











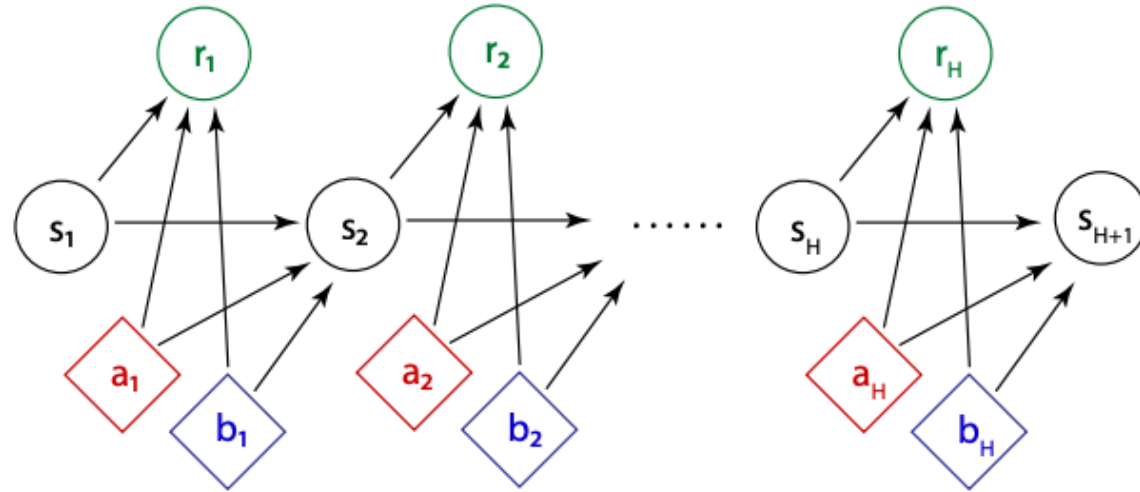


# Chapter 10: Multi-agent RL (Continued)

# Reminder: Course Project due next Tuesday

#	Team	Members	Score	Agents	Last	Join
1	Team GO		3000.8	2 	1d	
2	Team Q		2523.0	2 	1d	
3	Team S		2439.3	2 	17d	
4	Team S_1		2347.4	2 	22d	
5	MilesLiii		2019.9	2 	6h	
6	Team Lux		1464.8	2 	1d	

# 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.

# 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**  $(s, a, b)$ ,

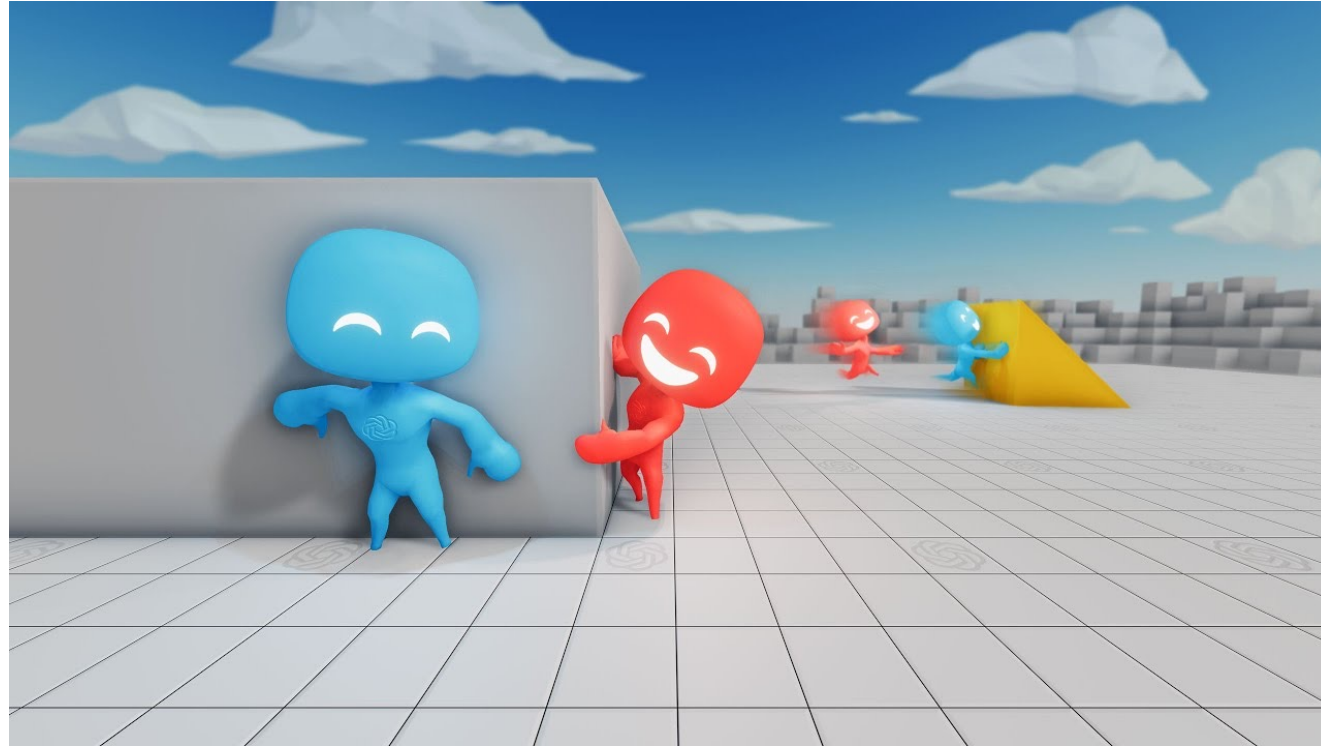
$$Q_h^*(s, a, b) \leftarrow r_h(s, a, b) + \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a, b)} V_{h+1}^*(s')$$

**for all**  $s$

$$(\pi_{1,h}^*(\cdot | s), \pi_{2,h}^*(\cdot | s)) \leftarrow \text{Nash}(Q_h^*(s, \cdot, \cdot)) \quad \text{NE for Normal-form Game}$$

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

# Today: Online Learning in Unknown MGs



How do we explore in an unknown Markov Game to learn an  $\epsilon$ -Nash strategy?

# Online Learning in Unknown MGs

## Interaction protocol

Environment samples initial state  $s_1$ .

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

two agents take their own **actions**  $(a_h, b_h)$  simultaneously.

both agents receive their own immediate **reward**  $\pm r_h(s_h, a_h, b_h)$ .

environment **transitions** to the next state  $s_{h+1} \sim \mathbb{P}_h(\cdot | s_h, a_h, b_h)$ .

# Recall UCBVI for Single-agent RL

**Inside iteration  $n$  :**

Use all previous data to estimate transitions  $\hat{P}^n$

Design reward bonus  $b_h^n(s, a), \forall s, a, h$

Optimistic planning with learned model:  $\pi^n = \text{Value-Iter} \left( \hat{P}^n, \{r_h + \underbrace{b_h^n}_{\text{Optimism}}\}_{h=1}^{H-1} \right)$

Collect a new trajectory by executing  $\pi^n$  in the real world  $P$  starting from  $s_0$

🤔 How do we achieve optimism in Two-Player Zero-sum MG?

# How do we modify Nash-VI?

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**  $s$

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

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



# Optimistic Nash-VI

## Optimistic Nash VI

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

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

for all  $(s, a, b)$ ,

$$\bar{Q}_h(s, a, b) \leftarrow r_h(s, a, b) + \mathbb{E}_{s' \sim \hat{\mathbb{P}}_h(\cdot | s, a, b)} \bar{V}_{h+1}(s') + \beta$$

$$\underline{Q}_h(s, a, b) \leftarrow r_h(s, a, b) + \mathbb{E}_{s' \sim \hat{\mathbb{P}}_h(\cdot | s, a, b)} \underline{V}_{h+1}(s') - \beta$$

for all  $s$

$$\pi_h(\cdot, \cdot | s) \leftarrow \text{CCE}(\bar{Q}_h(s, \cdot, \cdot), \underline{Q}_h(s, \cdot, \cdot)) \quad \text{CCE instead of NE}$$

$$\bar{V}_h(s) \leftarrow \langle \pi_h(\cdot, \cdot | s), \bar{Q}_h(s, \cdot, \cdot) \rangle$$

$$\underline{V}_h(s) \leftarrow \langle \pi_h(\cdot, \cdot | s), \underline{Q}_h(s, \cdot, \cdot) \rangle$$

execute policy  $\pi$ , collect samples, and update estimation  $\hat{\mathbb{P}}$ .

Reward bonus:

$$\beta = \mathcal{O}\left(\sqrt{\frac{1}{N_k(s, a, b)}}\right)$$

# Coarse Correlated Equilibria

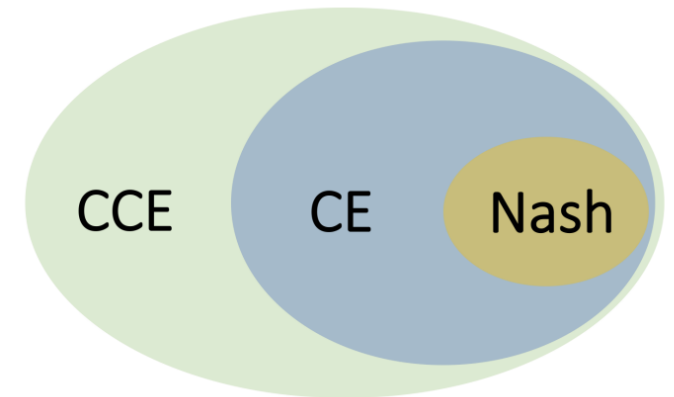
- **Coarse Correlated Equilibria (CCE):** A joint policy  $\pi: S \rightarrow A \times B$  is a CCE if

$$\max_{\pi': S \rightarrow A} V^{\pi', \pi_{-1}} \leq V^{\pi} \quad \text{and} \quad \max_{\pi': S \rightarrow B} V^{\pi_{-2}, \pi'} \geq V^{\pi}$$

- **CCE v.s. NE:**

- CCE allows correlated polices, e.g. traffic light.

	STOP	GO
STOP	(0,0)	(0,1)
GO	(1,0)	(-100,-100)



- CCE is efficiently computable for general-sum games, while NE isn't.

# Theoretical Guarantee of Nash-VI

**Theorem [Liu, Yu, Bai, Jin 2020]**

With high probability, optimistic Nash VI finds an  $\epsilon$ -Nash equilibrium in  $\tilde{O}(H^3 SAB/\epsilon^2)$  episodes.

$H$ : horizon;  $S$ : number of states;  $A, B$ : number of actions for each player.

**Optimistic Nash VI finds  $\epsilon$ -Nash in polynomial time and samples!**

# Drawbacks of Nash-VI

- Centralized learning: Requires keeping track of  $Q(s, a, b)$ .
- The algorithm can be generalized to the multi-agent setting:
- Nash-VI finds an  $\epsilon$ -CCE with  $\mathcal{O}(\text{poly}(S \prod_{i=1}^n A_i))$  sample and computational complexity.
- “The Curse of Multi-agent”:  $\prod_{i=1}^n A_i$  scaling

# The Curse of Multi-agent

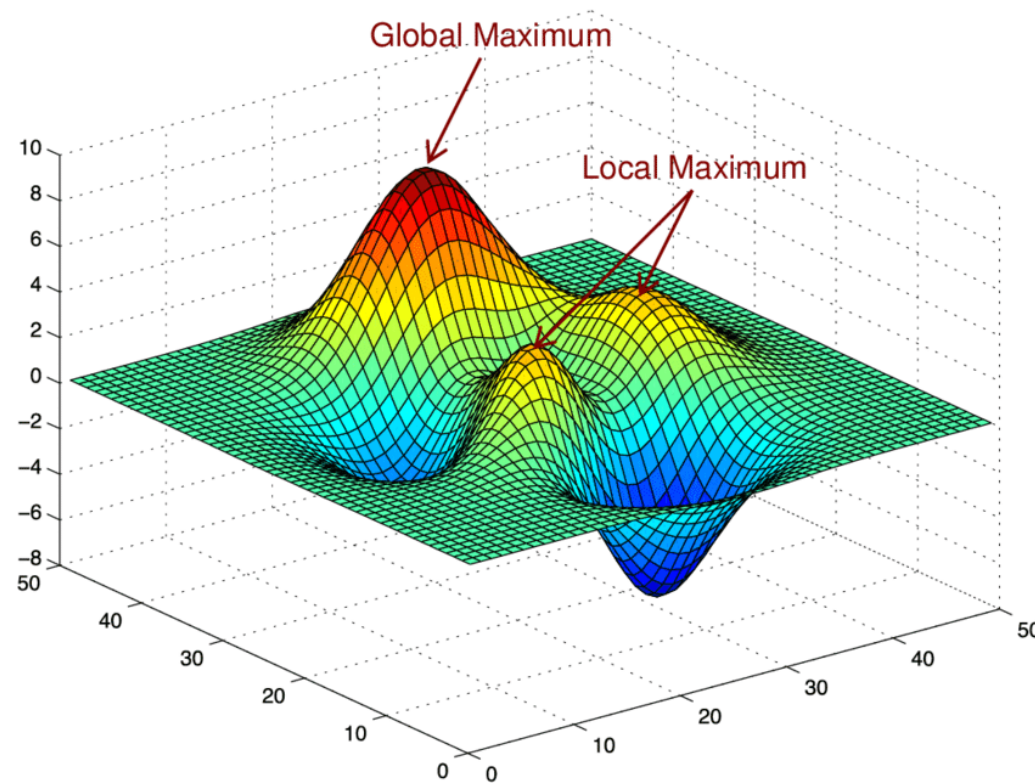
- Can we avoid the  $O(AB)$  scaling?

Information theoretical lower bound:  $\Omega(H^3 S \max\{A, B\} / \epsilon^2)$

- Observation: Nash-VI requires estimating the Q function with  $SAB$  entries, naturally resulting in the scaling with  $O(SAB)$ .

# The Curse of Multi-agent

- But why can we avoid trying each  $(s, a, b)$  tuple at least once?



# Simpler Setting: Normal-form Game

Each agent runs **no-regret algorithm** for adversarial bandit (e.g. EXP3) independently.

$$\sum_{t=1}^T \langle \mu_t, \ell_t \rangle - \min_{a \in \mathcal{A}} \sum_{t=1}^T \langle a, \ell_t \rangle \leq \text{poly}(A) T^{1-\alpha}.$$

- two-player zero-sum games:  $(\mathbb{E}_{t \sim \text{Unif}(T)} \mu_t^{(1)}) \times (\mathbb{E}_{t \sim \text{Unif}(T)} \mu_t^{(2)}) \rightarrow \text{Nash}.$
- sample complexity scales with  $\tilde{O}(A + B).$

Unfortunately, cannot run no-regret algorithm in MGs (recall from last lecture).

# V-Learning

**V-learning [Bai, Jin, Yu, 2020] [Jin, Liu, Wang, Yu, 2021]**

for  $k = 1, \dots, K$ , receive  $s_1$ ,

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

take action  $a_h \sim \pi_h(\cdot | s_h)$ , observe reward  $r_h$  and next state  $s_{h+1}$ .

$t = N_h(s_h) \leftarrow N_h(s_h) + 1$ .

$V_h(s_h) \leftarrow (1 - \alpha_t)V_h(s_h) + \alpha_t(r_h + V_{h+1}(s_{h+1}) + \beta_t)$ .

$\pi_h(\cdot | s_h) \leftarrow \text{Adv\_Bandit\_Update}(a_h, r_h + V_{h+1}(s_{h+1}))$

on the  $(s_h, h)^{\text{th}}$  adversarial bandit.

- Incremental updates of  $V$  instead of  $Q$ !
- Is a single-agent algorithm.



# Theoretical Guarantee

- Multiagent setting: **both agents run V-learning independently.**
- Adversarial bandit subroutine: **FTRL.**

## Theorem [Bai, **Jin**, Yu, 2020]

In two-player zero-sum Markov games, V-learning with FTRL finds  **$\epsilon$ -Nash** in  $\tilde{O}(H^5 S \max\{A, B\}/\epsilon^2)$  episodes.

V-learning is a **decentralized** algorithm that achieves **optimal  $O(\max\{A, B\})$**  sample complexity!

# Readily Generalize to Multi-agent MGs

**Theorem (CCE & CE) [Song et al. 2021][Jin, Liu, Wang, Yu, 2021]**

In general-sum Markov games,

(1) V-learning with **FTRL** finds  $\epsilon$ -**CCE** in  $\tilde{O}(H^5 S(\max_{i \in [m]} A_i)/\epsilon^2)$  episodes;

(2) V-learning with **FTRL\_swap** finds  $\epsilon$ -**CE** in  $\tilde{O}(H^5 S(\max_{i \in [m]} A_i)^2/\epsilon^2)$

episodes.

# Summary of Algorithms

Algorithm	Training	Main estimand	Sample complexity
Nash-VI	centralized	$\mathbb{P}_h(s' s, a, b)$	$\tilde{O}(H^3 SAB/\epsilon^2)$
Nash Q-Learning	centralized	$Q_h^*(s, a, b)$	$\tilde{O}(H^5 SAB/\epsilon^2)$
V-Learning	decentralized	$V_h^*(s)$	$\tilde{O}(H^5 S \max\{A, B\}/\epsilon^2)$
Lower bound	-	-	$\Omega(H^3 S \max\{A, B\}/\epsilon^2)$

# Lots of Future Work to be done

- Behavior of Decentralized Algorithms.
- Policy Gradient for Markov Games?
- Scalable algorithms? (closing theory-practice gap)
- Imperfect Information Markov Games.