

# When is an Offline Two-Player Zero-Sum Markov Game Solvable?

---

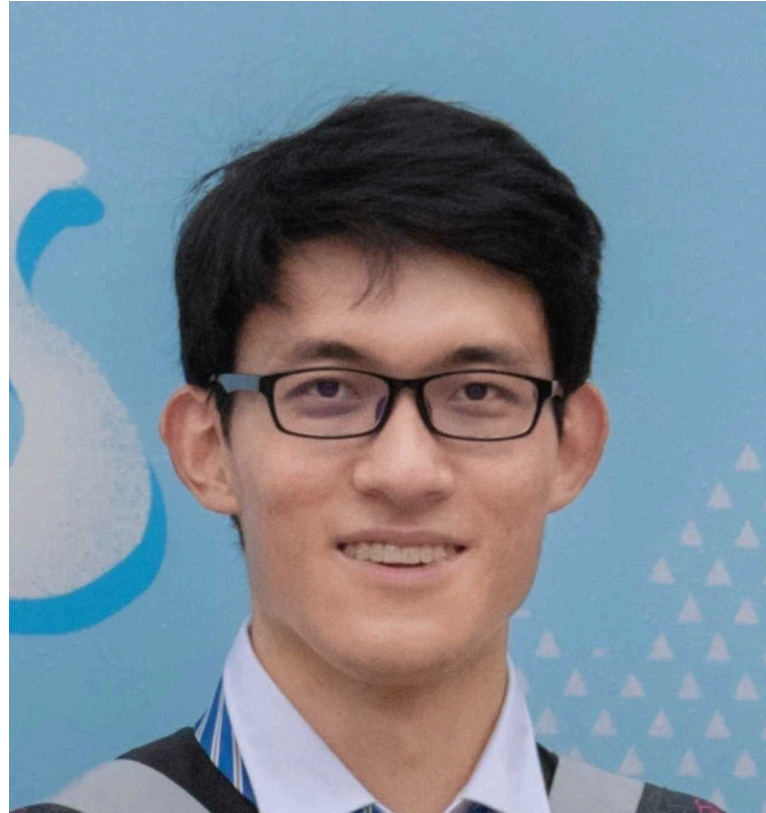
**Simon S. Du**

05-04-2022

Simons Institute

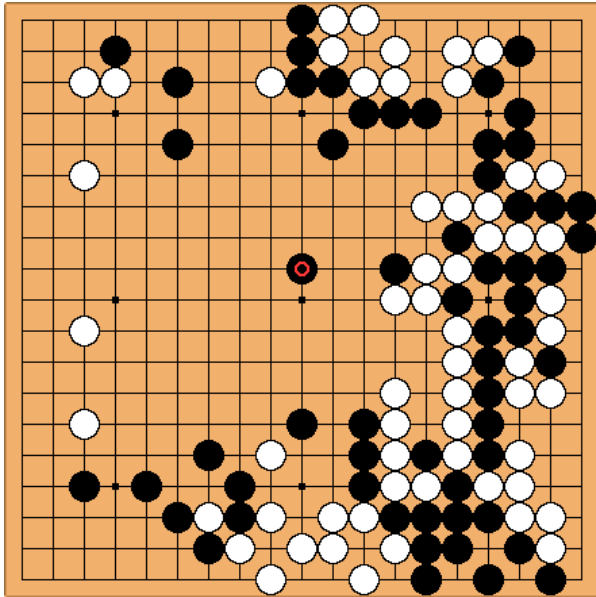
# Acknowledgement

---



Qiwen Cui  
University of Washington

# Two-Player Zero-Sum Markov Games



- Two players compete against each other. Each has a strategy.
- **Goal:** find a Nash Equilibrium
- Nash Equilibrium: a pair of strategies that no player can do better by unilaterally changing the policy.
- Applications: poker, Go, chess, computer games, investment, .....

# Offline Reinforcement Learning

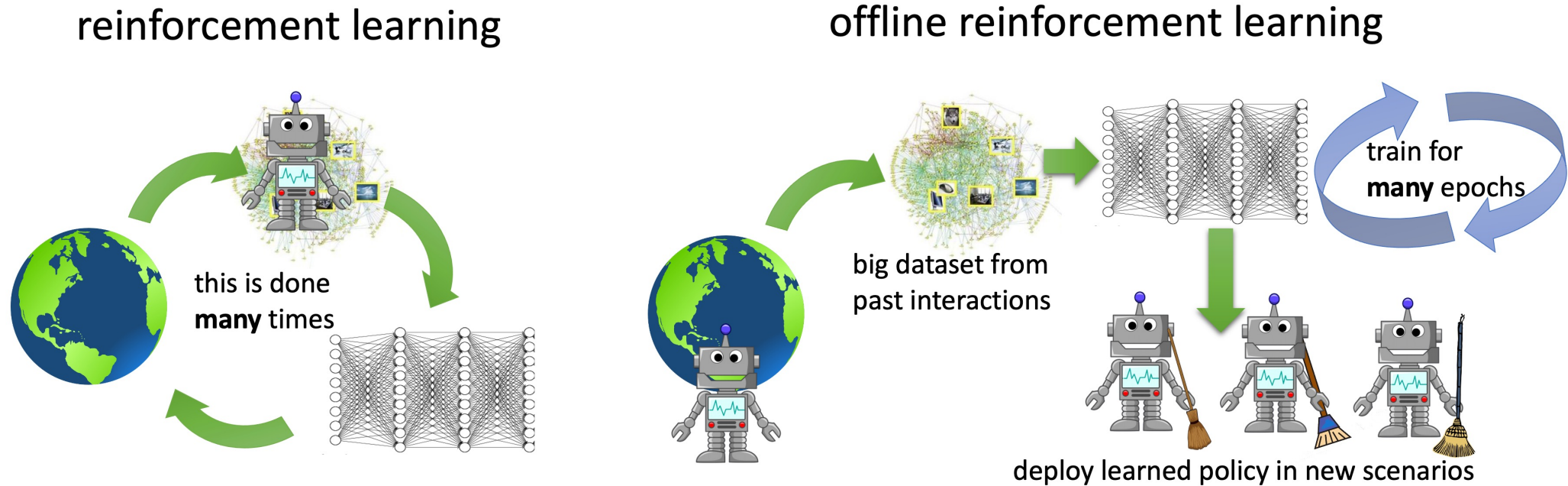
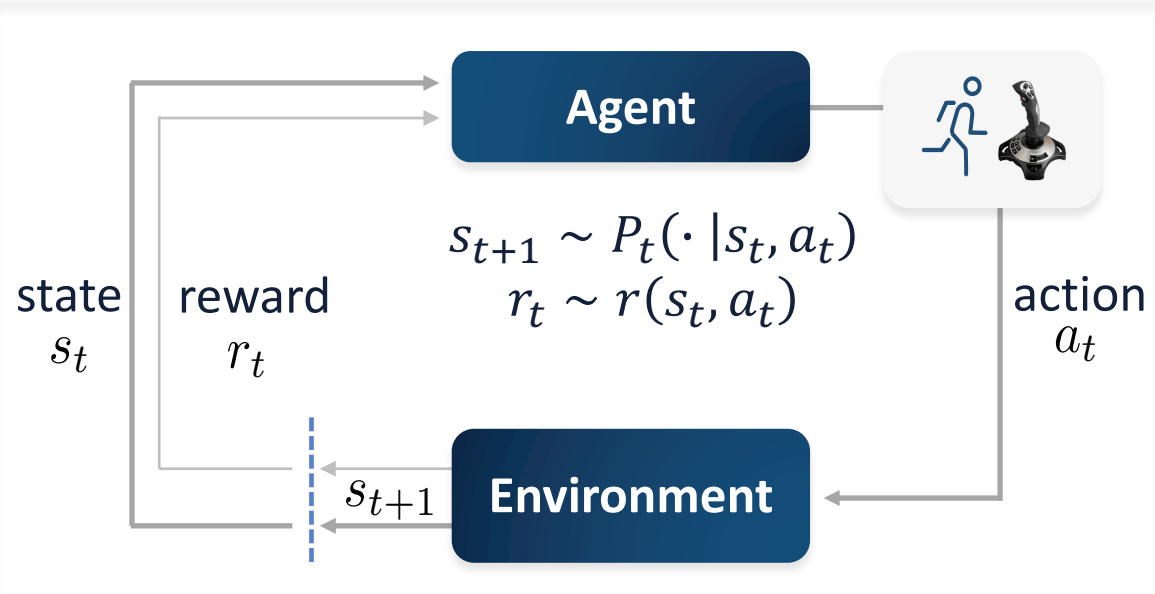


Figure credit: Berkeley AI Research Blog

- Lots of available **offline data** from prior experience. Fresh samples are expensive
- **This Talk:** When can we learn a Nash Equilibrium in offline two-player zero-sum Markov games?

# Single-Agent Reinforcement Learning



Repeat **H** times

**H**: planning horizon / Episode length

A policy  $\pi$  :

$\pi: \text{States}(S) \rightarrow \text{Actions}(A), a = \pi(s)$

Goal: maximize value function

$$V^\pi(s_1) = \mathbb{E}[r_1 + r_2 + \dots + r_H]$$

Near-optimal policy:

$$V^*(s_1) - V^\pi(s_1) \leq \epsilon$$

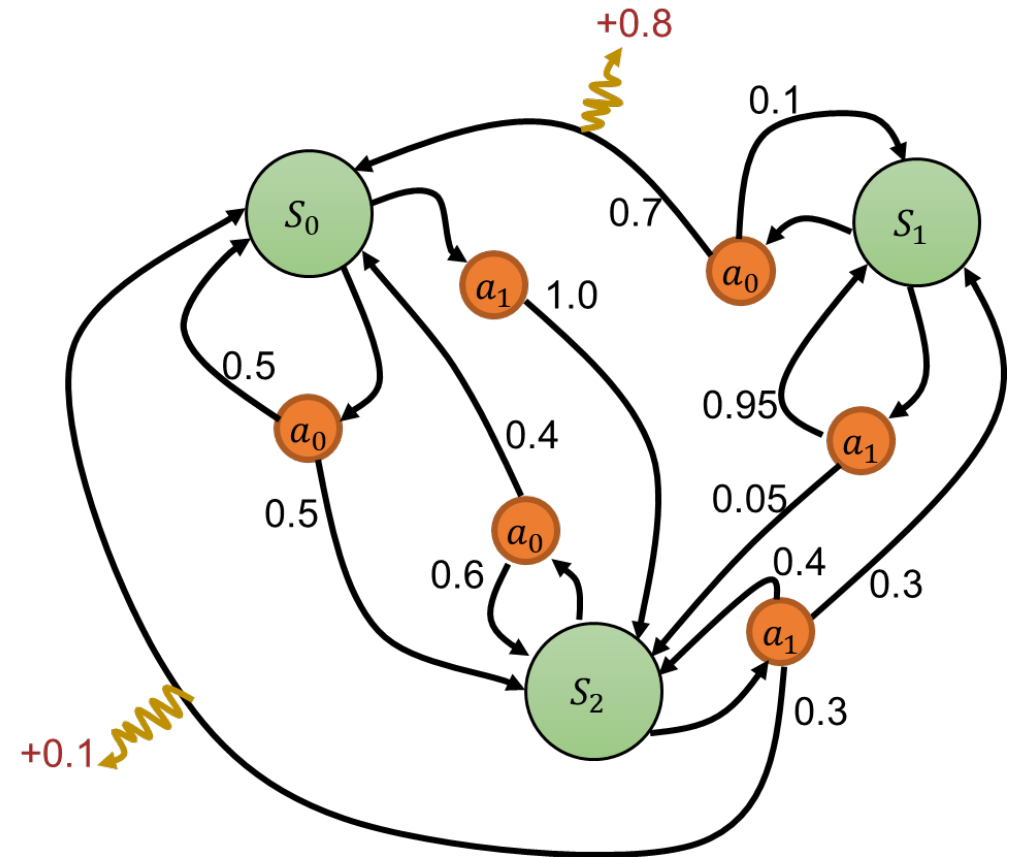
$V^* = V^{\pi^*}$  : value function of opt policy

# Tabular Markov Decision Process

## Assumptions:

1. # of States  $\mathcal{S} < \infty$
2. # of actions  $\mathbf{A} < \infty$
3. Bounded rewards:  
 $0 \leq r_h \leq 1, h = 1, \dots, H$

Sample complexity depends on  
 $(S, A, H, 1/\epsilon)$



# Offline Single-Agent Reinforcement Learning

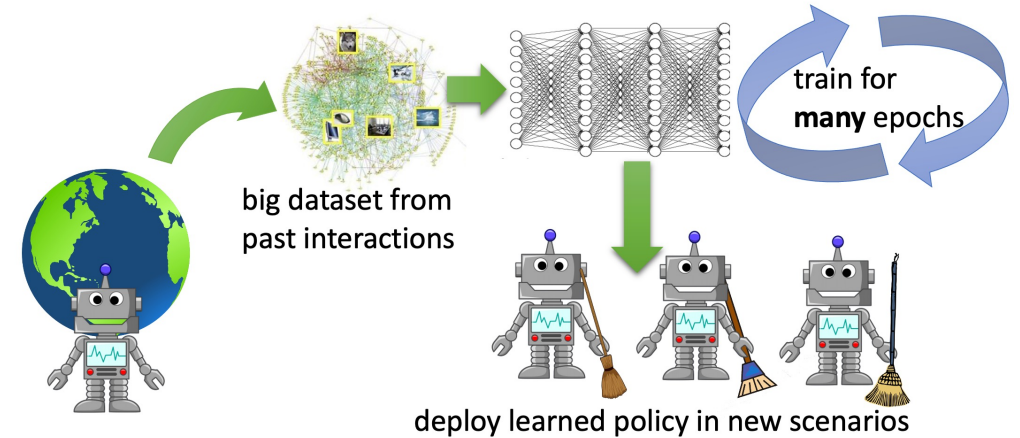
Offline Data:  $n$  (state, action, reward, next state) tuples:

$$D = \left\{ (s_h^i, a_h^i, r_h^i, s_{h+1}^i) \right\}_{h \in [H]}^{i \in [n]} \stackrel{i.i.d.}{\sim} d^\rho$$

- $\rho$  is the data-collection / behavior policy
- $d_h^\rho(s, a)$  is the state-action distribution induced by  $\rho$  and transition  $P$ .
- Goal: learn a policy  $\pi$  from  $D$  such that

$$V^*(s_1) - V^\pi(s_1) \leq \epsilon$$

offline reinforcement learning



Under what conditions on  $d^\rho$  we can learn a near-optimal policy?

# Dataset Coverage and Results

## Single Policy Coverage Assumption

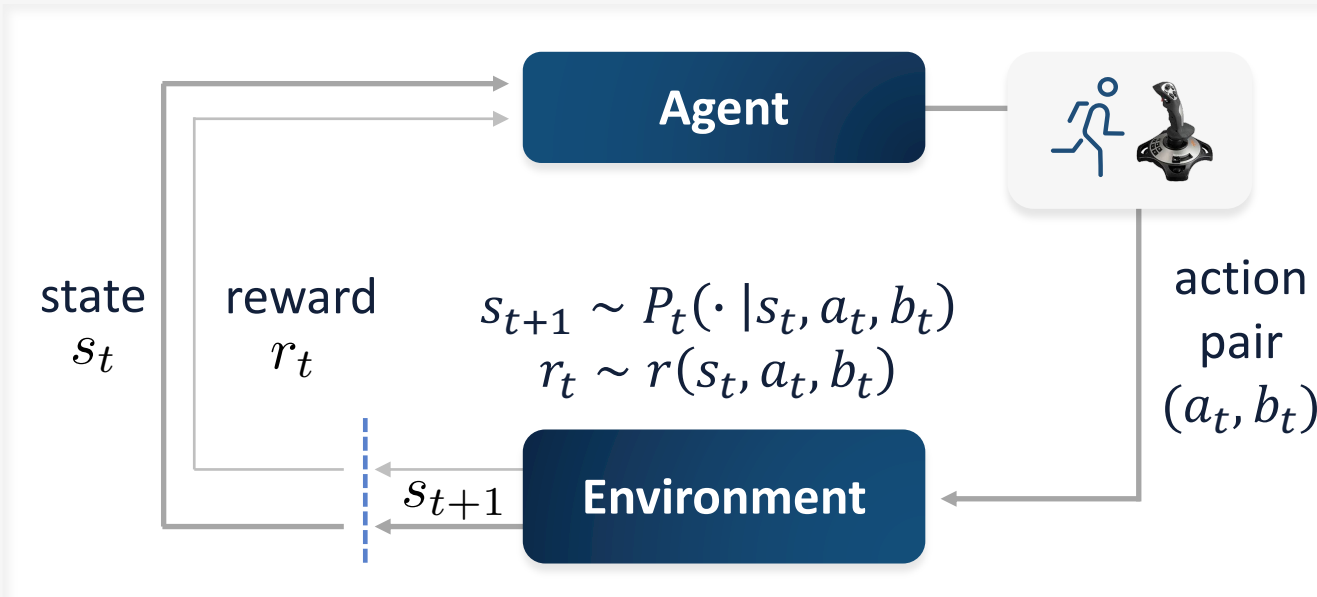
Necessary and Sufficient

- The behavior policy **only covers a single optimal policy**.
- There exists some constant  $\mathbf{C}_{\text{single}}$  such that  $\frac{d_{\mathbf{h}}^{\pi^*}(s,a)}{d_{\mathbf{h}}^{\rho}(s,a)} \leq \mathbf{C}_{\text{single}}$  for every  $(s, a)$  [LSAB19, JYW20].
- $1 \leq \mathbf{C}_{\text{single}} \leq \infty$
- Algorithmic idea: **Pessimism**. Penalize uncertain policies [JYW20, RZMJR21]. More later.
- Near-optimal bounds:  $\tilde{\Theta}\left(\frac{SH^3 \mathbf{C}_{\text{single}}}{\epsilon^2}\right)$  [XJWXB21].



# Two-Player Zero-Sum Markov Games

## Zero-Sum Markov Games



Repeat  $H$  times,  $H$ : planning horizon

Max player  $(a_1, a_2, \dots, a_H)$ :  $\max \mathbb{E}[r_1 + \dots + r_H]$

Min player  $(b_1, b_2, \dots, b_H)$ :  $\min \mathbb{E}[r_1 + \dots + r_H]$

## Zero-Sum Bandits



Special case of Markov games with  $H = 1$  and a fixed state.

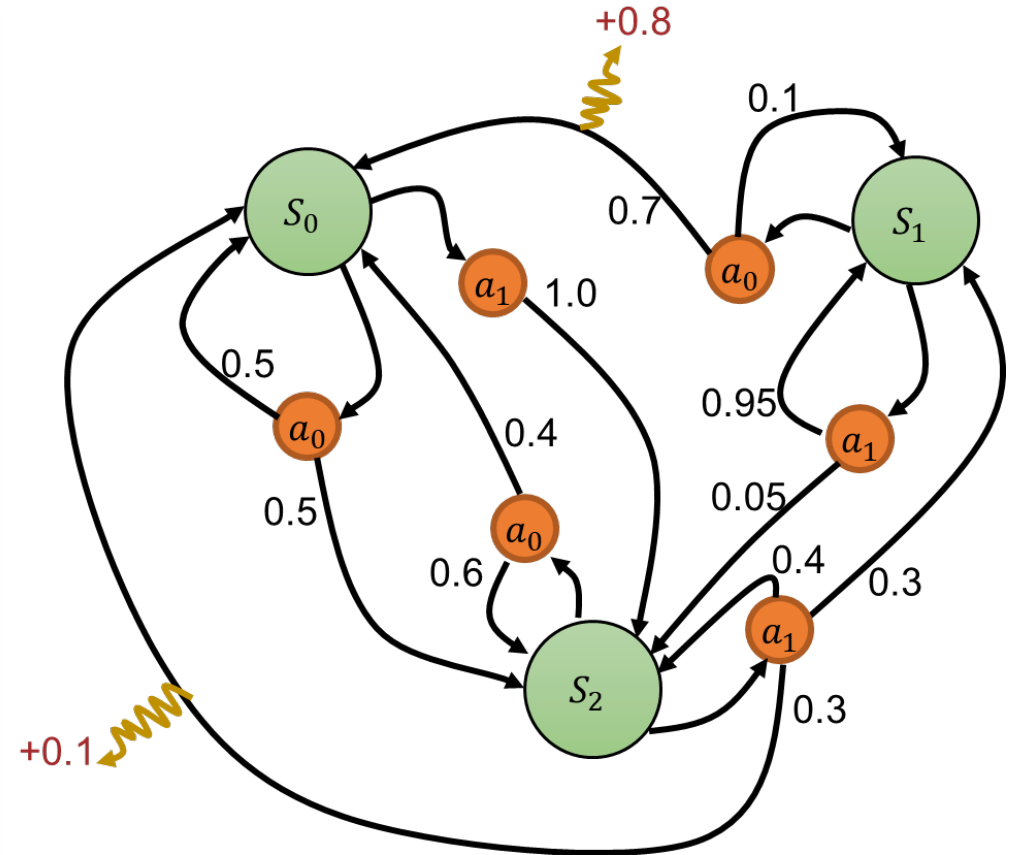
Only reward  $r(a, b)$  matters.

# Tabular Two-Player Zero-Sum Markov Games

## Assumptions:

1. # of States  $S < \infty$
2. Max player # of actions  $A < \infty$
3. Min player # of actions  $B < \infty$
4. Bounded rewards:  
 $0 \leq r_h \leq 1, h = 1, \dots, H$

Sample complexity depends on  
 $(S, A, B, H, 1/\epsilon)$



# Value Function, Best Response and Duality Gap

- **Policy pair:  $(\mu, \nu)$**

Max player policy  $\mu$  and min player policy  $\nu$ .  $\mu: \mathcal{S} \rightarrow \Delta(\mathcal{A}), \nu: \mathcal{S} \rightarrow \Delta(\mathcal{B})$ .

- **Q-function and Value Function:**

$$Q_h^{\mu, \nu}(s, a, b) = \mathbb{E}[r_h + r_{h+1} + \dots + r_H \mid s_h = s, a_h = a, b_h = b, \mu, \nu]$$
$$V_h^{\mu, \nu}(s) = \mathbb{E}[r_h + r_{h+1} + \dots + r_H \mid s_h = s, \mu, \nu]$$

- **Best response value for Max-player:** Given  $\mu$ ,  $V_h^{\mu, *}(s_h) = \min_{\nu} V_h^{\mu, \nu}(s_h)$
- **Best response value for Min-player:** Given  $\nu$ ,  $V_h^{*, \nu}(s_h) = \max_{\mu} V_h^{\mu, \nu}(s_h)$
- **Nash Equilibrium  $(\mu^*, \nu^*)$ :**  $V_h^{\mu^*, \nu^*}(s_h) = V_h^{\mu^*, *}(s_h) = V_h^{*, \nu^*}(s_h)$  [Shapley, 53].
- **Duality gap:**  $\text{Gap}(\mu, \nu) = V_1^{*, \nu}(s_1) - V_1^{\mu, *}(s_1)$

**Goal:** find  $(\mu, \nu)$  such that  $\text{Gap}(\mu, \nu) \leq \epsilon$

# Offline Two-Player Zero-Sum Markov Game

Offline Data:  $n$  (state, action, reward, next state) tuples:

$$D = \left\{ (s_h^i, a_h^i, b_h^i, r_h^i, s_{h+1}^i) \right\}_{h \in [H]}^{i \in [n]} \stackrel{i.i.d.}{\sim} d^\rho$$

- $\rho$ : data-collection /behavior policy pair
- $d_h^\rho(s, a, b)$  is the state-action distribution induced by  $\rho$  and transition  $P$ .
- Goal: learn a policy pair  $(\mu, \nu)$  from  $D$ :

$$\text{Gap}(\mu, \nu) \leq \epsilon$$



Under what conditions on  $d^\rho$  we can learn a near Nash Equilibrium?

What about single policy-pair coverage?

$$\frac{d_h^{(\mu^*, \nu^*)}(s, a, b)}{d_h^\rho(s, a, b)} \leq \mathbf{C}_{\text{single}}$$



NO

# Counter Example for Single Strategy Coverage

Min Player

Max Player

	$b_1$	$b_2$
$a_1$	0.5	1
$a_2$	0	0.5

Game 1

Min Player

Max Player

	$b_1$	$b_2$
$a_1$	0.5	0
$a_2$	1	0.5

Game 2

- NE for Game 1:  $(a_1, b_1)$ , NE for Game 2:  $(a_2, b_2)$
- Covers  $(a_1, b_1)$  and  $(a_2, b_2)$  with  $d^\rho(a_1, b_1) = d^\rho(a_2, b_2) = 0.5 \Rightarrow C_{\text{single}} = 2$ .
- We cannot differentiate Game 1 or Game 2!

Need to cover  $(a_1, b_2)$ ,  $(a_2, b_1)$

# Unilateral Coverage Assumption

Min Player

Max Player

	$b_1$	$b_2$	...	$b_B$
$a_1$	0.5	1	...	0.7
$a_2$	0	0.5	...	0.6
...	...	...	...	...
$a_A$	0.2	0.3	...	...

Nash Equilibrium:  $(a_1, b_1)$

- For a Nash Equilibrium  $(\mu^*, \nu^*)$ , the behavior policy covers  $(\mu^*, \nu)$  and  $(\mu, \nu^*)$  for all  $\mu$  and  $\nu$ .
- There exists some constant  $C_{\text{unilateral}}$  such that  $\frac{d_h^{\mu^*, \nu}(s, a, b)}{d_h^{\rho}(s, a, b)}, \frac{d_h^{\mu, \nu^*}(s, a, b)}{d_h^{\rho}(s, a, b)} \leq C_{\text{unilateral}}$  for every  $(s, a, b)$  and  $(\mu, \nu)$ .
- $A + B \leq C_{\text{unilateral}} \leq \infty$

Covered or not doesn't matter.

# A Weaker Assumption Than Unilateral Coverage?

Min Player

		$b_1$	$b_2$
Max Player	$a_1$	0.25	0.5
	$a_2$	0	0.75

Min Player

		$b_1$	$b_2$
Max Player	$a_1$	0.25	0.5
	$a_2$	1	0.75

Not Sufficient!

Game 1

$(a_2, b_1)$  not covered

Game 2

- A slightly weaker assumption: there exists **at most one** deterministic  $\mu$  or  $\nu$  such that the behavior policy  $\rho$  **does not cover**  $(\mu^*, \nu)$  or  $(\mu, \nu^*)$ .
- We cannot differentiate Game 1 or Game 2 without information of  $(a_2, b_1)$ .

# Algorithm for Two-Player Zero-Sum Bandits

Min Player

Max Player

	$b_1$	$b_2$	...	$b_B$
$a_1$	[0.4, 0.6]	[0.8, 1]	...	[0.7, 0.8]
$a_2$	[0, 0.1]	[0.4, 0.7]	...	[0.6, 0.7]
...	...	...	...	...
$a_A$	[0.1, 0.3]	[0.2, 0.4]	...	...

- Estimate  $r(a, b) \in [\underline{r}(a, b), \bar{r}(s, a)] \forall (a, b)$ .
- Compute NE  $(\underline{\mu}, \underline{\nu})$  for  $\underline{r}(\cdot, \cdot)$ .
- Compute NE  $(\bar{\mu}, \bar{\nu})$  for  $\bar{r}(\cdot, \cdot)$ .
- Output  $(\underline{\mu}, \bar{\nu})$ .

Pessimism




# Result for Two-Player Zero-Sum Bandits

## Theorem

- Sample complexity with unilateral coverage:  $\tilde{O}\left(\frac{ABC_{\text{unilateral}}}{\epsilon^2}\right)$
- Sample complexity with uniform coverage:  $\tilde{O}\left(\frac{C_{\text{unif}}}{\epsilon^2}\right)$
- Sample complexity for turn-based game with unilateral coverage:  $\tilde{O}\left(\frac{C_{\text{unilateral}}}{\epsilon^2}\right)$

- **Unilateral assumption is sufficient.**

- Lower bounds (from single-agent bandits)

- Sample complexity with unilateral coverage:  $\Omega\left(\frac{C_{\text{unilateral}}}{\epsilon^2}\right)$
  - Sample complexity with uniform coverage:  $\Omega\left(\frac{C_{\text{unif}}}{\epsilon^2}\right)$
  - Sample complexity for turn-based game with unilateral coverage:  $\Omega\left(\frac{C_{\text{unilateral}}}{\epsilon^2}\right)$
- Match
- 

# Algorithm for Markov Games

Min Player

Max Player

	$b_1$	$b_2$	...	$b_B$
$a_1$	[0.4, 0.6]	[0.8, 1]	...	[0.7, 0.8]
$a_2$	[0, 0.1]	[0.4, 0.7]	...	[0.6, 0.7]
...	...	...	...	...
$a_A$	[0.1, 0.3]	[0.2, 0.4]	...	...

Confidence for one state  $s$  at one step  $h$

- Estimate transition and reward using the dataset:  $\widehat{P}_h(s'|s, a, b), \widehat{r}(s, a, b)$
- Set  $\underline{V}_{H+1}(s) = \overline{V}_{H+1}(s) \leftarrow 0, \forall s.$
- For  $h=H, H-1, \dots, 1$ :
  - $\underline{Q}_h(s, a, b) \leftarrow \widehat{r}(s, a, b) + \langle \widehat{P}_h(\cdot | s, a, b), \underline{V}_{h+1}(\cdot) \rangle - \text{bonus}_h(s, a, b)$
  - Compute NE  $(\underline{\mu}_h, \underline{\nu}_h)$  for  $\underline{Q}_h(\cdot, \cdot, \cdot).$
  - $\underline{V}_h(s) \leftarrow \mathbb{E}_{(a,b) \sim (\underline{\mu}_h, \underline{\nu}_h)} [\underline{Q}_h(s, a, b)]$
  - Similarly get  $\overline{Q}_h$  with  $+\text{bonus}_h, \overline{V}_h, (\overline{\mu}_h, \overline{\nu}_h)$
- Output  $(\underline{\mu}, \overline{\nu}).$

DP Step




# Result for Two-Player Zero-Sum Markov Games

## Theorem

If the bonus is constructed using a **reference function** and **Bernstein bound**:

- with unilateral coverage:  $\tilde{O}\left(\frac{SABH^3 C_{\text{unilateral}}}{\epsilon^2}\right)$
- with uniform coverage:  $\tilde{O}\left(\frac{SH^3 C_{\text{unif}}}{\epsilon^2}\right)$
- for turn-based game with unilateral coverage:  $\tilde{O}\left(\frac{SH^3 C_{\text{unilateral}}}{\epsilon^2}\right)$

- **Unilateral assumption is sufficient for Markov games.**
- Lower bounds (from single-agent RL)

- with unilateral coverage:  $\Omega\left(\frac{SH^3 C_{\text{unilateral}}}{\epsilon^2}\right)$
  - with uniform coverage:  $\Omega\left(\frac{SH^3 C_{\text{unif}}}{\epsilon^2}\right)$
  - for turn-based game with unilateral coverage:  $\Omega\left(\frac{SH^3 C_{\text{unilateral}}}{\epsilon^2}\right)$
- Match**
- 

# Summary and Open Problems

## First theoretical study on two-player zero-sum Markov games

- Single-policy coverage not sufficient: **separation** between single-agent and two-player
- **Unilateral coverage**: sufficient and cannot be weakened.
- Algorithms based on pessimism for both players
  - Polynomial bound for unilateral coverage.
  - Near-optimal bounds for (1) uniform coverage, (2) unilateral coverage + turn-based games.
- Concurrent work also studied linear MDP [ZXTWZWY22].

## Future Directions

- Improve bound under unilateral coverage (now  $AB$  factor gap).
- General sum in multi-agent games (online setting [ZMB21, JLWY21, ...]).

Upcoming Work!



Thank You

# Analysis

- Confidence interval length:  $\text{bonus}(a, b) \approx \sqrt{\frac{1}{n(a,b)}} \approx \sqrt{\frac{1}{nd^\rho(a,b)}}$
- $r(\underline{\mu}^*, \underline{\nu}^*) \leq r(\underline{\mu}^*, \underline{\nu})$  (by the defn of  $\underline{\nu}^*$ )
- $r(\underline{\mu}, *) \geq \underline{r}(\underline{\mu}, *) \geq \underline{r}(\underline{\mu}, \underline{\nu}) \geq \underline{r}(\underline{\mu}^*, \underline{\nu})$  (by the defns of  $\underline{r}$  and  $\underline{\nu}$ )
- $r(\underline{\mu}^*, \underline{\nu}^*) - r(\underline{\mu}, *) \leq r(\underline{\mu}^*, \underline{\nu}) - \underline{r}(\underline{\mu}^*, \underline{\nu}) \leq \mathbb{E}_{(a,b) \sim (\underline{\mu}^*, \underline{\nu})}[\text{bonus}(a, b)]$
- Similarly,  $r(*, \bar{\nu}) - r(\underline{\mu}^*, \underline{\nu}^*) \leq \mathbb{E}_{(a,b) \sim (\bar{\mu}, *)}[\text{bonus}(a, b)]$
- $\text{Gap}(\underline{\mu}, \bar{\nu}) \leq \mathbb{E}_{(a,b) \sim (\underline{\mu}^*, \underline{\nu})}[\text{bonus}(a, b)] + \mathbb{E}_{(a,b) \sim (\bar{\mu}, *)}[\text{bonus}(a, b)]$
- Then use Cauchy-Schwartz