

Chapter 10: Multi-agent RL (Continued)

Recap: Normal-form Game

A **normal-form game** is a tuple $(n, \mathcal{A}_{1\dots n}, R_{1\dots n})$,

- n is the number of players,
- \mathcal{A}_i is the set of actions available to player i
 - \mathcal{A} is the joint action space $\mathcal{A}_1 \times \dots \times \mathcal{A}_n$,
- R_i is player i 's payoff function $\mathcal{A} \rightarrow \mathfrak{R}$.

$$R_1 = \left(\begin{array}{c} a_2 \\ \vdots \\ \vdots \\ \vdots \\ a_1 \left(\begin{array}{ccc} \cdots & R_1(a) & \cdots \end{array} \right) \\ \vdots \\ \vdots \\ \vdots \end{array} \right)$$
$$R_2 = \left(\begin{array}{c} a_2 \\ \vdots \\ \vdots \\ \vdots \\ a_1 \left(\begin{array}{ccc} \cdots & R_2(a) & \cdots \end{array} \right) \\ \vdots \\ \vdots \\ \vdots \end{array} \right)$$

Minimax Optimal Solution

- Play strategy with the best worst-case outcome.

$$\operatorname{argmax}_{\sigma_i \in \Delta(\mathcal{A}_i)} \min_{a_{-i} \in \mathcal{A}_{-i}} R_i(\langle \sigma_i, \sigma_{-i} \rangle)$$

- How to compute it?
- Linear programming [Whiteboard Example].

Nash Equilibria

- A **best response set** is the set of all strategies that are optimal given the strategies of the other players.

$$\text{BR}_i(\sigma_{-i}) = \{\sigma_i \mid \forall \sigma'_i \quad R_i(\langle \sigma_i, \sigma_{-i} \rangle) \geq R_i(\langle \sigma'_i, \sigma_{-i} \rangle)\}$$

- A **Nash equilibrium** is a joint strategy, where all players are playing best responses to each other.

$$\forall i \in \{1 \dots n\} \quad \sigma_i \in \text{BR}_i(\sigma_{-i})$$

- **Nash = Minimax** in Two-Player **Zero-sum** games, but not always [Whiteboard Example].

Existence of Nash Equilibria

- All finite normal-form games have at least one Nash equilibrium. (Nash, 1950)
- In zero-sum games...
 - Equilibria all have the same value and are interchangeable.

$\langle \sigma_1, \sigma_2 \rangle, \langle \sigma'_1, \sigma'_2 \rangle$ are Nash $\Rightarrow \langle \sigma_1, \sigma'_2 \rangle$ is Nash.

- Equilibria correspond to minimax optimal strategies.

Computation of Nash Equilibria

- The exact complexity of computing a Nash equilibrium is an open problem. (Papadimitriou, 2001)

The Complexity of Computing a Nash Equilibrium*

Constantinos Daskalakis
Computer Science Division,
UC Berkeley
costis@cs.berkeley.edu

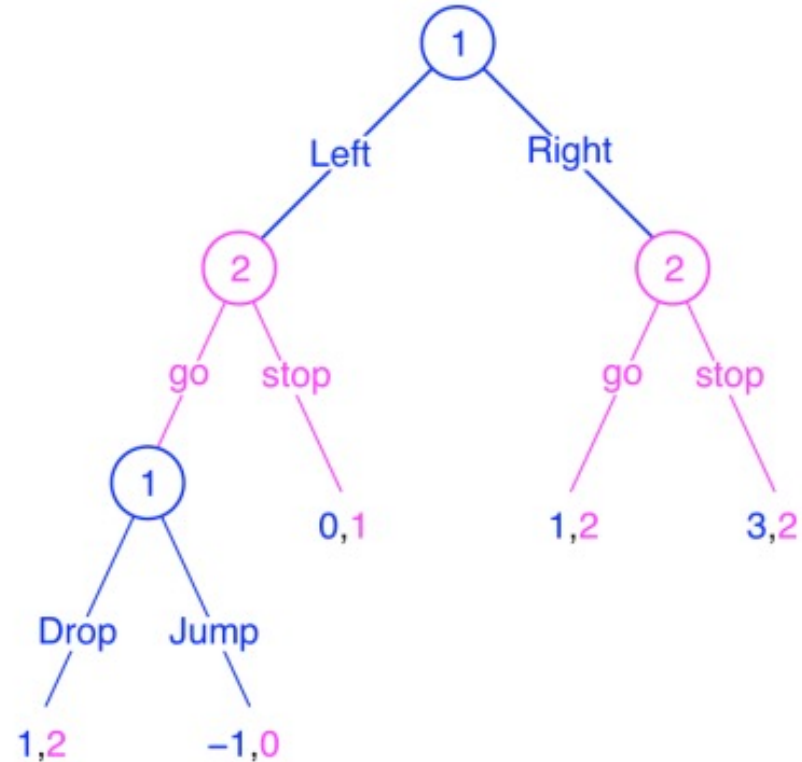
Paul W. Goldberg
Dept. of Computer Science,
University of Liverpool
P.W.Goldberg@liver-
pool.ac.uk

Christos H. Papadimitriou
Computer Science Division,
UC Berkeley
christos@cs.berkeley.edu

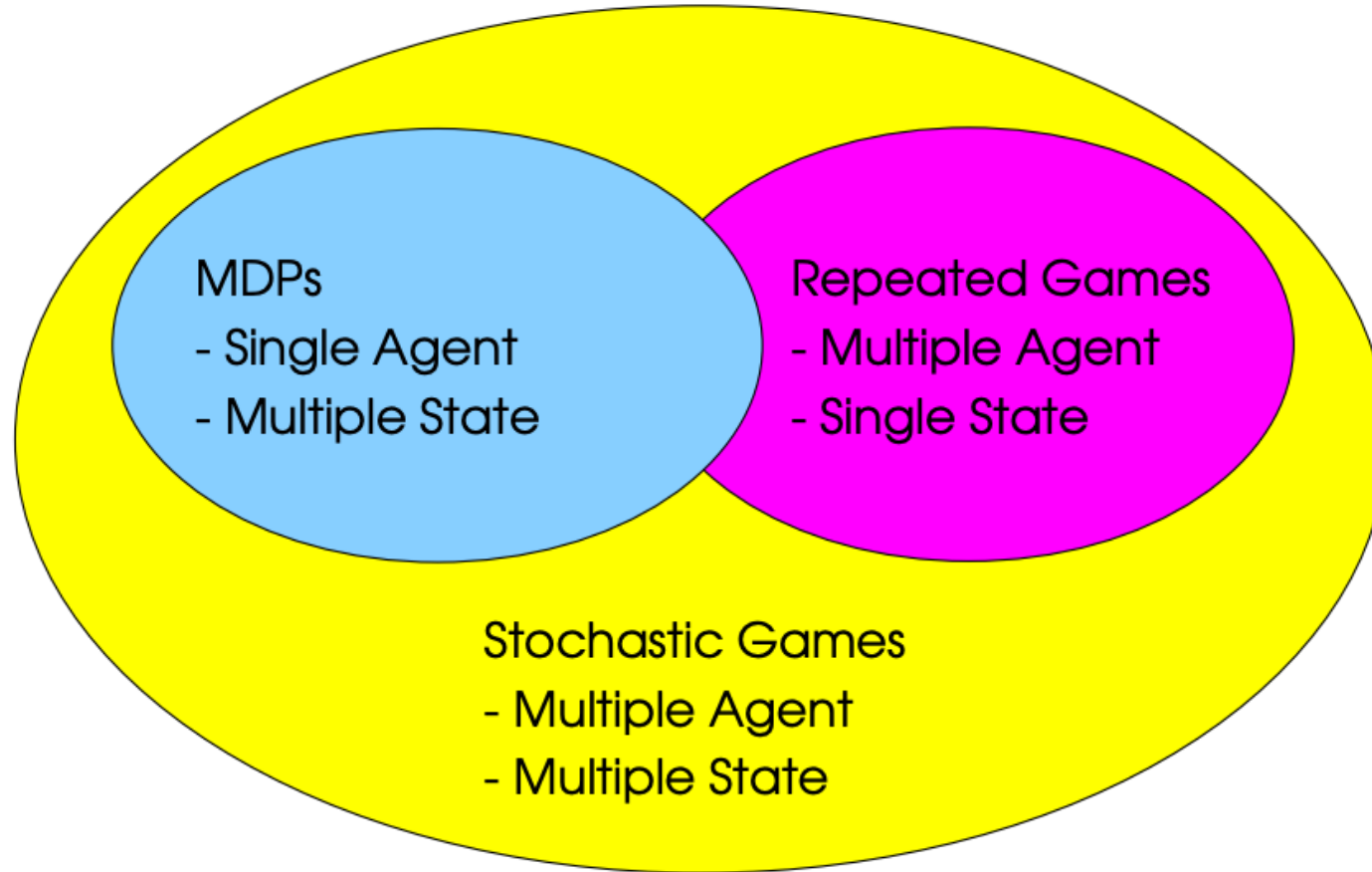
- Nash-equilibrium is PPAD-hard [2008].

Extensive-form Game

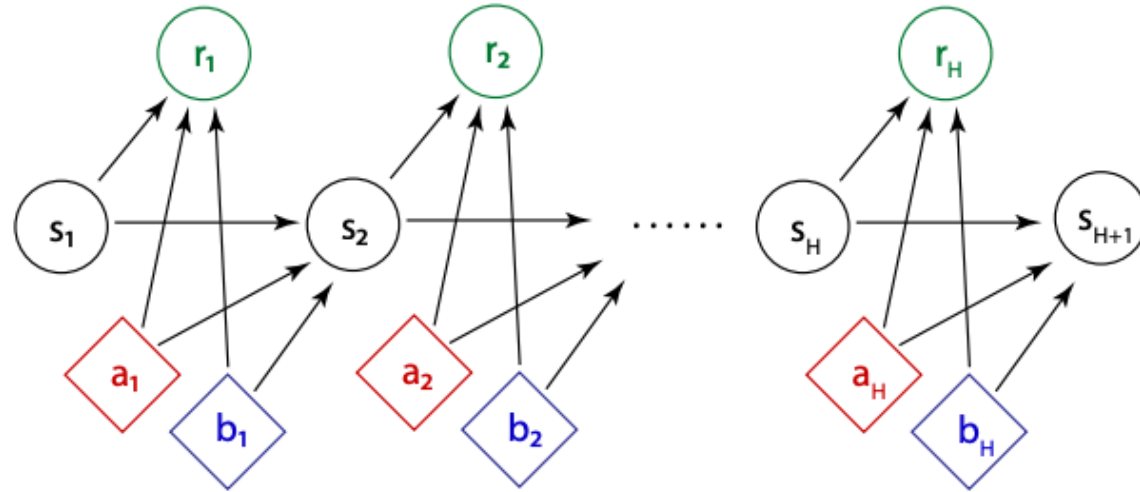
- Example: any full-observation turn-based games, e.g. Chess, Go.



Stochastic/Markov Games



Stochastic/Markov Games



Two-player zero-sum Markov Game $(\mathcal{S}, \mathcal{A}, \mathcal{B}, \mathbb{P}, r, H)$ [Shapley 1953].

- \mathcal{S} : set of **states**; \mathcal{A}, \mathcal{B} : set of **actions** for the max-player/the min-player.
- $\mathbb{P}_h(s_{h+1}|s_h, a_h, b_h)$: **transition** probability.
- $r_h(s_h, a_h, b_h) \in [0, 1]$: **reward** for the max-player (**loss** for the min-player).
- H : horizon/the length of the game.

Our Setup

- **Fully observable**: joint actions and states are revealed to both agents.
- **Tabular**: the size of $\mathcal{S}, \mathcal{A}, \mathcal{B}$ is finite and small.

Policy and Value

- **General policy** for the max-player (depends on the **entire history**):

$$\pi_{1,h} : (\mathcal{S} \times \mathcal{A} \times \mathcal{B})^{h-1} \times \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$$

- **Markov policy** for the max-player (depends on the **current state**):

$$\pi_{1,h} : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$$

Policy of the min-player can be defined by symmetry.

- **Value** V^π for joint policy $\pi = (\pi_1, \pi_2)$: the expected cumulative reward received by the max-player if both agents follow the joint policy π :

$$V^\pi = \mathbb{E}_\pi \left[\sum_{h=1}^H r_h(s_h, a_h, b_h) \right]$$

Nash Equilibria

Nash Equilibria

The policies (π_1^*, π_2^*) is a **Nash equilibrium** if no player has incentive to deviate from her current policy. That is, for any π_1, π_2

$$V^{\pi_1, \pi_2^*} \leq V^{\pi_1^*, \pi_2^*} \leq V^{\pi_1^*, \pi_2}$$

In two-player zero-sum Markov games, **minimax theorem** holds:

$$\max_{\pi_1} \min_{\pi_2} V^{\pi_1, \pi_2} = \min_{\pi_2} \max_{\pi_1} V^{\pi_1, \pi_2}$$

Nash Equilibria

The optimal strategy if always facing best responses.

“We may not win by a large margin, but no one beats us.”

Objective: find ϵ -approximate Nash equilibria $(\hat{\pi}_1, \hat{\pi}_2)$ using a small number of samples with mild dependency on S, A_1, A_2, ϵ, H .

$$\max_{\pi_1} V^{\pi_1, \hat{\pi}_2} - \min_{\pi_2} V^{\hat{\pi}_1, \pi_2} \leq \epsilon.$$

Technical Challenges

To name a few:

- Large size of **policy space**:

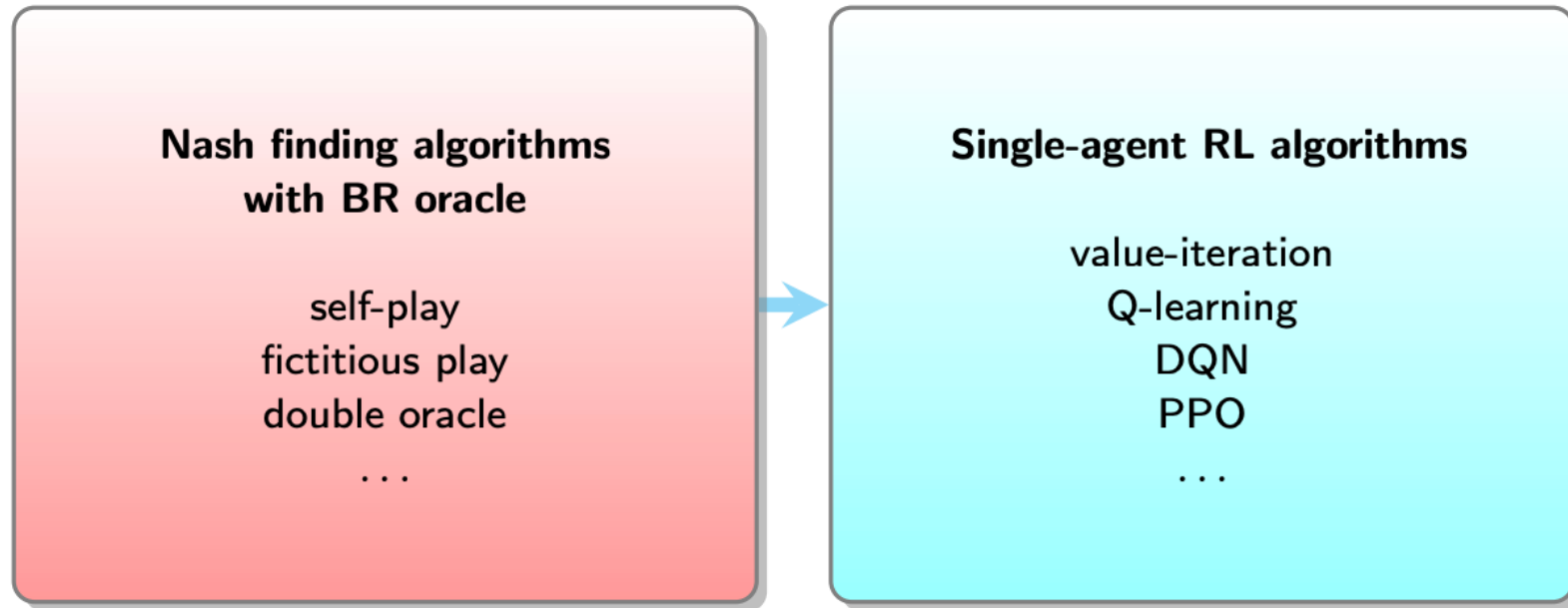
$\Omega((1/\epsilon)^{HSA})$ **Markov** policies in the **tabular** setting

- **Nash equilibrium policy is Markov**, but the best response may **not** be.
- **MGs do not allow efficient no-regret learning** [Bai, Jin, Yu, 2020].

$$\max_{\pi_1} \sum_{t=1}^T V_1^{\pi_1 \times \pi_2^t} - \sum_{t=1}^T V_1^{\pi_1^t \times \pi_2^t} \leq \text{poly}(H, S, A, B) T^{1-\alpha}.$$

Computing NE in Zero-sum Markov Games: “anecdotal Recipe”

Key observation: **given a fixed opponent, computing best response (BR) is a single-agent RL problem.**



commonly used in practice.

Computing NE in Zero-sum Markov Games

Fictitious play [Brown, 1949]

for $k = 1, \dots, K$,

$$\pi_1^{k+1} = BR[(1/k) \cdot (\pi_2^1 + \dots + \pi_2^k)].$$

$$\pi_2^{k+1} = BR[(1/(k+1)) \cdot (\pi_1^1 + \dots + \pi_1^{k+1})].$$

π_i^k : the policy of the i^{th} player at the k^{th} iteration

Computing the best response to the average policy of the opponent.

Computing NE in Zero-sum Markov Games

Asymptotic convergence of fictitious play [Robinson 1951]

Fictitious play indeed converges to Nash equilibrium!

However, how **fast**?

- inspecting the proof of [Robinson 1951], it requires $(1/\epsilon)^{\Omega(A)}$ iterations to converge to ϵ -Nash equilibrium for a normal-form game with A actions.
- Karlin conjectured in 1959 that this rate can be improved to $\mathcal{O}(1/\epsilon^2)$.
- Daskalakis and Pan [2014] **refute** the conjecture, and prove that $(1/\epsilon)^{\Omega(A)}$ **is real** in the worst case.

Drawbacks of Direct Combinations

- Algorithms are designed based on black-box usage of single-agent RL, which **does not exploit** the **detailed structure of MGs**.
- Converting a MG into a norm-form game gives a number of action $A = (1/\epsilon)^{HSA'}$.
- Finding BR is **NOT** a easy single-agent RL problem:
 - When the min-player deploys a fixed **non-Markovian** policy, the game is **NOT** an MDP from the perspective of the max-player.
 - Existing single-agent RL results do not apply.

Planning in Markov Games

We start with the setting of known transition \mathbb{P} and reward r .

A Nash equilibrium of a MG is a Markov policy.

We define $V_h^*(s)$, $Q_h^*(s, a, b)$ which satisfies the **Bellman optimality equation**:

$$Q_h^*(s, a, b) = r_h(s, a, b) + \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a, b)} V_{h+1}^*(s')$$
$$V_h^*(s) = \max_{\mu \in \Delta_{\mathcal{A}}} \min_{\nu \in \Delta_{\mathcal{B}}} \sum_{a, b} \mu(a) \nu(b) Q_h^*(s, a, b)$$
$$:= \text{Nash_Value}(Q_h^*(s, \cdot, \cdot))$$

Planning in Markov Games

A dynamical programming approach to find a Nash equilibrium.

Nash Value Iteration (Nash VI)

Initialize $V_{H+1}^*(s) = 0$ for all s .

for $h = H, \dots, 1$,

for all (s, a, b) ,

$$Q_h^*(s, a, b) \leftarrow r_h(s, a, b) + \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a, b)} V_{h+1}^*(s')$$

for all s

$$(\pi_{1,h}^*(\cdot | s), \pi_{2,h}^*(\cdot | s)) \leftarrow \text{Nash}(Q_h^*(s, \cdot, \cdot))$$

$$V_h^*(s) \leftarrow \langle \pi_{1,h}^*(\cdot | s) \times \pi_{2,h}^*(\cdot | s), Q_h^*(s, \cdot, \cdot) \rangle$$

Nash VI computes the Nash equilibrium of MGs in $\text{poly}(H, S, A, B)$ steps!