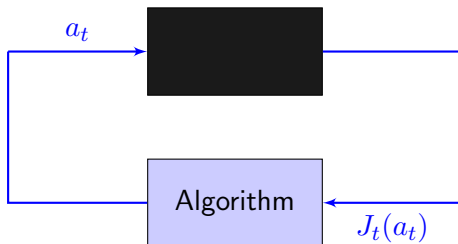


In a game:  $J_t(\cdot) := J^i(\cdot, a_t^{-i})$  for player  $i$

Benchmark: no-regret

► Regret: 
$$R(T) = \underbrace{\sum_{t=0}^T J_t(a_t)}_{\text{incurred cost}} - \min_a \underbrace{\sum_{t=0}^T J_t(a)}_{\text{best cost}}$$

## Game as an adversarial bandit problem



In a game:  $J_t(\cdot) := J^i(\cdot, a_t^{-i})$  for player  $i$

Benchmark: no-regret

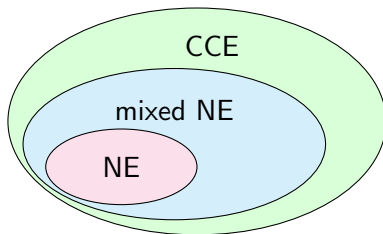
► Regret: 
$$R(T) = \underbrace{\sum_{t=0}^T J_t(a_t)}_{\text{incurred cost}} - \min_a \underbrace{\sum_{t=0}^T J_t(a)}_{\text{best cost}}$$

Algorithm is no-regret:  $R(T)/T \rightarrow 0$

## No-regret learning and equilibria

Let each player adopt a no-regret algorithm

- ▶ empirical distribution of actions  $\rightarrow$  *coarse-correlated equilibrium*



Remark

- ▶ CCEs may have better efficiency but
- ▶ CCEs can have weight on strictly dominated actions

# Multiplicative weight algorithms for no-regret

Player  $i$ 's actions  $\{1, 2, \dots, n\}$ , unknown cost:  $J_t(\cdot)$

Probability distribution on actions:  $w_t$

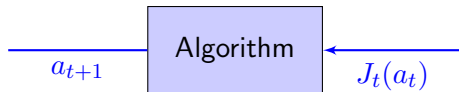
- ▶ sample:  $a_t \sim w_t$
- ▶ play the action:  $J_t(a_t)$
- ▶ update probabilities  $w_{t+1}$ , based on  $J_t(a_t)$ 
  - ▶ bandit feedback:  $w_{t+1}(k) = w_t(k) \exp^{-\eta_t J_t(k)}$ , for  $k = a_t$
  - ▶ full feedback:  $w_{t+1}(k) = w_t(k) \exp^{-\eta_t J_t(k)}$ , for  $\forall k$



# Optimal regret rates based on player's feedback

$n$ : number of actions for player,  $T$ : number of iterations

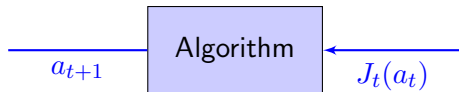
- ▶ Bandit feedback [Auer et al. 2003]



# Optimal regret rates based on player's feedback

$n$ : number of actions for player,  $T$ : number of iterations

- ▶ Bandit feedback [Auer et al. 2003]



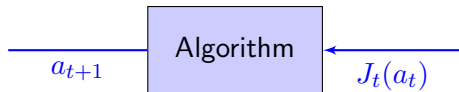
- ▶ Regret  $R(T)$  grows as  $\sqrt{Tn \log n}$



# Optimal regret rates based on player's feedback

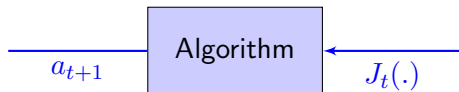
$n$ : number of actions for player,  $T$ : number of iterations

- ▶ Bandit feedback [Auer et al. 2003]



- ▶ Regret  $R(T)$  grows as  $\sqrt{Tn \log n}$

- ▶ Full feedback [Freund et al.1997]

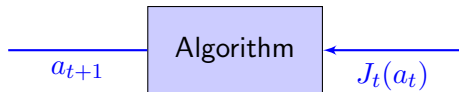


- ▶ Regret  $R(T)$  grows as  $\sqrt{T \log n}$

# Optimal regret rates based on player's feedback

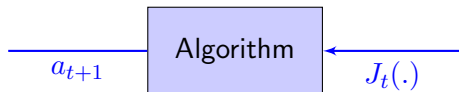
$n$ : number of actions for player,  $T$ : number of iterations

- ▶ Bandit feedback [Auer et al. 2003]



- ▶ Regret  $R(T)$  grows as  $\sqrt{Tn \log n}$

- ▶ Full feedback [Freund et al. 1997]



- ▶ Regret  $R(T)$  grows as  $\sqrt{T \log n}$

**Can we improve the dependence on  $n$ ?**

# Idea: mimic full feedback

Notice:  $J^i(\cdot, a_t^{-i})$  is a static function

Algorithms achieving optimal regret rate:

- ▶ bandit:  $w_{t+1}(k) = w_t(k) \exp^{-\eta_t J_t(k)}$
- ▶ full:  $w_{t+1}(\cdot) = w_t(\cdot) \exp^{-\eta_t J_t(\cdot)}$

Player  $i$  estimates its cost from past data  $\hat{J}_t^i(a_t^i, a_t^{-i}), a_t^i, a_t^{-i}$

- ▶ mimic full:  $w_{t+1}(\cdot) = w_t(\cdot) \exp^{-\eta_t \hat{J}_t^i(\cdot, a_t^{-i})}$

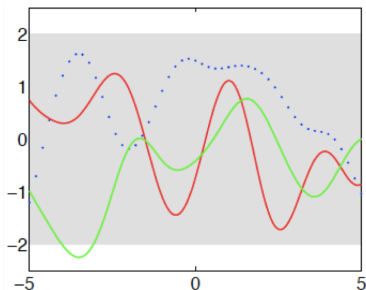


## Modeling class for cost function

$J$  has a bounded norm in a reproducing Kernel space  $\implies$   
 $J$  can be modeled by a Gaussian process

- ▶  $J(\mathbf{a}) \sim \mathcal{GP}(\mu(\mathbf{a}), k(\mathbf{a}, \mathbf{a}'))$
- ▶  $\mu$ : mean,  $k$ : covariance (kernel)
- ▶ examples of covariance function:

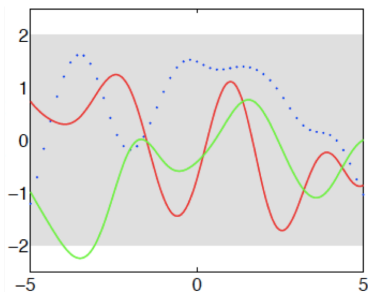
$$k_{poly}(\mathbf{a}, \mathbf{a}') = (l + \mathbf{a}^\top \mathbf{a}')^d, \quad k_{SE}(\mathbf{a}, \mathbf{a}') = \exp\left(-\frac{\|\mathbf{a} - \mathbf{a}'\|^2}{l^2}\right)$$



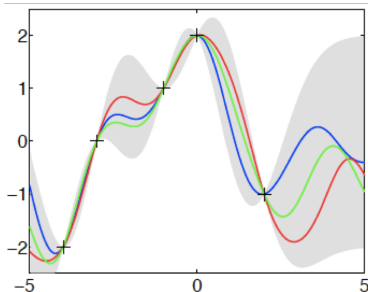
# Estimating the cost function distribution

$$J(\mathbf{a}) \sim \mathcal{GP}(\mu(\mathbf{a}), k(\mathbf{a}, \mathbf{a}'))$$

- ▶ observe: costs  $J(\mathbf{a}_l)$ , actions  $\mathbf{a}_l$ ,  $l = 1, \dots, t$
- ▶ obtain posterior distribution of  $J(\cdot)$ 
  - ▶ analytic formula for updating mean  $\mu_t(\cdot)$  and variance  $\sigma_t(\cdot)$



prior

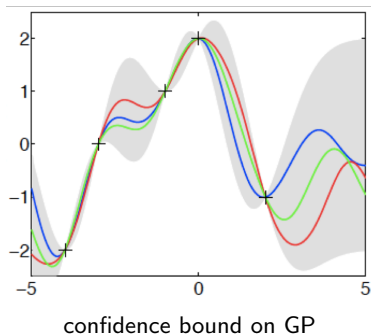


posterior

## Confidence bounds on the estimated cost

$$\hat{J}_t(\mathbf{a}) := \mu_t(\mathbf{a}) - \beta_t \sigma_t(\mathbf{a})$$

- ▶  $\hat{J}_t(\mathbf{a})$  small  $\implies$  cost low or uncertainty high
- ▶  $\beta_t > 0$  chosen to ensure  $\hat{J}_t(\mathbf{a}) \leq J(\mathbf{a})$  with high probability



# Gaussian process multiplicative weight algorithm (GPMW)

Player  $i$ 's actions  $\{1, 2, \dots, n\}$ , unknown cost:  $J(a^i, a^{-i})$

Optimistic cost estimate at time  $t$ :  $\hat{J}_t^i(\mathbf{a}) := \mu_t(\mathbf{a}) - \beta_t \sigma_t(\mathbf{a})$

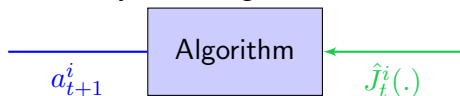
Probability distribution on actions:  $w_t$

- ▶ sample:  $a \sim w_t$
- ▶ observe:  $J^i(a, a_t^{-i})$  and  $a_t^{-i}$
- ▶ update  $\hat{J}_t^i(\cdot)$ 
  - ▶  $w_{t+1}(k) = w_t(k) \exp^{-\eta_t \hat{J}_t^i(\mathbf{a})}, \forall k$



## GPMW regret rates

- ▶ Mimic full feedback by observing others' actions

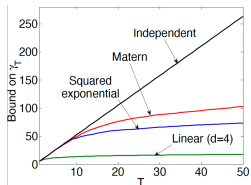


### Theorem

Assume: player's cost from a GP prior

- ▶ Regret grows as:  $(\sqrt{T \log n} + \gamma_T \sqrt{T})$

[Sessa, Bogunovic, MK, Krause, NeurIPS 2019]

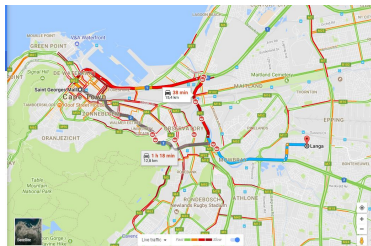


bound on  $\gamma_T$  based on the kernel [Srinivas et al. 2010]

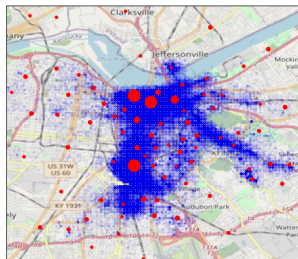


# Extensions of GP multi-agent learning

- ▶ contextual games [NeurIPS 2020],[AISTATS 2024], equilibria efficiency and game design [AISTAT2019, ICML2021]
- ▶ transportation network, resource allocation, electricity auctions, autonomous driving, energy management



Reducing congestion on road networks [ NeurIPS 2020]



Balancing bike distribution to maximize utility [ ICML 2021]

# Extension to multi-agent reinforcement learning (MARL)

- ▶ Dynamics  $s_{h+1} = f(s_h, a_h^1, a_h^2 \dots, a_h^N) + \omega_h$ 
  - ▶  $s_h \in S \subset \mathbb{R}^p$ ,  $a_h^i \in A^i \subset \mathbb{R}^q$
- ▶ Objective  $V^i(\pi^i, \pi^{-i}) = \mathbb{E}[\sum_{h=0}^{H-1} r^i(s_h, \pi^i(s_h), \pi^{-i}(s_h))]$

Approach: estimate the transition function  $f$  via its posterior mean  $\mu_t(s, \mathbf{a}) \in \mathbb{R}^p$  and confidence functions  $\Sigma_t(s, \mathbf{a}) \in \mathbb{R}^{p \times p}$

# Approach: model-based learning of equilibrium distribution

- ▶ Initialize  $\mathcal{P}_0$ . For  $t = 0, 1, \dots$

- ▶ sample  $(\pi_t^1, \dots, \pi_t^N) \sim \mathcal{P}_t$



- ▶ estimate  $P(\cdot | s_h, a_h^1, \dots, a_h^N) \rightarrow \{\bar{V}_t^i(\boldsymbol{\theta})\}_{i=1}^N$
    - ▶ compute  $\mathcal{P}_{t+1}$  as the equilibrium distribution of  $\{\bar{V}_t^i(\boldsymbol{\theta})\}_{i=1}^N$
- ▶  $\bar{V}_t^i(\cdot)$ : optimistic estimate of  $V^i(\cdot)$  at iteration  $t$

# Regret of the MARL algorithm

Dynamic regret

$$R^i(T) := \sum_{t=1}^T \max_{\pi \in \Pi^i} \mathbb{E}_{\pi_t^{-i}} [V^i(\pi, \pi_t^{-i})] - \mathbb{E}_{\pi_t} [V^i(\pi_t)]$$

## Theorem

Under Lipschitz continuity of  $f$ ,  $\{r^i, \pi^i\}_{i=1}^N$

$$R^i(T) = \mathcal{O}(LH^{1/2} \sqrt{T\mathcal{I}_T}) + \sum_{t=0}^T \epsilon_t$$

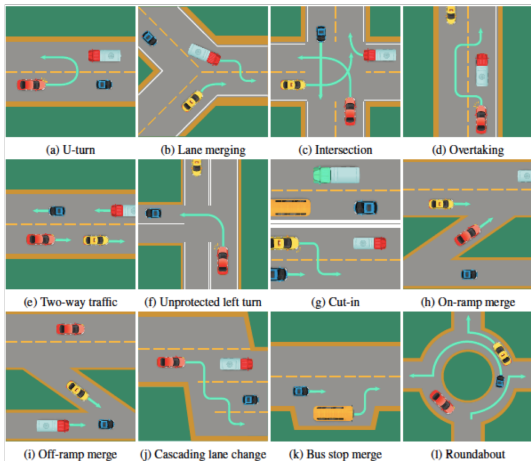
- ▶  $\mathcal{I}_T(p, H, \gamma_{HT})$ : information gain
- ▶  $\epsilon_t$ : approximate CCE for  $\{\bar{V}_t^i(\boldsymbol{\theta})\}_{i=1}^N$

[Sessa, MK, Krause, ICML 2022]

# Example: Multi-agent RL in autonomous driving

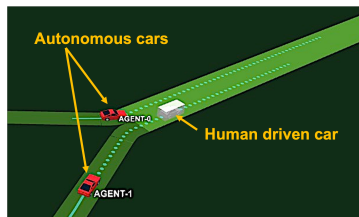
SMARTS autonomous car simulation environment [Zhou et al. 2021]

- ▶ testing multi-agent RL algorithms for autonomous driving
- ▶ realistic traffic data and car dynamics



# Multiagent reinforcement learning for autonomous driving

- ▶ Objective: progress towards the goal, avoid collision
- ▶ Dynamics:  $P(\cdot | s_h, a_h^1, a_h^2)$ 
  - ▶  $s$ : positions and velocities of cars
  - ▶  $a^i$ : heading and speed,  $i = 1, 2$
  - ▶  $\pi_{\theta^i}(s)$ : parametrized by neural networks,  $i = 1, 2$

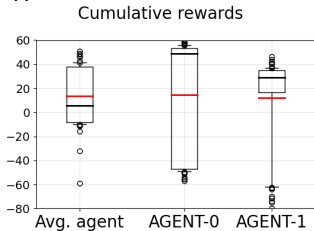


The autonomous cars can coordinate and overtake the human-driven car

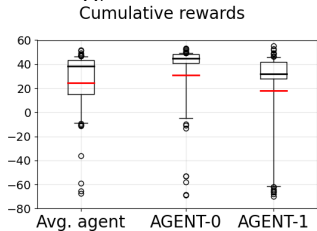
Implementation on multiagent autonomous car simulation environment [Zhou et al. 2021]

# Learning to coordinate

Learning to coordinate  $\implies$  less breaking, more successful merges



Single-agent optima



Multi-agent equilibrium

Average rewards for the agents

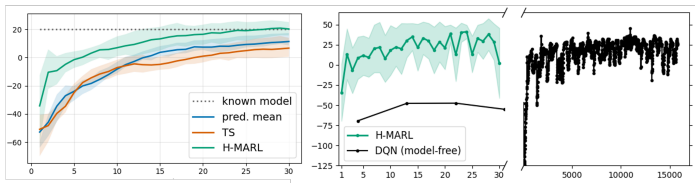


Figure: left: value of optimism, right: value of learning the model

# Outline

Introduction

Learning Nash equilibria

No-regret learning

Conclusions

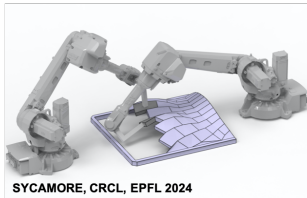
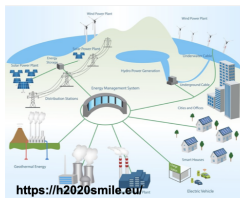


# Summary

- ▶ Payoff-based learning of Nash equilibria
  - ▶ require assumptions on pseudo-gradient or the equilibrium
  - ▶ challenging to extend to Markov games
  
- ▶ No-regret learning
  - ▶ tractable and ensure convergence to CCEs
  - ▶ can improve rates using a model-based approach

# Outlook

- ▶ Learning equilibria in Markov games under coupling constraint
- ▶ Provable algorithms under partial and asymmetric information
- ▶ Learning of “good” equilibria, mechanism design
- ▶ Applications: power markets, robotics, autonomous driving



# Acknowledgements

- ▶ Former and current students and postdocs: O Karaca, L Furieri, P Giuseppe Sessa, A Maddux, G Salizzoni, S Hosseinirad, R Ouhamma
- ▶ Collaborators: T Tatarenko, A Krause, Bugonovic
- ▶ Funding : ERC, NSERC Canada, Swiss National Fund, NCCR Automation



<https://www.epfl.ch/labs/sycamore/>