

Lecture 6: Stochastic Approximation & Learning Algorithms

Last Lecture Recap

Last time: Repeated N -player finite games

- P_i uses mixed strategy \mathbf{x}_i^k at iteration k , depending on
 - Observation w_i^k
 - Internal variable z_i^k
 - Update rule f_i
- **Goal:** Learning process that iteratively updates z_i^k and \mathbf{x}_i^k

Game Elements in This Lecture

Players N players playing repeatedly at $k = 0, 1, 2, \dots$

Strategy Use action j -th or e_{ij} with $\mathbb{P}(e_i^k = e_{ij}) = x_{ij}^k$ where $\mathbf{x}_i^k \in \Delta_i$ at iteration k

Cost $\bar{J}_i(\mathbf{x}^k)$ evaluated at each iteration

Focus on:

- Intuition (origin of alg.)
- Informational reqs
- Steady-state is at NE (or related to NE)
- Idea of convergence (CT-ODE)

6.1 Stochastic approximation - ODE method (Appx A, Ch 9)

DT - Discrete-Time Stochastic Process. A generic learning process has the form:

$$\begin{cases} \mathbf{z}_i^{k+1} = \mathbf{z}_i^k + \gamma^k f(\mathbf{z}_i^k, w_i^k) \\ \mathbf{x}_i^k = \sigma(\mathbf{z}_i^k) \end{cases}$$

where:

- \mathbf{z}_i^k : internal state/variable at iteration k
- $\gamma^k > 0$: step size (learning rate) at iteration k
- $f_i(\cdot, \cdot)$: update function (depends on state and observation)
- w_i^k : information/observation at iteration k
- $\sigma(\cdot)$: strategy mapping from internal state to mixed strategy

How to study the long-run behavior ($k \rightarrow \infty$)?.

Problem: Stochastic processes are complex to analyze directly!

Solution: Use **Stochastic Approximation**

- Connect DT stochastic process to CT-ODE (continuous-time ordinary differential equation)
- Analyze the simpler ODE instead

DT Stochastic Process. General form:

$$\mathbf{z}^{k+1} = \mathbf{z}^k + \gamma^k [f(\mathbf{z}^k) + \xi^k] \longrightarrow \text{deterministic} + \text{noise}$$

where:.

1. $\{\mathbf{z}^k\}$: stochastic process
2. $\{\gamma^k\}$: diminishing step size
 - $\gamma^k \geq 0$ for all k
 - $\sum_{k=0}^{\infty} \gamma^k = \infty$ (infinite travel)
 - $\lim_{k \rightarrow \infty} \gamma^k = 0$ (vanishing step size)
3. $\{\xi^k\}$: perturbations with **martingale difference property**

$$\mathbb{E}[\xi^k \mid \mathcal{F}_k] = 0 \quad (\text{zero-mean, conditioned on all past info})$$

4. $f(\mathbf{z}^k)$: mean update direction

$$f(\mathbf{z}^k) = \mathbb{E} \left[\frac{1}{\gamma^k} (\mathbf{z}^{k+1} - \mathbf{z}^k) \mid \mathcal{F}_k \right]$$

Under some further assumptions, long run behaviour of the DT stochastic process can be described by the long run behaviour of the CT-ODE

$$\frac{\mathbf{z}^{k+1} - \mathbf{z}^k}{\gamma_k} \approx \dot{\mathbf{z}}(t_k) = f(\mathbf{z}(t_k)), \quad t_k := \sum_{s=0}^{k-1} \gamma^s$$

as if DT is a perturbation of Euler discretization with variable step size.

Theorem.

Let $\bar{\mathbf{z}}$ be an **(asymptotically) stable equilibrium** for ODE.

If $\{\gamma^k\}$ goes to 0 at a **suitable rate**,

then $\{\mathbf{z}^k\}$ sequence **converges almost surely (a.s.)** to $\bar{\mathbf{z}}$.

Note: For CT-ODE analysis \rightarrow use **Linearization** and **Lyapunov theory**

6.2 Best-response and Perturbed (smooth) best-response

Problem: NE is a Fixed Point of BR, which is a set-valued map.

$$\boxed{1} \quad \mathbf{x}_i^* \in \text{BR}_i(\mathbf{x}_{-i}^*) \quad \forall i \in \mathcal{I}, \quad \Leftrightarrow \quad \mathbf{x}^* \in \text{BR}(\mathbf{x}^*)$$

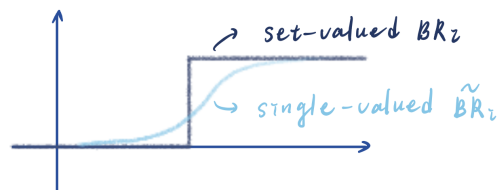
where

$$\underbrace{\text{BR}_i(\mathbf{x}_{-i})}_{\text{set-valued map}} = \arg \min_{\mathbf{x}_i \in \Delta_i} \bar{J}_i(\mathbf{x}_i, \mathbf{x}_{-i}) = \arg \max_{\mathbf{x}_i \in \Delta_i} U_i(\mathbf{x}_i, \mathbf{x}_{-i})$$

To avoid working with set-valued maps, introduce perturbed best-response:

$$\boxed{2} \quad \widetilde{\text{BR}}_i(\mathbf{x}_{-i}) = \arg \max_{\mathbf{x}_i \in \Delta_i} \underbrace{\left[\bar{U}_i(\mathbf{x}_i, \mathbf{x}_{-i}) - \varepsilon v_i(\mathbf{x}_i) \right]}_{\tilde{U}_i(\mathbf{x}_i, \mathbf{x}_{-i})}$$

- $\varepsilon > 0$ (small): perturbation
- $v_i(\mathbf{x}_i)$: **strictly convex** in \mathbf{x}_i
- $\tilde{U}_i(\mathbf{x}_i, \mathbf{x}_{-i}) = \bar{U}_i(\mathbf{x}_i, \mathbf{x}_{-i}) - \varepsilon v_i(\mathbf{x}_i)$



Smooth / perturbed best response.

$\widetilde{\text{BR}}_i$ works like a softmax function:

it transforms the set-valued fixed-point condition [1] into a single-valued fixed point [2]

$$\boxed{1} \quad \mathbf{x}^* \in \text{BR}(\mathbf{x}^*) \quad \Leftrightarrow \quad \boxed{2} \quad \mathbf{x}^* = \widetilde{\text{BR}}(\mathbf{x}^*)$$

NE distribution.

Let $\mathbf{x}^*(\varepsilon)$ be a Nash equilibrium of the perturbed game, i.e.

$$\mathbf{x}_i^*(\varepsilon) = \widetilde{\text{BR}}_i(\mathbf{x}_{-i}^*(\varepsilon)).$$

As $\varepsilon \rightarrow 0$, the NE distribution $\mathbf{x}^*(\varepsilon)$ converges to a Nash equilibrium \mathbf{x}^{NE} of the original (unperturbed) game:

$$\mathbf{x}^*(\varepsilon) \rightarrow \mathbf{x}^{NE}.$$

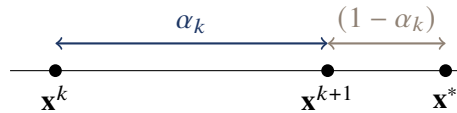
6.3 Two iterative algorithms for Repeated Finite Games:

Iterative Algorithm Design for repeated games.

- Intuition / origin
- Information: w_i^k
- Goal I: steady state of the algorithm is an NE (or closely related to an NE)
- Goal II: convergence to an NE

Question: Given that NE are fixed points of perturbed $\widetilde{\text{BR}}$, how to design a learning/update rule from \mathbf{x}^k to \mathbf{x}^{k+1} that will converge to such a fixed point?

Key idea: If \mathbf{x}^* is a fixed point of a map F (e.g. $\text{BR}, \widetilde{\text{BR}}$), then a relaxation step with step size α_k moves you along the segment from \mathbf{x}^k to \mathbf{x}^* :



$$\mathbf{x}^{k+1} = (1 - \alpha_k)\mathbf{x}^k + \alpha_k\mathbf{x}^*$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k(\mathbf{x}^* - \mathbf{x}^k)$$

Use this intuition to get iterative algorithms for $\mathbf{x}^* = f(\mathbf{x}^*)$ to get iterative algorithms based on

- [1] Best-Response Play
- [2] Perturbed Best-Response Play

6.3.1 Best-Response Play and its variants (Ch 9.2)

BR Play Algorithm.

$$\boxed{1} \quad \mathbf{x}_i^* \in \text{BR}_i(\mathbf{x}_{-i}^*)$$

$$\Rightarrow \mathbf{x}_i^{k+1} \in \mathbf{x}_i^k + \gamma^k (\text{BR}_i(\mathbf{x}_{-i}^k) - \mathbf{x}_i^k)$$

$$(\text{ODE-SA}) \Rightarrow \dot{\mathbf{x}}_i \in \underbrace{\text{BR}_i(\mathbf{x}_{-i}) - \mathbf{x}_i}_{\text{(set-valued CT dynamics)}} \quad \boxed{1^*}$$

Perturbed / Smooth BR Play Algorithm.

$$\begin{aligned}
 \boxed{2} \quad \mathbf{x}_i^* &= \widetilde{\text{BR}}_i(\mathbf{x}_{-i}^*) = \arg \max_{\mathbf{x}_i \in \Delta_i} [\bar{U}_i(\mathbf{x}_i, \mathbf{x}_{-i}) - \varepsilon v_i(\mathbf{x}_i)] \\
 \Rightarrow \quad \mathbf{x}_i^{k+1} &= \mathbf{x}_i^k + \gamma^k (\widetilde{\text{BR}}_i(\mathbf{x}_{-i}^k) - \mathbf{x}_i^k) \\
 (\text{ODE-SA}) \Rightarrow \quad \dot{\mathbf{x}}_i &= \underbrace{\widetilde{\text{BR}}_i(\mathbf{x}_{-i}) - \mathbf{x}_i}_{\text{single-valued CT dynamics}} \quad \boxed{2^*}:
 \end{aligned}$$

Note: This ODE is **coupled** across players. e.g. when $N = 2$, $\mathcal{I} = \{1, 2\}$

$$\begin{cases} \dot{\mathbf{x}}_1 = \widetilde{\text{BR}}_1(\mathbf{x}_2) - \mathbf{x}_1, \\ \dot{\mathbf{x}}_2 = \widetilde{\text{BR}}_2(\mathbf{x}_1) - \mathbf{x}_2. \end{cases}$$

Steady-State Value of the Algorithm.

$$\mathbf{x}_i^k = \bar{\mathbf{x}}_i = \widetilde{\text{BR}}_i(\bar{\mathbf{x}}_{-i}), \quad \forall i.$$

DT view from $\boxed{1^*}$

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^k + \gamma_k (\widetilde{\text{BR}}_i(\mathbf{x}_{-i}^k) - \mathbf{x}_i^k).$$

At steady state, $\mathbf{x}_i^k = \bar{\mathbf{x}}_i$ and $\mathbf{x}_{-i}^k = \bar{\mathbf{x}}_{-i}$, so

$$\bar{\mathbf{x}}_i = \bar{\mathbf{x}}_i + \gamma_k (\widetilde{\text{BR}}_i(\bar{\mathbf{x}}_{-i}) - \bar{\mathbf{x}}_i).$$

CT view from $\boxed{2^*}$:

$$\dot{\mathbf{x}}_i = \widetilde{\text{BR}}_i(\mathbf{x}_{-i}) - \mathbf{x}_i.$$

At equilibrium $\dot{\mathbf{x}}_i = 0$:

$$0 = \widetilde{\text{BR}}_i(\bar{\mathbf{x}}_{-i}) - \bar{\mathbf{x}}_i \iff \bar{\mathbf{x}} \text{ is an NE distribution.}$$

6.3.2 Fictitious Play and its Variants (Ch 9.3)

Idea of Fictitious Play.

Instead of knowing others' mixed strategies (as in BR play):

$$\mathbf{x}_{-i}^k = \{\mathbf{x}_{i'}^k, i' \neq i\}$$

player P_i will approximate them using empirical averages (frequencies) of play.

Strategy level. Instead of knowing the whole vector

$$\mathbf{x}_{-i}^k = \{\mathbf{x}_{i'}^k, i' \neq i\},$$

player P_i approximates the mixed strategy of every $i' \neq i$ from the empirical average (frequency) of past actions.

Action level. Instead of the probabilities $x_{i'j}^k = \Pr(P_{i'} \text{ selects action } j \text{ at iteration } k)$, player P_i uses the empirical average (frequency) of how many times player $P_{i'}$ has used action j up to iteration k .

$$\Rightarrow P_i \text{ denotes this approximation by } \hat{\mathbf{x}}_{i'}^k \approx \mathbf{x}_{i'}^k, \quad \forall i' \neq i.$$

Empirical Mixed Strategy (Information w_i^k).

Information available to P_i at iteration k :

$$w_i^k = \{e_{i'}^t : i' \neq i, 0 \leq t \leq k\} \quad (\text{history of all other players' actions}).$$

For each $i' \neq i$, define the empirical mixed strategy at iteration $k + 1$:

$$\boxed{3} \quad \hat{\mathbf{x}}_{i'}^{k+1} = \frac{1}{k+1} \sum_{t=0}^k e_{i'}^t,$$

where $e_{i'}^t$ is the unit-vector of action j of player i' at iteration t .

write it recursively so we don't need to track the whole history:

$$\boxed{4} \quad \hat{\mathbf{x}}_{i'}^{k+1} = \hat{\mathbf{x}}_{i'}^k + \frac{1}{k+1} (e_{i'}^k - \hat{\mathbf{x}}_{i'}^k), \quad \forall i' \neq i.$$

Now P_i use $\{\hat{\mathbf{x}}_{i'}^k\}$ as its internal variable \mathbf{z}_i^k (beliefs about others):

$$\{\hat{\mathbf{x}}_{i'}^k\}_{i' \neq i} =: \hat{\mathbf{x}}_{-i}^k =: \mathbf{z}_i^k$$

Fictitious Play Algorithm for Player i .

P_i plays a (perturbed) best response to the fictitious/approximated mixed strategy profile $\hat{\mathbf{x}}_{-i}^k$.

Internal processing (belief update from $\boxed{4}$):

$$\hat{\mathbf{x}}_{i'}^{k+1} = \hat{\mathbf{x}}_{i'}^k + \frac{1}{k+1} (e_{i'}^k - \hat{\mathbf{x}}_{i'}^k), \quad \forall i' \neq i.$$

Strategy update (best response to beliefs):

$$\mathbf{x}_i^k = \widetilde{\text{BR}}_i(\hat{\mathbf{x}}_{-i}^k), \quad k = 0, 1, 2, \dots$$

All Other Players Also Use Fictitious Play.

Assume every player i' uses the Fictitious Play strategy update above. Then, at iteration k , the probability of selecting action j

$$\boxed{5} \quad \mathbb{P}(e_{i'}^k = e_{i'j} \mid \hat{\mathbf{x}}^k) = x_{i'j}^k$$

The conditional expectation of the action vector is

$$\begin{aligned} \boxed{6} \quad \mathbb{E}[e_{i'}^k \mid \hat{\mathbf{x}}^k] &= \sum_j e_{i'j} \mathbb{P}(e_{i'}^k = e_{i'j} \mid \hat{\mathbf{x}}^k) \\ &= \sum_j x_{i'j}^k e_{i'j} = \mathbf{x}_{i'}^k = \widetilde{\text{BR}}_{i'}(\hat{\mathbf{x}}_{-i'}^k). \end{aligned}$$

Recall: stochastic approximation.

$$\mathbf{z}^{k+1} = \mathbf{z}^k + \gamma^k (f(\mathbf{z}^k) + \xi^k) \quad \text{where} \quad \begin{cases} \gamma^k > 0, \sum_k \gamma^k = \infty, \gamma^k \rightarrow 0, \\ \mathbb{E}[\xi^k | \mathcal{F}_k] = 0. \end{cases}$$

$$\underbrace{\mathbb{E} \left[\frac{\mathbf{z}^{k+1} - \mathbf{z}^k}{\gamma^k} \middle| \mathcal{F}_k \right]}_{\text{compute it for [4]}} = f(\mathbf{z}^k) \quad \rightsquigarrow \quad \dot{\mathbf{z}} = f(\mathbf{z}).$$

Let

$$\gamma^k = \frac{1}{k+1}, \quad \mathbf{z}^k = \hat{\mathbf{x}}^k,$$

and recall the belief update [4] (for each player i):

$$\hat{\mathbf{x}}_i^{k+1} = \hat{\mathbf{x}}_i^k + \frac{1}{k+1} (e_i^k - \hat{\mathbf{x}}_i^k).$$

Then

$$\frac{\hat{\mathbf{x}}_i^{k+1} - \hat{\mathbf{x}}_i^k}{\gamma^k} = e_i^k - \hat{\mathbf{x}}_i^k.$$

Taking conditional expectation (given all past info, summarized by $\hat{\mathbf{x}}^k$):

$$\begin{aligned} \mathbb{E} \left[\frac{\hat{\mathbf{x}}_i^{k+1} - \hat{\mathbf{x}}_i^k}{\gamma^k} \middle| \hat{\mathbf{x}}^k \right] &= \mathbb{E}[e_i^k | \hat{\mathbf{x}}^k] - \hat{\mathbf{x}}_i^k \\ &= \underbrace{\mathbb{E}[e_i^k | \hat{\mathbf{x}}^k]}_{\text{[6]} = \mathbf{x}_i^k = \widetilde{\text{BR}}_i(\hat{\mathbf{x}}_{-i}^k)} - \hat{\mathbf{x}}_i^k \\ &= \widetilde{\text{BR}}_i(\hat{\mathbf{x}}_{-i}^k) - \hat{\mathbf{x}}_i^k =: f_i(\hat{\mathbf{x}}^k). \end{aligned}$$

By the SA–ODE correspondence,

$$\dot{\hat{\mathbf{x}}}_i = \widetilde{\text{BR}}_i(\hat{\mathbf{x}}_{-i}) - \hat{\mathbf{x}}_i, \quad \forall i.$$

Comparison with BR Play.

Perturbed BR Play:

$$\dot{\mathbf{x}}_i = \widetilde{\text{BR}}_i(\mathbf{x}_{-i}) - \mathbf{x}_i, \quad \forall i.$$

Fictitious Play:

$$\dot{\hat{\mathbf{x}}}_i = \widetilde{\text{BR}}_i(\hat{\mathbf{x}}_{-i}) - \hat{\mathbf{x}}_i, \quad \forall i.$$

- BR Play: state \mathbf{x}_i = actual mixed strategy.
- FP: state $\hat{\mathbf{x}}_i$ = belief / empirical frequency.

In both cases the CT–ODE is a relaxation toward a fixed point $\mathbf{x}^* = \widetilde{\text{BR}}(\mathbf{x}^*)$, i.e. a (perturbed) NE, but BR play moves the strategies, while FP moves the beliefs.