

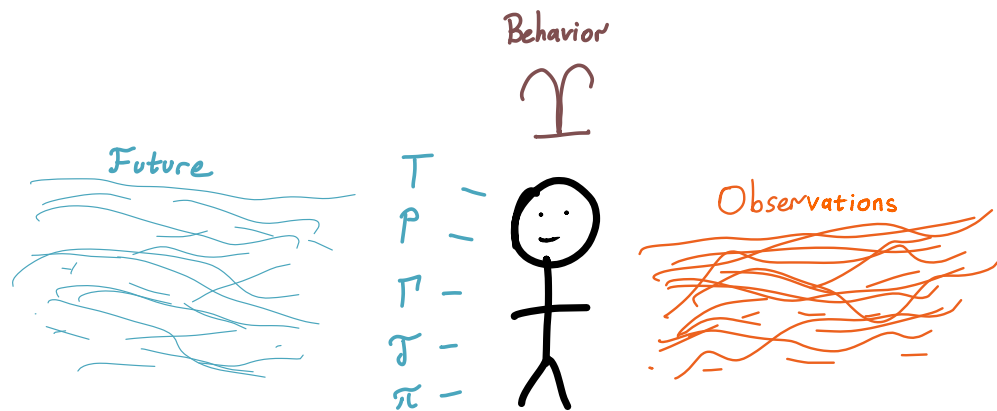
The Learning Approach to Games

Melih İşeri

Joint work with Erhan Bayraktar

- University of Michigan ●





* Players know each other's strategies and don't deviate?

* Central Planner announce and design compatible/stable policies for individuals

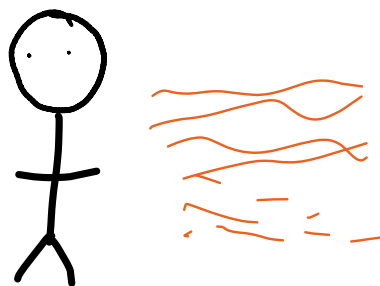
- Environmental Regulations
- Traffic Management
- Public Health Initiatives

Central Planner

α^1 ↗

α^2 →

α^3 ↘



* Central Planner cannot model players up to the detail of their Learning Algorithms.

Definition [Player]

We say $(O, L_1, \dots, L_k, \Upsilon)$ is a player in the environment $(\Omega^u, \mathcal{F}^u, \mathbb{P}^u)$

O : Observation : $\Omega^u \times \mathbb{N} \rightarrow \mathcal{E}$

L_i : Learning Algorithm : $\mathcal{E} \rightarrow \mathcal{M}_i$ [spaces of estimations w/ domain \mathcal{D}]

Υ : Behavior : $\mathcal{M}_1 \times \dots \times \mathcal{M}_k \rightarrow (\mathcal{D} \rightarrow \mathcal{P}(A))$

where \mathcal{E} is the set of finite sequences of observables

A is the set of finite sequences of actions

Furthermore, we call ${}^n \Upsilon : \Omega^u \times \mathbb{N} \rightarrow (\mathcal{D} \rightarrow \mathcal{P}(A))$ the planned behavior of the player
at age n .
 ${}^n \Upsilon \equiv \Upsilon(L_1({}^n O), \dots, L_k({}^n O))$

Definition [Recurrent Behavior]

We say $\mathcal{T}^*: D \rightarrow \mathcal{P}(A)$ is a (r, δ) -recurrent behavior if

$$\mathbb{P}^u \left(\liminf_{n \rightarrow \infty} d(\mathcal{T}^*, {}^n \mathcal{T}) > r \right) \leq \delta$$

Lemma Suppose $\exists (\varphi_1^*, \dots, \varphi_k^*) \in (\mathcal{M}_1, \dots, \mathcal{M}_k)$ such that

$$\mathbb{P}^u \left(\liminf_{n \rightarrow \infty} \max_{1 \leq \ell \leq k} d_\ell(\varphi_\ell^*, \mathcal{L}_\ell({}^n o)) = 0 \right) = 1$$

If $\mathcal{T}: \mathcal{M}_1 \times \dots \times \mathcal{M}_k \rightarrow (D \rightarrow \mathcal{P}(A))$ is a continuous mapping, then

$$\mathcal{T}^* := \mathcal{T}(\varphi_1^*, \dots, \varphi_k^*)$$

is almost surely a recurrent behavior.

Chess As a player, we need to learn a lot!

- How many steps ahead I can analyze?
- What are the values of those states in the future?
- Is my opponent playing aggressively or defensively?
- How can I trick/deceive my opponent?

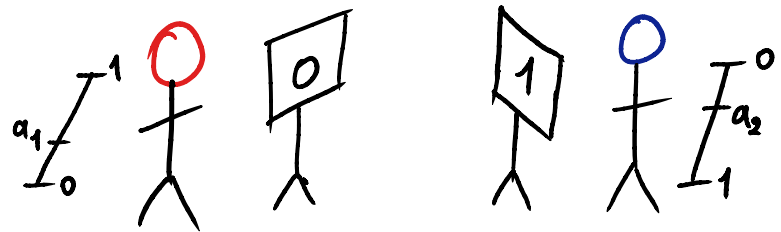
Each piece might
have own positional values
Multi-Dimensional values

At the very late stages, trained player does the same move almost surely!

Openings of players might have a fixed distribution too.

Learning might continually evolve for many other configurations.

Two-player [Simple yet dynamic]



* Lose \$1 if Player 2 show 1. Gain $c \cdot a_1$.

* Lose \$1 if Player 1 \neq Player 2

Player 2 : Observe opponent and if noise is acceptable, do the same.
Otherwise, explore other actions with rationale to penalize noise.

Player 1 : If opponent appears at 0, explore larger actions to reduce the cost
If opponent appears at 1, keep exploring until cost is consistent with expectation

* They don't announce their strategies.. How are they gonna behave?

* The key is to design the players! Who is playing?

Designing the Player

- Q-Learning. Simple, converges, not dynamic.

- Keep tables. " " "

- Predict Opponent (Γ) & Randomize Cost (F) •

$${}^n \Gamma^i = \mathcal{L}_{\Gamma}^i ({}^n \Theta^i) = \frac{1}{K} \sum_{k=1}^K \int \mathcal{N}_{\Gamma}^{i,k}$$

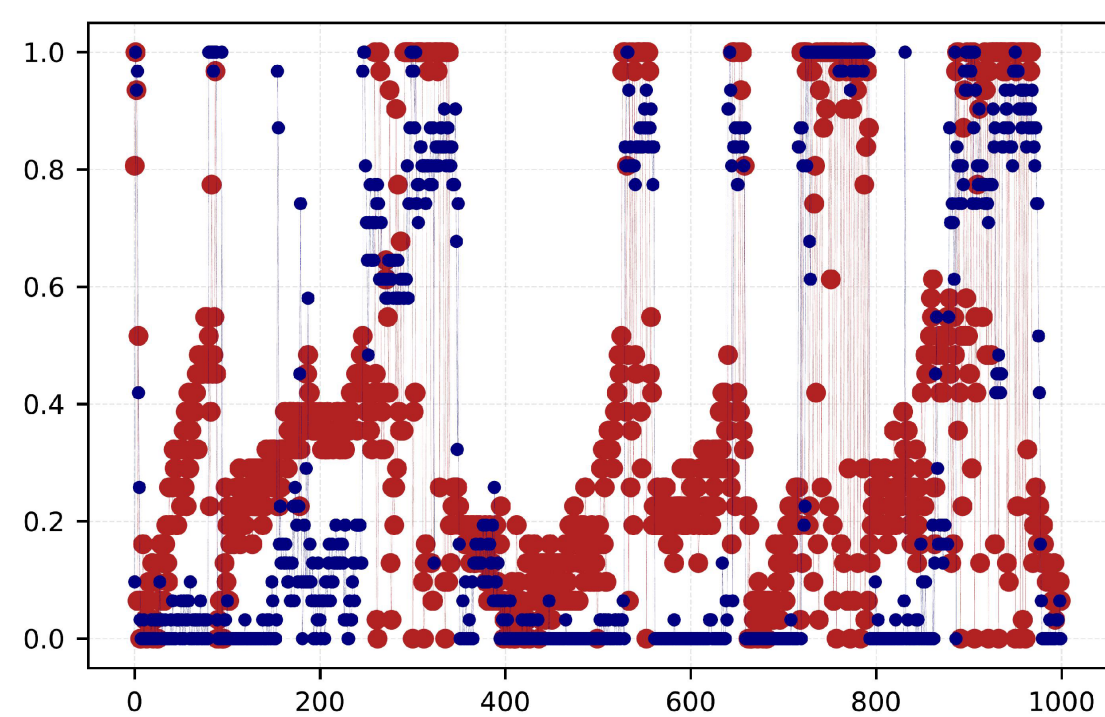
$$F^i(\ell, w, a) = \mathcal{N}_F^{i,\ell}(w)(a)$$

$(\ell, w) \in \{1, \dots, K\} \times \Omega'$, w/ unif. first marginal

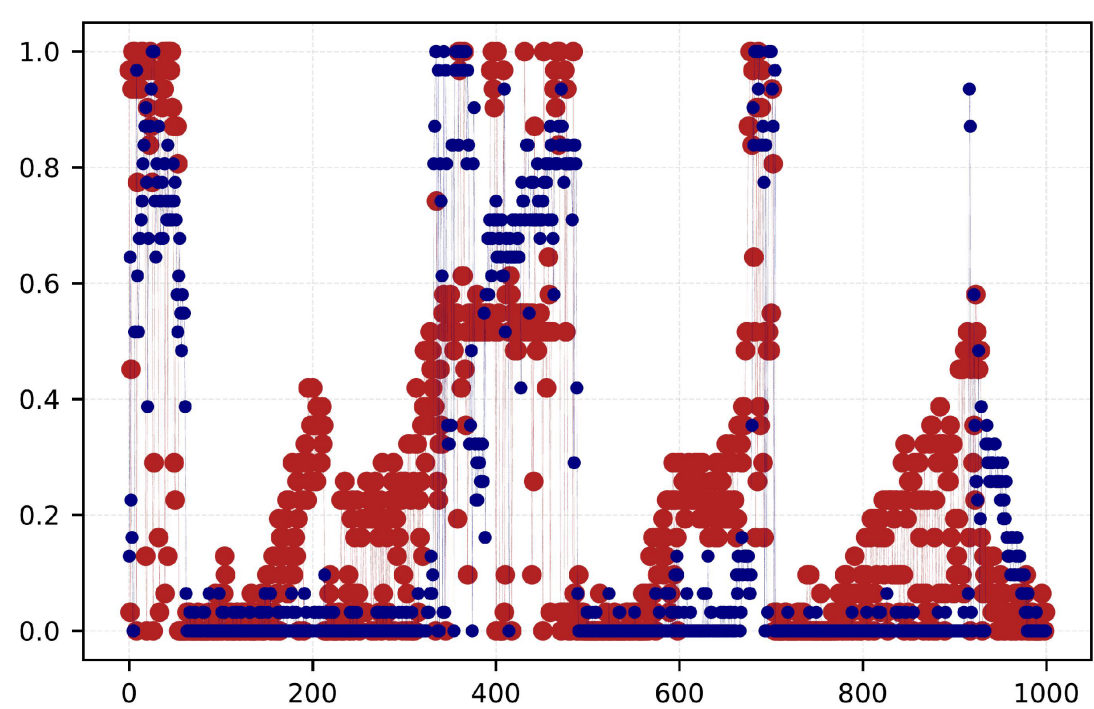
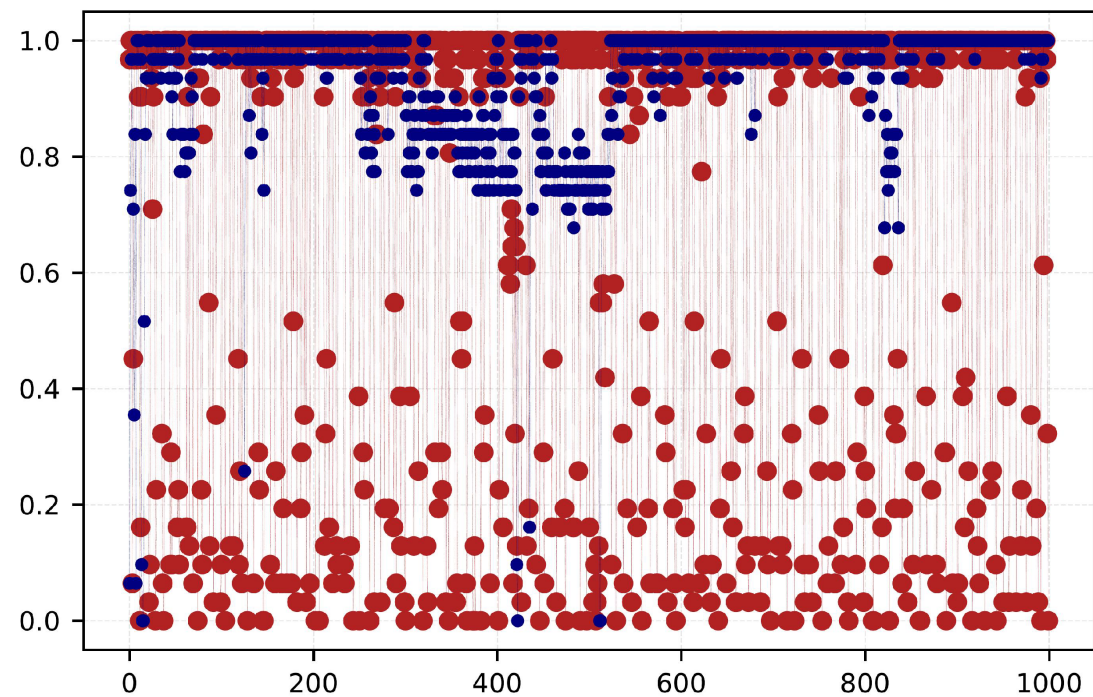
$$\mathcal{N}_{\Gamma}^{i,k} : [0,1] \rightarrow [0,1]$$

$$\mathcal{N}_F^{i,\ell} : \Omega' \times [0,1] \rightarrow \mathbb{R}$$

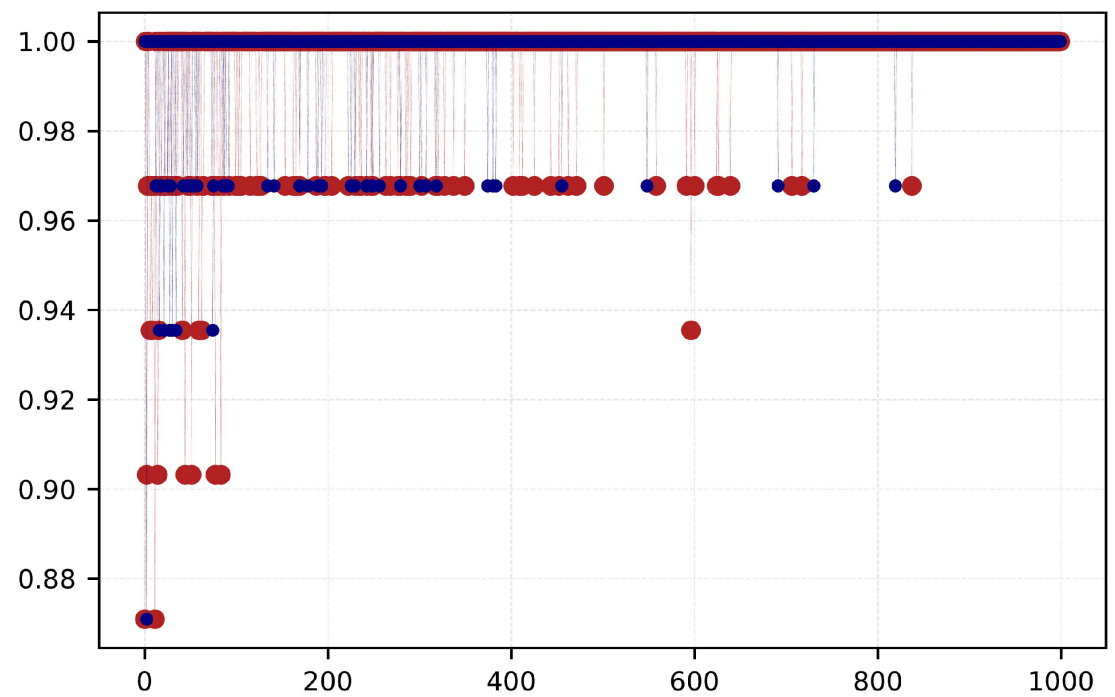
- Players will predict only one step ahead, yet they will be dynamic!
- It is crucial to define the player to predict the behavior
How do they process memories? What are their expectations!?



$$\hat{c} = \frac{3}{10}, B^1 = \frac{1}{10} \quad / \quad c=1, B^1=1 \quad \Downarrow$$



$$\hat{c} = \frac{1}{20}, B^1 = -\frac{1}{10} \quad / \quad c=1, B^1=0 \quad \Downarrow$$



Discrete Games

T : time, $\Omega \doteq \prod_{t \in T} S_t$: states, A : actions, \bar{A} : controls

\hat{T}	:	$T \times \Omega \rightarrow T$:	horizon
\hat{p}	:	$T \times \Omega \times \bar{A} \times S \rightarrow \mathbb{R}^+$:	transition
$\hat{F}/\hat{\phi}$:	$\hat{\Omega} \times T \times \Omega \times A \rightarrow \mathbb{R}$:	value
$\hat{\pi}$:	$\hat{\Omega} \times T \times \Omega \rightarrow \mathcal{P}(A^i)$:	optimal control
$\hat{\Gamma}$:	$A^i \rightarrow \mathcal{P}(\bar{A})$:	opponent's strategy

$$\mathcal{J}(t, x, \alpha) \doteq \int_{\bar{A}} \mathcal{J}(t, x, \vec{\alpha}) \hat{\Gamma}_{\alpha}(d\vec{\alpha}), \text{ where}$$

$$\mathcal{J}(t, x, \vec{\alpha}) \doteq \mathbb{E}^{t, x, \vec{\alpha}} \left[\hat{\phi}(t + \hat{T}, X_{t + \hat{T}}) + \sum_{s=t}^{t + \hat{T} - 1} \hat{F}(s, X_s, \alpha_s) \right]$$

Behavior

$$\mathcal{U}^{t, x}(d\alpha) \doteq \int_{\hat{\Omega}} \hat{\pi}(\hat{w}, t, x)(d\alpha) d\hat{P}(\hat{w})$$

Definition [Uncertain Equilibrium] We say $\{\vec{T}, \vec{p}, \vec{F}, \vec{\phi}, \vec{\pi}, \vec{\Gamma}\}$ is (ϵ, ν, δ) -uncertain equilibrium at $(t, x) \in \mathbb{T} \times S_t$ under the Learning Algorithms $\vec{\mathcal{L}}$ if,

(i) $\{\vec{T}, \vec{p}, \vec{F}, \vec{\phi}, \vec{\pi}, \vec{\Gamma}\}$ are the priors of the players

$$(ii) \int_{\hat{\Omega}} \int_{\mathcal{A}^i} \left(\sup_{\tilde{\alpha} \in \mathcal{A}^i} {}^n \mathcal{T}^i(\hat{w}, t, x, \tilde{\alpha}) - {}^n \mathcal{T}^i(\hat{w}, t, x, \alpha) \right) {}^n \pi^i(\hat{w}, t, x)(d\alpha) d\hat{\mathbb{P}}(\hat{w}) \leq \epsilon \quad \forall i, n$$

$$(iii) \mathbb{P}^\nu \left(\liminf_{n \rightarrow \infty} \sup_{i \in \mathbb{N}_0} d^{t, x, i}({}^0 \Gamma_{\vec{\mathcal{L}}}^{t, x, i}, {}^n \Gamma_{\vec{\mathcal{L}}}^{t, x, i}) > \nu \right) \leq \delta$$

Correlated Equilibrium

$$\rho \in \mathcal{P}(\bar{\mathcal{A}}) ; \quad \rho(d\vec{\alpha}) = \rho^{-i}(d\vec{\alpha} | \alpha^i) \rho^i(d\alpha^i)$$

Nash - type $\int_{\mathcal{A}^i} \int_{\bar{\mathcal{A}}} \sup_{\tilde{\alpha}^i \in \mathcal{A}^i} \mathcal{J}^i(\tilde{\alpha}^i, \vec{\alpha}^{-i}) \rho^{-i}(d\vec{\alpha} | \alpha^i) \rho^i(d\alpha^i)$

Correlated $\int_{\mathcal{A}^i} \sup_{\tilde{\alpha}^i \in \mathcal{A}^i} \int_{\bar{\mathcal{A}}} \mathcal{J}^i(\tilde{\alpha}^i, \vec{\alpha}^{-i}) \rho^{-i}(d\vec{\alpha} | \alpha^i) \rho^i(d\alpha^i)$

Uncertain $\int_{\hat{\Omega}} \sup_{\tilde{\alpha}^i \in \mathcal{A}^i} \int_{\bar{\mathcal{A}}} \mathcal{J}^i(\hat{\omega}, \tilde{\alpha}^i, \vec{\alpha}^{-i}) \Pi_{\tilde{\alpha}^i}^i(d\vec{\alpha}) \hat{\mathbb{P}}(d\hat{\omega})$

C. Correlated $\sup_{\tilde{\alpha}^i \in \mathcal{A}^i} \int_{\mathcal{A}^i} \int_{\bar{\mathcal{A}}} \mathcal{J}^i(\tilde{\alpha}^i, \vec{\alpha}^{-i}) \rho^{-i}(d\vec{\alpha} | \alpha^i) \rho^i(d\alpha^i)$

More Estimations [Don't let anyone stop you.]

- Communication

- Embedding of Raw Observations

- Best Expected $\hat{B}: \mathbb{T} \times \Omega \rightarrow \mathbb{R}$

• (ii') ${}^n \mathcal{K}^i(t, x) > K \quad \forall i, n$ •

$${}^n \mathcal{K}^i(t, x) \doteq \hat{\mathbb{P}} \left(\int_{\mathcal{X}^i} {}^n \mathcal{J}^i(\hat{w}, t, x, \alpha) {}^n \hat{\pi}^i(\hat{w}, t, x)(d\alpha) > {}^n \hat{B}^i(t, x) \right)$$

• Desperate - Discouraged - Doubtful - Cautious - Hopeful - Determined - Confident - Optimistic - Euphoric •

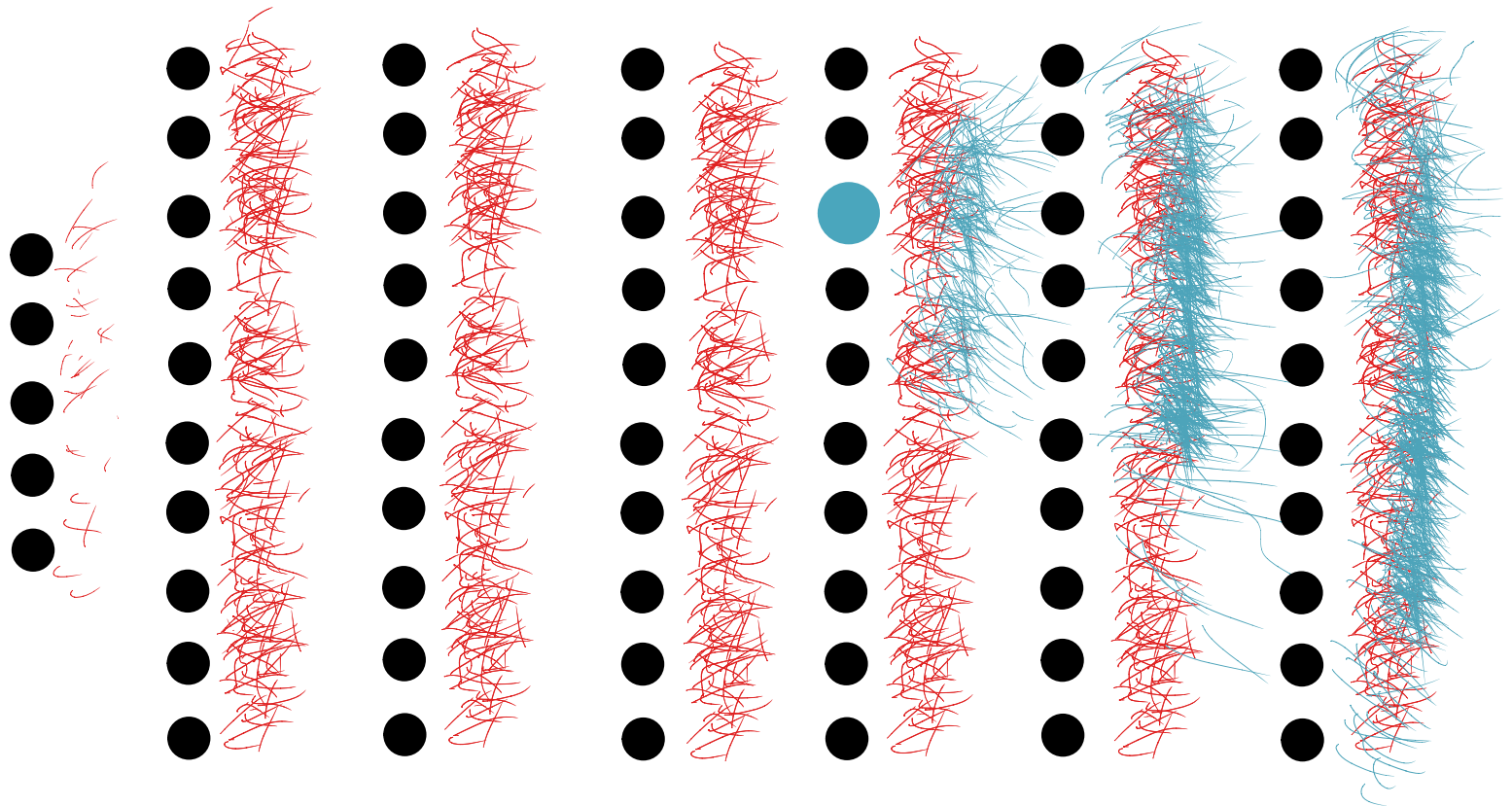
$$(ii'') \quad \text{supp} \left({}^n \hat{\pi}_{t,x,\alpha}^i \right) \subset \text{supp} \left({}^n \Upsilon^{t,x,1} \times {}^n \Upsilon^{t,x,2} \times \dots \right)$$

Time-consistency (DPP) We say $\{\hat{T}, \hat{p}, \hat{\Gamma}, \hat{F}, \hat{\phi}, \hat{\pi}\}$ yields time-consistent

Value : $\int_{\mathcal{U}^i} \mathcal{J}(T_0; \hat{w}, t, x, \alpha) \hat{\pi}(\hat{w}, t, x)(d\alpha) = \int_{\mathcal{U}^i} \mathcal{J}(\hat{w}, t, x, \alpha) \hat{\pi}(\hat{w}, t, x)(d\alpha), 0 \leq T_0 \leq \hat{T}(t, x)$

In particular, $\hat{\phi}(t, x) = \sup_{\alpha} \int_{\mathcal{U}} E^{t, x, \vec{\alpha}} \left[\hat{\phi}(t + \hat{T}, X_{t + \hat{T}}) + \sum_{s=t}^{t + \hat{T} - 1} \hat{F}(s, X_s, \vec{\alpha}) \right] \hat{\Gamma}_{\alpha}(d\vec{\alpha})$

Strategy :



Stated Mean Field Games Observations can be generated by symmetries!

• Stated: $T, p(t, x, \mu, a, \cdot), \phi(x, \mu), F(t, x, \mu, a)$ are given (constant \mathcal{L}) and $\hat{\pi}$ yields optimal control.

• Player estimates $\hat{\Gamma} \in \mathcal{P}(\mathcal{P}(S \times \mathcal{A}))$ • $\mathcal{P}(S \times \mathcal{A}) \leftrightarrow \vec{\mathcal{A}}$ •

• $\mathbb{E}_{s+1}(dy, d\alpha) = \int_S p(s, x, \mu_s^{\mathbb{E}}, \alpha(s, x, \mu_s^{\mathbb{E}}); dy) \mathbb{E}_s(dx, d\alpha)$ where $\mu_s^{\mathbb{E}} \doteq \mathbb{E}_s(\cdot, \mathcal{A})$

– $\mathcal{J}(t, \mu, x, \alpha) \doteq \int_{\mathcal{P}(S \times \mathcal{A})} \mathcal{J}(t, \mathbb{E}, x, \alpha) d\hat{\Gamma}(\mathbb{E}); \mathcal{J}(t, \mathbb{E}, x, \alpha) \doteq \mathbb{E}^{t, \mathbb{E}, x, \alpha} [\phi(\cdot) + \mathcal{L} \cdot]$ –

• Observables $\mathcal{P}(\mathcal{P}(S \times \mathcal{A}))$ and a Learning Algorithm (homogeneous):

$$\mathcal{L}_{\Gamma}({}^n \mathcal{O}) \doteq c \delta(\mu, \delta_{n+1} \alpha) + (1-c) {}^n \mathcal{O}$$

where ${}^{n+1} \alpha$ is the optimal control under ${}^n \mathcal{O}$.

Example $S = [0, 1]$, $T = \{0, 1\}$, $A = [0, 1]$ and

$$p(0, x, a, \mu; dy) = \delta_a$$

Introduce a discontinuous cost as

$$\mathcal{J}(\mathbb{E}; \alpha) \doteq \mathbb{E}^{\mathbb{E}, \alpha} \left[X_1^\alpha \mathbb{1}_{\{\bar{\mu}_1^\alpha \in [0, 1/2]\}} - X_1^\alpha \mathbb{1}_{\{\bar{\mu}_1^\alpha \in (1/2, 1]\}} \right] \text{ where } \bar{\mu}^\alpha \doteq \int_{[0, 1]} x d\mu^\alpha$$

Whereas there exists no relax equilibrium, \mathcal{L}_T oscillates around

$$\frac{1}{2} (\delta_{\delta_0} + \delta_{\delta_1}) \in \mathcal{P}(\mathcal{P}(A))$$

and induces an action distribution δ_0 and δ_1 infinitely often.

Terminologies in Reinforcement Learning

\hat{T} : discount factor / stopping-time / options

\hat{p} : model based methods

\hat{F} : rewards

$\hat{\phi}$: value

$\hat{\pi}$: policy learning

$\hat{\Gamma}$: behavior prediction

Algorithmic Collusion (+ Neil Mascarenhas)

$Q(x, a)$ + arbitrary randomization
of behavior

→ 400.000 to millions of steps!

$$\int_{\hat{\Omega}} \left(\int_{\mathcal{A}} T(\hat{w}, x; \alpha) \hat{\pi}(\hat{w}, x) (d\alpha | \alpha(x) = a) \right) d\hat{P}(\hat{w}) \rightarrow 20-30 \text{ steps}$$

Thank

You