

Fictitious Play in Markov Games with Single Controller



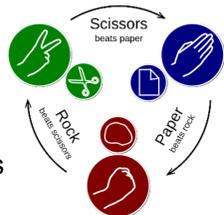
Muhammed O. Sayin, Dep. Electrical & Electronics Eng., Bilkent University, Ankara, Turkey
 Kaiqing Zhang, Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA
 Asuman Ozdaglar, Dep. Electrical Eng. & Computer Sci., Massachusetts Institute of Technology, Cambridge, MA

Frontiers in AI:
Multi-agent Learning in Dynamic Environments

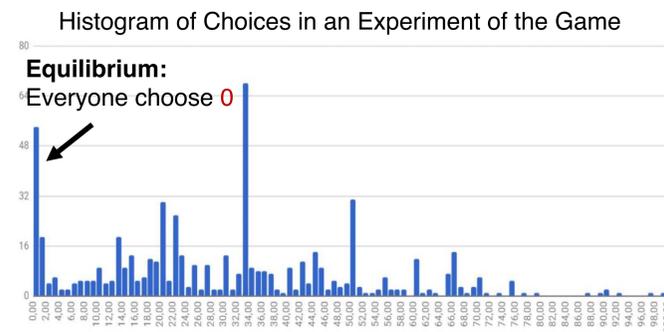


- Focus:**
- **Individual objective optimization** (beyond equilibrium seeking)
 - **Dynamic environments** (beyond repeated games)
 - Systematic guarantees

Game Theory:
Nash Equilibrium: Players do not have any **incentive** to change their strategies unilaterally.
Always exists in finite games if players can **randomize** their actions.



How **predictive** is a Nash equilibrium?
 Let's play a game!
 1. Everyone will pick an **integer between [0,100]** simultaneously.
 2. The one **closest to the 2/3 of the average** will share the price.



Humans are **NOT** necessarily seeking for an equilibrium!

Self-interested players **reach to an equilibrium** if they play the same game **repeatedly** and receive **feedback**

Fictitious Play:

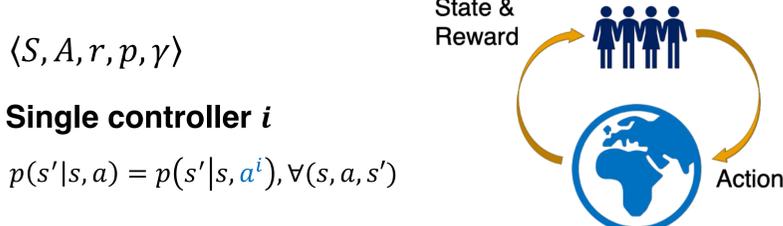
- **Model the opponent** as playing a **stationary strategy**
- **Form a belief** about the **opponent strategy**
- **Respond** to the belief by taking a **greedy best response**



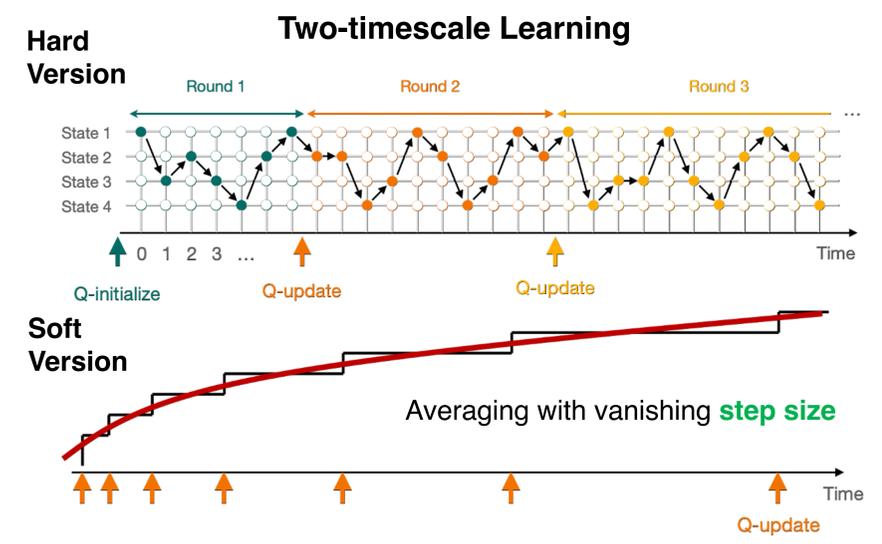
What about **BEYOND** the **games with repeated play**?

Markov Games (a.k.a. Stochastic Games):

- **Multi state** environment
- Markovian **state transitions**
- Markovian rewards
- **Objective:** **discounted** sum of rewards collected over **infinite horizon**



Challenges	Solutions
Trade-off between now and (ambiguous) future	(Standard) Define Q-function for state-action values
Non-stationarity of how the future is perceived	(Key Idea) Two-timescale learning
Deviation from the identical-interest structure in the stage games	(Key Idea) Single-controller Markov games for strategic equivalence to identical interest games!



Two-timescale Fictitious Play:

- **Model the opponent** as playing a **stationary strategy specific to each state**
- **Form a belief** about the **opponent strategy** (**Large** steps)
- **Form a belief** about the **Q-function** (**Small** steps)
- **Respond** to the beliefs by taking a **greedy best response**

For each state s

Player i

$$\pi_s^j \leftarrow \pi_s^j + \alpha_k (a_s^j - \pi_s^j), \forall j \neq i$$

$$a_s^i \in \operatorname{argmax}_{a^i} \{E_{a^{-i} \sim \pi_s^{-i}} [Q_s^i(a^i, a^{-i})]\}$$

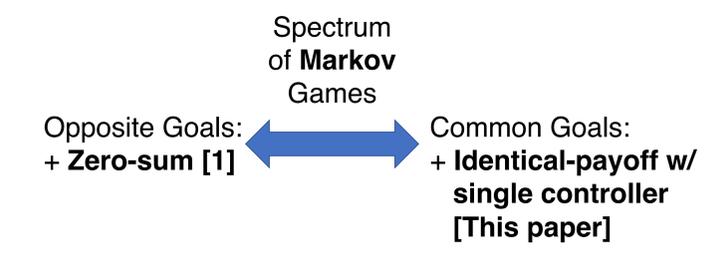
$$Q_s^i \leftarrow Q_s^i + \beta_k (R_s + \gamma E_{s'} [v_s^i] - Q_s^i)$$

$$v_s^i \leftarrow \max_{a^i} \{E_{a^{-i} \sim \pi_s^{-i}} [Q_s^i(a^i, a^{-i})]\}$$

Two-timescale step-sizes: $\lim_{k \rightarrow \infty} \frac{\beta_k}{\alpha_k} = 0$

Theorem:
 Given a **multi-player identical-interest discounted Markov game**, suppose that players follow the two-timescale fictitious play dynamics with non-summable vanishing step sizes. If each state gets visited infinitely often, then we have the beliefs π and Q **converge to a stationary equilibrium** and the associated Q-function **almost surely**.

Corollary (Potential-game-like extension):
 Under the conditions in Theorem, suppose that $r^j(s, \tilde{a}^j, a^{-j}) - r^j(s, a) = r^i(s, \tilde{a}^j, a^{-j}) - r^i(s, a)$ for all $(s, a), \tilde{a}^j$ and $j \neq i$ given that player i is the single controller. Then, we have the beliefs π and Q **converge to a stationary equilibrium** and the associated Q-function **almost surely**.



Sketch of the Proof:
 Show **accumulated monotonicity-like condition:**

$$\liminf_{k_1 \rightarrow \infty} \inf_{k_2 \geq k_1} \{Q_{k_2}^i(s, a) - Q_{k_1}^i(s, a)\} \geq 0,$$

implying converges since Q-functions are bounded.

Alternatively, define

$$\underline{u} = \min_{(s,a)} \{r + \gamma E_{s'} [E_{a^{-i} \sim \pi_k(s)} [Q_k^i(s', a')]] - Q_k^i\}$$

And show $\liminf_{k_1 \rightarrow \infty} \inf_{k_2 \geq k_1} \sum_{k=k_1}^{k_2} \beta_k \underline{u}_k \geq 0.$

We can show it because \underline{u}_k^i satisfies

$$\underline{u}_{k+1}^i \geq \underline{u}_k^i (1 - (1 - \gamma)\beta_k) - e_k$$

Absolutely summable error term (for **single controller case**)

Given that Q-function estimates converge, we can show they indeed converge to the Q-function associated with an equilibrium. ■

[1] M. O. Sayin, F. Parise, and A. Ozdaglar, "Fictitious play in zero-sum stochastic games," *SIAM J. Cont. Opt.*, in print.