

# Learning in Zero-sum games

## Reinforcement Learning Seminar

---

Yingru Li

May 12, 2019

[yingruli@link.cuhk.edu.cn](mailto:yingruli@link.cuhk.edu.cn)

The Chinese University of Hong Kong, Shenzhen



## Motivation: a Long-Standing Goal of AI...



Figure 1: Deep Blue

## Motivation: a Long-Standing Goal of AI...



Figure 2: AlphaGo

# Motivation: a Long-Standing Goal of AI...



Figure 3: Libratus

# Motivation: a Long-Standing Goal of AI...



**Figure 4:** StarCraft II: A New Challenge for Reinforcement Learning. DeepMind AlphaStar Jan. 2019; Tencent AI Lab TStarBots Sep. 2018

## ...with Potential Applications in Real-World Environments

- Security
- Negotiation
- Diplomatic and Military Strategy
- Financial Market
- E-Commerce
- Distributed Cooperated and Competitive Robotics
- Game AI
- .....



## Learning in Two-Player Zero-Sum Games

Regret Minimization and Nash Equilibrium

The Exp3 Algorithms

## From Normal Form to Extensive Form Imperfect Information Games

Regret Minimization and Nash Equilibria

Counterfactual Regret Minimization



## Learning in Two-Player Zero-Sum Games

Regret Minimization and Nash Equilibrium

The Exp3 Algorithms

## From Normal Form to Extensive Form Imperfect Information Games

Regret Minimization and Nash Equilibria

Counterfactual Regret Minimization





## *The Game*

- Set of players  $N = \{1, \dots, n\}$
- Action sets  $A_i$ , **joint** action set  $A = A_1 \times \dots \times A_n$
- Joint action  $a \in A$ , player  $i$ 's action  $a_i$ , all other players'  $a_{-i}$
- **Utility (payoff/reward)** function  $u : A \rightarrow \mathbb{R}^n$ ,
- Player  $i$ 's utility  $u_i : A \rightarrow \mathbb{R}$

## *Mixed strategies*

- Joint strategy  $\sigma \in \mathcal{D}(A)$  is distribution over  $A$ , such that

$$\sigma(a) = \prod_{i=1}^n \sigma_i(a_i)$$

- Utility of a strategy for player  $i$  (**expected utility**):

$$u_i(\sigma) = \sum_{a_i} \sum_{a_{-i}} \sigma_i(a_i) \sigma_{-i}(a_{-i}) u_i(a_i, a_{-i})$$

## *The Game*

- Best response:

$$\sigma_i^* \in BR(\sigma_{-i}) \text{ iff } \forall \sigma_i \in \mathcal{D}(A_i), u_i(\sigma_i^*, \sigma_{-i}) \geq u_i(\sigma_i, \sigma_{-i})$$

- Nash equilibrium:  $\sigma$  is a Nash equilibrium iff  $\forall i, \sigma_i \in BR(\sigma_{-i})$
- Every finite game has a Nash equilibrium! [Nash, 1950]

# Finite Two-Player Zero-Sum Games

## *The Game*

- Set of players  $N = \{1, 2\} = \{i, j\}$
- Action sets  $A_i$ , joint action set  $A = A_1 \times A_2$
- Joint action  $a \in A$ , player  $i$ 's action  $a_i$ , all other players'  $a_j$
- Utility (payoff/reward) function  $u : A \rightarrow \mathbb{R}^n$ , player  $i$ 's utility  $u_i : A \rightarrow \mathbb{R}$

$$\forall a \in A, \quad u_1(a) = -u_2(a)$$

## *Mixed strategies*

- Nash equilibrium [Minimax theorem (von Neumann, 1928)]

$$\begin{aligned}(\sigma_1^*, \sigma_2^*) &= \arg \max_{\sigma_1} \min_{\sigma_2} u_1(\sigma_1, \sigma_2) \\ &= \arg \min_{\sigma_1} \max_{\sigma_2} u_2(\sigma_1, \sigma_2)\end{aligned}$$

- Value of the game

$$V = \max_{\sigma_1} \min_{\sigma_2} u_1(\sigma_1, \sigma_2) = \min_{\sigma_2} \max_{\sigma_1} u_1(\sigma_1, \sigma_2)$$

# Rock-Paper-Scissors The Game

Action set  $A_1 = A_2 = \{(R)ock, (P)aper, (S)cissor\}$

	<i>R</i>	<i>P</i>	<i>S</i>
<i>R</i>	0, 0	-1, 1	1, -1
<i>P</i>	1, -1	0, 0	-1, 1
<i>S</i>	-1, 1	1, -1	0, 0



# Rock-Paper-Scissors The Solution

Action set  $A_1 = A_2 = \{(R)ock, (P)aper, (S)cissor\}$

	<i>R</i>	<i>P</i>	<i>S</i>
<i>R</i>	0, 0	-1, 1	1, -1
<i>P</i>	1, -1	0, 0	-1, 1
<i>S</i>	-1, 1	1, -1	0, 0

- if  $(\sigma_1^*, \sigma_2^*)$  is a Nash equilibrium, then

$$\begin{aligned}\sigma_1^* &= \text{BR}(\sigma_2^*) = \arg \max_{\sigma_1} u_1(\sigma_1, \sigma_2^*) \\ &= \arg \max_{\sigma_1} \sum_{a_1 \in A_1} \sigma_1(a_1) u_1(a_1, \sigma_2^*)\end{aligned}$$



# Rock-Paper-Scissors The Solution

Action set  $A_1 = A_2 = \{(R)ock, (P)aper, (S)cissor\}$

	<i>R</i>	<i>P</i>	<i>S</i>
<i>R</i>	0, 0	-1, 1	1, -1
<i>P</i>	1, -1	0, 0	-1, 1
<i>S</i>	-1, 1	1, -1	0, 0

- if  $(\sigma_1^*, \sigma_2^*)$  is a Nash equilibrium, then

$$\begin{aligned}\sigma_1^* &= \text{BR}(\sigma_2^*) = \arg \max_{\sigma_1} u_1(\sigma_1, \sigma_2^*) \\ &= \arg \max_{\sigma_1} \sum_{a_1 \in A_1} \sigma_1(a_1) u_1(a_1, \sigma_2^*)\end{aligned}$$

$$\Rightarrow \forall a_1 \in A, \quad u_1 = u_1(a_1, \sigma_2^*)$$

# Rock-Paper-Scissors The Solution (sketch)

	<i>R</i>	<i>P</i>	<i>S</i>
<i>R</i>	0, 0	-1, 1	1, -1
<i>P</i>	1, -1	0, 0	-1, 1
<i>S</i>	-1, 1	1, -1	0, 0

- Let  $\sigma_2 = (\sigma_2(R), \sigma_2(P), \sigma_2(S))$  the strategy of player **column**,

$$u_1 = u_1(R, \sigma_2) = 0\sigma_2(R) - 1\sigma_2(P) + 1\sigma_2(S)$$

$$u_1 = u_1(P, \sigma_2) = 1\sigma_2(R) + 0\sigma_2(P) - 1\sigma_2(S)$$

$$u_1 = u_1(S, \sigma_2) = -1\sigma_2(R) + 1\sigma_2(P) + 0\sigma_2(S)$$

$$1 = \sigma_2(R) + \sigma_2(P) + \sigma_2(S)$$



# Rock-Paper-Scissors The Solution (sketch)

	<i>R</i>	<i>P</i>	<i>S</i>
<i>R</i>	0, 0	-1, 1	1, -1
<i>P</i>	1, -1	0, 0	-1, 1
<i>S</i>	-1, 1	1, -1	0, 0

- Let  $\sigma_2 = (\sigma_2(R), \sigma_2(P), \sigma_2(S))$  the strategy of player **column**,

$$u_1 = u_1(R, \sigma_2) = 0\sigma_2(R) - 1\sigma_2(P) + 1\sigma_2(S)$$

$$u_1 = u_1(P, \sigma_2) = 1\sigma_2(R) + 0\sigma_2(P) - 1\sigma_2(S)$$

$$u_1 = u_1(S, \sigma_2) = -1\sigma_2(R) + 1\sigma_2(P) + 0\sigma_2(S)$$

$$1 = \sigma_2(R) + \sigma_2(P) + \sigma_2(S)$$

- Solving for all variables gives  $\sigma_2^* = (1/3, 1/3, 1/3)$  and  $u_1 = 0$





# Rock-Paper-Scissors The Solution (sketch)

	<i>R</i>	<i>P</i>	<i>S</i>
<i>R</i>	0, 0	-1, 1	1, -1
<i>P</i>	1, -1	0, 0	-1, 1
<i>S</i>	-1, 1	1, -1	0, 0

- Let  $\sigma_2 = (\sigma_2(R), \sigma_2(P), \sigma_2(S))$  the strategy of player column,

$$u_1 = u_1(R, \sigma_2) = 0\sigma_2(R) - 1\sigma_2(P) + 1\sigma_2(S)$$

$$u_1 = u_1(P, \sigma_2) = 1\sigma_2(R) + 0\sigma_2(P) - 1\sigma_2(S)$$

$$u_1 = u_1(S, \sigma_2) = -1\sigma_2(R) + 1\sigma_2(P) + 0\sigma_2(S)$$

$$1 = \sigma_2(R) + \sigma_2(P) + \sigma_2(S)$$

- Solving for all variables gives  $\sigma_2^* = (1/3, 1/3, 1/3)$  and  $u_1 = 0$
- Repeating for player row gives  $\sigma_1^* = (1/3, 1/3, 1/3)$  and  $u_2 = 0$
- $(\sigma_1^*, \sigma_2^*)$  is a Nash equilibrium and the value of the game is  $V = 0$



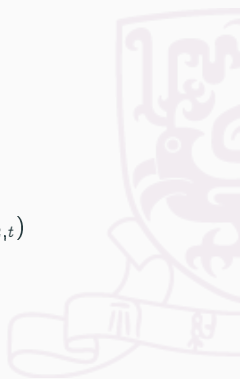
# A Single-Player Perspective

## *Sequential game*

- For  $t = 1, \dots, n$ 
  - Player 1 chooses  $\sigma_{1,t}$
  - Player 2 chooses  $\sigma_{2,t}$
  - Players play actions  $a_{1,t} \sim \sigma_{1,t}$  and  $a_{2,t} \sim \sigma_{2,t}$
  - Players receive payoffs  $u_1(a_{1,t}, a_{2,t})$  and  $u_2(a_{1,t}, a_{2,t})$

*Solution:* Nash equilibrium

$$(\sigma_1^*, \sigma_2^*) = \arg \max_{\sigma_1} \min_{\sigma_2} u_1(\sigma_1, \sigma_2)$$



# A Single-Player Perspective

*Sequential game*  $\Rightarrow$  *Single-player game*

- For  $t = 1, \dots, n$ 
  - Player 1 chooses  $\sigma_{1,t}$
  - ~~Player 2 chooses  $\sigma_{2,t}$~~
  - Players play actions  $a_{1,t} \sim \sigma_{1,t}$  and  ~~$a_{2,t} \sim \sigma_{2,t}$~~
  - Players receive payoffs  $u_1(a_{1,t}, a_{2,t})$  and  ~~$u_2(a_{1,t}, a_{2,t})$~~

*Solution: Nash equilibrium*  $\Rightarrow$  *Maximize the (average) utility*

$$(\sigma_1^*, \sigma_2^*) = \arg \max_{\sigma_1} \min_{\sigma_2} u_1(\sigma_1, \sigma_2)$$

$$\begin{aligned}(a_{1,1}^*, \dots, a_{1,n}^*) &= \arg \max_{(a_{1,1}, \dots, a_{1,n})} \frac{1}{n} \sum_{t=1}^n u_1(a_{1,t}, a_{2,t}) \\ &= \arg \max_{(a_{1,1}, \dots, a_{1,n})} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t})\end{aligned}$$

# The (Multi-Armed Bandit) Problem

*A learning problem*

- For  $t = 1, \dots, n$ 
  - Player 1 chooses  $\sigma_{1,t}$
  - Player 1 plays action  $a_{1,t} \sim \sigma_{1,t}$
  - Player 1 receives payoff  $u_{1,t}(a_{1,t})$

*Remarks*

- No information about  $a_{2,t}$  and utility  $u_2$
- Utility function  $u_{1,t}$  is only observed for  $a_{1,t}$  (i.e., bandit feedback  $u_{1,t}(a_{1,t})$ )



# The (Multi-Armed Bandit) Problem

- **Regret in hindsight** w.r.t. any fixed action  $a_1$

$$R_n(a_1) = \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t})$$

- Objective: find actions  $(a_{1,1}, \dots, a_{1,n})$  that maximize average utility  $\approx$  **minimize the regret** w.r.t. the **best action  $a_1$  in hindsight**

$$\text{Utility: } \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t})$$

$$\text{Regret: } R_n = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t})$$

# Regret Minimization and Nash Equilibrium

## Theorem

A learning algorithm is **Hannan's consistent** if

$$\limsup_{n \rightarrow \infty} R_n = 0 \quad a.s.$$

Given a two-player zero-sum game with **value**  $V$ , if players choose strategies  $\sigma_{1,t}$  and  $\sigma_{2,t}$  using a Hannan's consistent algorithm, then

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_1(a_{1,t}, a_{2,t}) = V$$

Furthermore, let empirical frequency strategies be

$$\hat{\sigma}_{1,n}(a_1) = \frac{1}{n} \sum_{t=1}^n \mathbb{I}\{a_{1,t} = a_1\} \quad \text{and} \quad \hat{\sigma}_{2,n}(a_2) = \frac{1}{n} \sum_{t=1}^n \mathbb{I}\{a_{2,t} = a_2\}$$

then the joint empirical strategy

$$\hat{\sigma}_{1,n} \times \hat{\sigma}_{2,n} \xrightarrow{n \rightarrow \infty} \{(\sigma_1^*, \sigma_2^*)\}_{\text{Nash}}$$

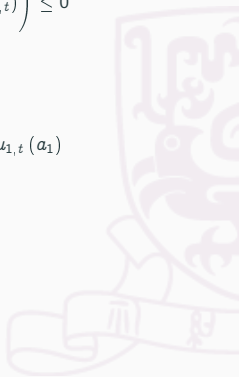
# Regret Minimization and Nash Equilibria [proof]

- Hannan's consistency

$$\limsup_{n \rightarrow \infty} R_n \leq 0 \iff \limsup_{n \rightarrow \infty} \left( \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}) \right) \leq 0$$

- linearity of utility function

$$\max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(\sigma_1) = \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n \sum_{a_1 \in A_1} \sigma_1(a_1) u_{1,t}(a_1) = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1)$$



# Regret Minimization and Nash Equilibria [proof]

- Hannan's consistency

$$\limsup_{n \rightarrow \infty} R_n \leq 0 \iff \limsup_{n \rightarrow \infty} \left( \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}) \right) \leq 0$$

- linearity of utility function

$$\max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(\sigma_1) = \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n \sum_{a_1 \in A_1} \sigma_1(a_1) u_{1,t}(a_1) = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1)$$

- definition  $u_{1,t}(\sigma_1) = u_1(\sigma_1, a_{2,t}) \Rightarrow$

$$\frac{1}{n} \sum_{t=1}^n u_{1,t}(\sigma_1) = \frac{1}{n} \sum_{t=1}^n \sum_{a_2 \in A_2} \mathbb{I}\{a_{2,t} = a_2\} u_1(\sigma_1, a_2) = \sum_{a_2 \in A_2} u_1(\sigma_1, a_2) \underbrace{\frac{1}{n} \sum_{t=1}^n \mathbb{I}\{a_{2,t} = a_2\}}_{\widehat{\sigma}_{2,n}(a_2)}$$



# Regret Minimization and Nash Equilibria [proof]

- Hannan's consistency

$$\limsup_{n \rightarrow \infty} R_n \leq 0 \iff \limsup_{n \rightarrow \infty} \left( \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}) \right) \leq 0$$

- linearity of utility function

$$\max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(\sigma_1) = \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n \sum_{a_1 \in A_1} \sigma_1(a_1) u_{1,t}(a_1) = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1)$$

- definition  $u_{1,t}(\sigma_1) = u_1(\sigma_1, a_{2,t}) \Rightarrow$

$$\frac{1}{n} \sum_{t=1}^n u_{1,t}(\sigma_1) = \frac{1}{n} \sum_{t=1}^n \sum_{a_2 \in A_2} \mathbb{I}\{a_{2,t} = a_2\} u_1(\sigma_1, a_2) = \sum_{a_2 \in A_2} u_1(\sigma_1, a_2) \underbrace{\frac{1}{n} \sum_{t=1}^n \mathbb{I}\{a_{2,t} = a_2\}}_{\widehat{\sigma}_{2,n}(a_2)}$$

- one-side of the result

$$\max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(\sigma_1) = \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \widehat{\sigma}_{2,n}) \geq \max_{\sigma_1} \min_{\sigma_2} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_2) = V$$

- one-side of the result

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}, a_{2,t}) \geq \max_{\sigma_1} \min_{\sigma_2} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_2) = V$$



- one-side of the result

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}, a_{2,t}) \geq \max_{\sigma_1} \min_{\sigma_2} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_2) = V$$

- If player 2 also plays Hannan consistent strategies, then we get

$$\max_{\sigma_2} \frac{1}{n} \sum_{t=1}^n u_{2,t}(\sigma_2) \geq \max_{\sigma_2} \min_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_2(\sigma_1, \sigma_2) = V$$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}, a_{2,t}) \leq \min_{\sigma_2} \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_2) = V$$

- 

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}, a_{2,t}) = V \quad a.s.$$

## Remark

The joint empirical strategy converges to the set of correlated equilibrium almost surely as  $n \rightarrow \infty$ .

In particular, for any (finite) two-person zero-sum game, for each player, the empirical distribution of play converges to the set of optimal mixed actions.

$$\hat{\sigma}_{1,n} \times \hat{\sigma}_{2,n} \xrightarrow{n \rightarrow \infty} \{(\sigma_1^*, \sigma_2^*)\}_{\text{Nash}} \quad a.s.$$

Note that approaching to a set does not imply convergence to particular point.

## Corollary

If

$$R_n \leq \epsilon$$

then the joint empirical strategy is  $\epsilon$ -Nash (more precisely, correlated  $\epsilon$ -equilibrium), i.e.,

$$u_1(\hat{\sigma}_{1,n} \times \hat{\sigma}_{2,n}) \geq V - \epsilon$$

*A learning problem*

- For  $t = 1, \dots, n$ 
  - Player 1 chooses  $\sigma_{1,t}$
  - Player 1 plays action  $a_{1,t} \sim \sigma_{1,t}$
  - Player 1 receives payoff  $u_{1,t}(a_{1,t})$

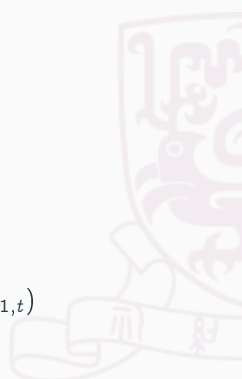
*Objective*

- Regret

$$R_n = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t})$$

- Hannan's consistent algorithm

$$\limsup_{n \rightarrow \infty} R_n \leq 0 \quad a.s.$$



# Learning the Nash Equilibrium

*Version 1:* fictitious play full information (aka follow-the-leader)

- For  $t = 1, \dots, n$ 
  - Compute greedy action

$$a_t^* = \arg \max_{a \in A_1} \sum_{s=1}^{t-1} u_{1,t}(a)$$

- Player chooses  $\sigma_{1,t} = \delta(a_t^*)$
- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$

*Remarks*

- This strategy is easily exploitable  $R_n = O(1)$
- E.g. Opponents set  $u_{1,t}(a = a_{1,t}) = -1$  and  $u_{1,t}(a \neq a_{1,t}) = 1$



# Learning the Nash Equilibrium

*Version 1: fictitious play full information (aka follow-the-leader)*

- For  $t = 1, \dots, n$ 
  - Compute greedy action

$$a_t^* = \arg \max_{a \in A_1} \sum_{s=1}^{t-1} u_{1,t}(a)$$

- Player chooses  $\sigma_{1,t} = \delta(a_t^*)$
- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$

*Remarks*

- This strategy is easily exploitable  $R_n = O(1)$
- Self play **does not converge** in general [Recall Hannan's consistency]





# Learning the Nash Equilibrium

Version 2: [Randomization]



# Learning the Nash Equilibrium

Version 2: [Randomization] exponentially weighted forecaster (EWF)

- Initialize weights  $w_0(a) = 1$  for all  $a \in A_1$
- For  $t = 1, \dots, n$ 
  - Player chooses

$$\sigma_{1,t}(a) = \frac{w_{t-1}(a)}{\sum_{b \in A_1} w_{t-1}(b)} \quad [\text{prop. to weights}]$$

- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$  and  $u_{1,t}(a)$  for all  $a$  [full info]
- Update weights  $w_{1,t}(a_{1,t})$

$$w_t(a) = w_{t-1}(a) \exp(\eta_t u_{1,t}(a)) \quad [\text{exponentiated utility}]$$

# Learning the Nash Equilibrium

## Theorem

If EWF is run over  $n$  steps with  $\eta_t = \eta$ , then with probability  $1 - \delta$

$$R_n = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}) \leq \frac{\log(A_1)}{n\eta} + \frac{\eta}{8} + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$$

Setting  $\eta = \sqrt{8 \log(A_1) / n}$  we obtain

$$R_n \leq \sqrt{\frac{\log(A_1)}{2n}} + \sqrt{\frac{1}{2n} \log(1/\delta)}$$

## Remarks

- $\limsup_{n \rightarrow \infty} R_n \leq 0 \Rightarrow$  Hannan's consistency
- Rate of convergence  $O(1/\sqrt{n})$
- In self-play EWF converges to the Nash equilibrium

# Learning the Nash Equilibrium

Version 2: [Randomization] exponentially weighted forecaster (EWF)

- Initialize weights  $w_0(a) = 0$  for all  $a \in A_1$
- For  $t = 1, \dots, n$ 
  - Player chooses

$$\sigma_{1,t}(a) = \frac{w_{t-1}(a)}{\sum_{b \in A_1} w_{t-1}(b)} \quad \text{[prop. to weights]}$$

- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$  and  $u_{1,t}(a)$  for all  $a$  [full info]
- Update weights  $w_{1,t}(a_{1,t})$

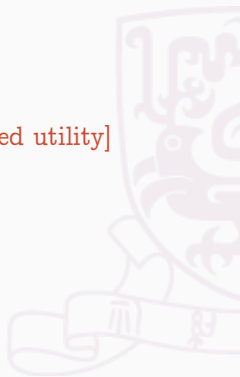
$$w_t(a) = w_{t-1}(a) \exp(\eta_t u_{1,t}(a)) \quad \text{[exponentiated utility]}$$

# Learning the Nash Equilibrium

*Problem:*

- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$
- Update weights  $w_{1,t}(a_{1,t})$

$$w_t(a) = w_{t-1}(a) \exp(\eta_t u_{1,t}(a)) \quad \text{[exponentiated utility]}$$



# Learning the Nash Equilibrium

*Problem:*

- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$
- Update weights  $u_{1,t}(a_{1,t})$

$$w_t(a) = w_{t-1}(a) \exp(\eta_t u_{1,t}(a)) \quad \text{[exponentiated utility]}$$

*Solution:*

- Importance sampling

$$\tilde{u}_{1,t}(a) = \begin{cases} \frac{u_{1,t}(a_{1,t})}{\sigma_{1,t}(a_{1,t})} & \text{if } a = a_{1,t} \\ 0 & \text{otherwise} \end{cases}$$

- Unbiased estimator

$$\forall a \in A_1 \quad \mathbb{E}_{a \sim \sigma_{1,t}} [\tilde{u}_{1,t}(a)] = \sigma_{1,t}(a) \frac{u_{1,t}(a)}{\sigma_{1,t}(a)} + (1 - \sigma_{1,t}(a)) \times 0 = u_{1,t}(a)$$

# Learning the Nash Equilibrium

## Version 3: EWF for Exploration-Exploitation (EXP3)

- Initialize weights  $w_0(a) = 0$  for all  $a \in A_1$
- For  $t = 1, \dots, n$ 
  - Player chooses

$$\sigma_{1,t}(a) = \frac{w_{t-1}(a)}{\sum_{b \in A_1} w_{t-1}(b)} \quad [\text{prop. to weights}]$$

- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$
- Compute pseudo-payoffs

$$\tilde{u}_{1,t}(a) = \begin{cases} \frac{u_{1,t}(a_{1,t})}{\sigma_{1,t}(a_{1,t})} & \text{if } a = a_{1,t} \\ 0 & \text{otherwise} \end{cases}$$

- Update weights  $w_{1,t}(a_{1,t})$

$$w_t(a) = w_{t-1}(a) \exp(\eta_t \tilde{u}_{1,t}(a))$$



## Theorem

If EXP3 is run over  $n$  steps with  $\eta_t = \sqrt{2 \log(A_1) / (nA_1)}$ , then its **pseudo-regret** is bounded as

$$\bar{R}_n = \max_{a_1} \frac{1}{n} \sum_{t=1}^n \mathbb{E} [u_{1,t}(a_1)] - \frac{1}{n} \sum_{t=1}^n \mathbb{E} [u_{1,t}(a_{1,t})] \leq \sqrt{\frac{2A_1 \log(A_1)}{n}}$$



# Learning the Nash Equilibrium

## Theorem

If EXP3 is run over  $n$  steps with  $\eta_t = \sqrt{2 \log(A_1) / (nA_1)}$ , then its **psuedo-regret** is bounded as

$$\bar{R}_n = \max_{a_1} \frac{1}{n} \sum_{t=1}^n \mathbb{E} [u_{1,t}(a_1)] - \frac{1}{n} \sum_{t=1}^n \mathbb{E} [u_{1,t}(a_{1,t})] \leq \sqrt{\frac{2A_1 \log(A_1)}{n}}$$

## Remarks

- $\limsup_{n \rightarrow \infty} \bar{R}_n \leq 0 \Rightarrow$  Hannan's consistency?
- Rate of convergence  $O(1/\sqrt{n})$
- Regret larger by a factor  $\sqrt{A_1}$  (observing 1 vs  $A_1$  payoffs)

# Rock-Paper-Scissors– The Simulation

Action set  $A_1 = A_2 = \{(R)ock, (P)aper, (S)cissor\}$

	R	P	S
R	0, 0	-1, 1	5, -5
P	1, -1	0, 0	-1, 1
S	-1, 1	1, -1	0, 0

- Equilibrium  $\sigma_1^* = (1/7, 11/21, 1/3)$
- Value of the game  $V = 4/21 (\approx 0.1904)$



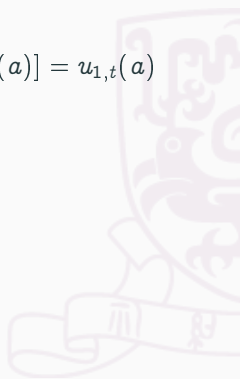
*Problem:*

- Importance sampling is unbiased

$$\tilde{u}_{1,t}(a) = \begin{cases} \frac{u_{1,t}(a_{1,t})}{\sigma_{1,t}(a_{1,t})} & \text{if } a = a_{1,t} \\ 0 & \text{otherwise} \end{cases} ; \quad \mathbb{E}_{a \sim \sigma_{1,t}} [\tilde{u}_{1,t}(a)] = u_{1,t}(a)$$

- Variance

$$\mathbb{V}_{a \sim \sigma_{1,t}} [\tilde{u}_{1,t}(a)] \xrightarrow{\sigma_{1,t}(a) \rightarrow 0} \infty$$



# Learning the Nash Equilibrium

*Problem:*

- Importance sampling is unbiased

$$\tilde{u}_{1,t}(a) = \begin{cases} \frac{u_{1,t}(a_{1,t})}{\sigma_{1,t}(a_{1,t})} & \text{if } a = a_{1,t} \\ 0 & \text{otherwise} \end{cases} ; \quad \mathbb{E}_{a \sim \sigma_{1,t}} [\tilde{u}_{1,t}(a)] = u_{1,t}(a)$$

- Variance

$$\mathbb{V}_{a \sim \sigma_{1,t}} [\tilde{u}_{1,t}(a)] \xrightarrow{\sigma_{1,t}(a) \rightarrow 0} \infty$$

*Solution:*

- **Bias** both pseudo-payoff

$$\tilde{u}_{1,t}(a) = \frac{u_{1,t}(a_{1,t}) \mathbb{I}\{a = a_{1,t}\} + \beta_t}{\sigma_{1,t}(a_{1,t})}$$

- Mix strategy with **uniform** exploration (**now bounded below**)

$$\sigma_{1,t}(a) = (1 - \gamma_t) \frac{w_{1,t}(a)}{\sum_{b \in A_1} w_{1,t}(b)} + \frac{\gamma_t}{A_1}$$

# Learning the Nash Equilibrium

Version 3: EWF for Exploration-Exploitation w.h.p. (EXP3.P)

- Initialize weights  $w_0(a) = 0$  for all  $a \in A_1$
- For  $t = 1, \dots, n$ 
  - Player chooses

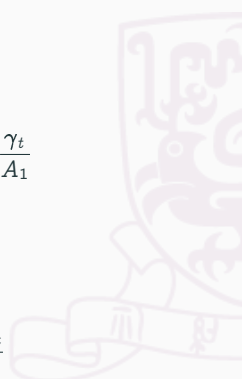
$$\sigma_{1,t}(a) = (1 - \gamma_t) \frac{w_{1,t}(a)}{\sum_{b \in A_1} w_{1,t}(b)} + \frac{\gamma_t}{|A_1|}$$

- Player plays action  $a_{1,t} \sim \sigma_{1,t}$
- Player receives payoff  $u_{1,t}(a_{1,t})$
- Compute pseudo-payoffs

$$\tilde{u}_{1,t}(a) = \frac{u_{1,t}(a_{1,t}) \mathbb{I}\{a = a_{1,t}\} + \beta_t}{\sigma_{1,t}(a_{1,t})}$$

- Update weights  $w_{1,t}(a_{1,t})$

$$w_t(a) = w_{t-1}(a) \exp(\eta_t \tilde{u}_{1,t}(a))$$



## Lemma

For  $\beta_t \leq 1$ , let

$$\tilde{u}_{1,t}(a) = \frac{u_{1,t}(a_{1,t}) \mathbb{I}\{a = a_{1,t}\} + \beta_t}{\sigma_{1,t}(a_{1,t})}$$

Then, w.p. at least  $1 - \delta$ ,

$$\sum_{t=1}^n u_{i,t}(a) \leq \sum_{t=1}^n \tilde{u}_{i,t}(a) + \frac{\log \delta^{-1}}{\beta_t}$$

# Learning the Nash Equilibrium

## Theorem

If EXP3 is run over  $n$  steps with  $\beta_t \approx \eta_t = \sqrt{2 \log(A_1) / (nA_1)}$ ,  $\gamma_t = \sqrt{A_1 \log(A_1) / n}$ , then with probability  $1 - \delta$  its regret is bounded as

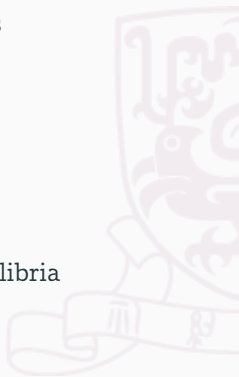
$$R_n = \max_{a_1} \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_1) - \frac{1}{n} \sum_{t=1}^n u_{1,t}(a_{1,t}) \leq 6 \sqrt{\frac{A_1 \log(A_1/\delta)}{n}}$$

## Remarks

- $\lim_{n \rightarrow \infty} R_n \leq 0 \Rightarrow$  Hannan's consistency!
- EXP3.P in self-play converges to Nash equilibrium

# Summary

- + EXP3.P minimizes regret in adversarial environments
  - + EXP3.P converges to Nash equilibria in self-play
  - + No need to know
    - Utility function (i.e., the rules of the game)
    - Actions performed by the adversary
- ≈ Some of this can be extended to learn correlated equilibria
- Exponential may be tricky to manage
  - Convergence is only in the empirical frequency
  - Convergence is relatively slow





## Learning in Two-Player Zero-Sum Games

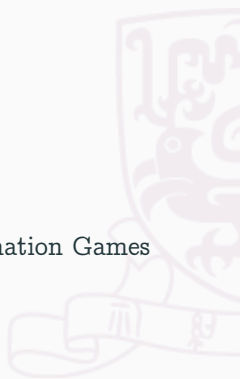
Regret Minimization and Nash Equilibrium

The Exp3 Algorithms

## From Normal Form to Extensive Form Imperfect Information Games

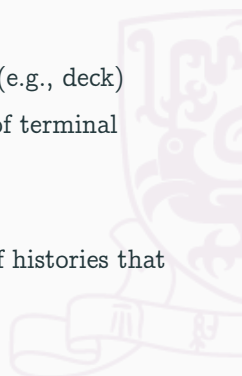
Regret Minimization and Nash Equilibria

Counterfactual Regret Minimization



## The game

- Set of players  $N = \{1, \dots, n\}$  and  $c$  chance player (e.g., deck)
- Set of possible sequences of actions  $H, Z \subseteq H$  set of terminal histories
- Player function  $P : H \rightarrow N \cup \{c\}$
- Set of information sets  $\mathcal{I} = \{I\}$  (i.e.,  $I$  is a subset of histories that are not distinguishable)
- Utility of a terminal history  $u_i : Z \rightarrow \mathbb{R}$
- Strategy  $\sigma_i : \mathcal{I} \rightarrow \mathcal{D}(A)$  (in all  $h \in I$  such that  $P(h) = i$ )



## Histories

- Prob. of reaching history  $h \in H$  following joint strategy  $\sigma$ ,  $\pi^\sigma(h)$
- Prob. of reaching information set  $I \in \mathcal{I}$  following joint strategy  $\sigma$ ,  $\pi^\sigma(I) = \sum_{h \in I} \pi^\sigma(h)$
- Prob. of reaching history  $h \in H$  following joint strategy  $\sigma_{-i}$ , except player  $i$  following actions in  $h$  w.p. 1,  $\pi_{-i}^\sigma(h)$
- Prob. of reaching history  $h \in H$  following player  $i$ 's actions, except others,  $\pi_i^\sigma(h)$
- Replacement of  $\sigma(I)$  to  $\delta(a)$ ,  $\sigma_{I \rightarrow a}$

## Solution concept

- Nash equilibrium  $(\sigma_1^*, \sigma_2^*) = \arg \max_{\sigma_1} \min_{\sigma_2} u_1(\sigma_1, \sigma_2)$
- Value of the game  $V = \max_{\sigma_1} \min_{\sigma_2} u_1(\sigma_1, \sigma_2)$
- **Remark:** other concepts exist in this case, NE

- Regret in hindsight w.r.t. any fixed strategy  $\sigma_1$

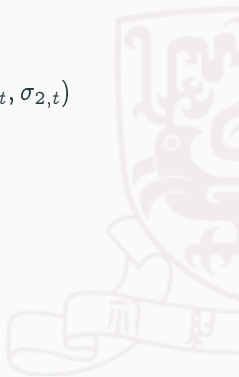
$$R_n(\sigma_1) = \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_{2,t}) - \frac{1}{n} \sum_{t=1}^n u_1(\sigma_{1,t}, \sigma_{2,t})$$

- Regret against the best strategy in hindsight

$$R_n = \max_{\sigma_1} R_n(\sigma_1)$$

- Empirical strategy:

$$\hat{\sigma}_{1,n}(I, a) = \frac{\sum_{t=1}^n \pi_i^{\sigma_t}(I) \sigma_t(I, a)}{\sum_{t=1}^n \pi_i^{\sigma_t}(I)}$$



## Theorem

A learning algorithm is Hannan's consistent if

$$\limsup_{n \rightarrow \infty} R_n \leq 0 \quad a.s.$$

Given a two-player zero-sum extensive-form game with value  $V$ , if players choose strategies  $\sigma_{1,t}$  and  $\sigma_{2,t}$  using a Hannan's consistent algorithm, then

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_{1,t}, \sigma_{2,t}) = V$$

Furthermore, the joint empirical strategy

$$\hat{\sigma}_{1,n} \times \hat{\sigma}_{2,n} \xrightarrow{n \rightarrow \infty} \{(\sigma_1^*, \sigma_2^*)\}_{Nash}$$

# Regret Matching Algorithm

- Back to Rock-Paper-Scissors
- Let  $a_1 = \text{rock}$  and  $a_2 = \text{paper}$
- Then the counterfactual regret

$$r(a_1 \rightarrow \text{rock}) = u_1(\text{rock}, a_{2,t}) - u_1(a_{1,t}, a_{2,t}) = -1 - (-1) = 0$$

$$r(a_1 \rightarrow \text{paper}) = u_1(\text{paper}, a_{2,t}) - u_1(a_{1,t}, a_{2,t}) = 0 - (-1) = 1$$

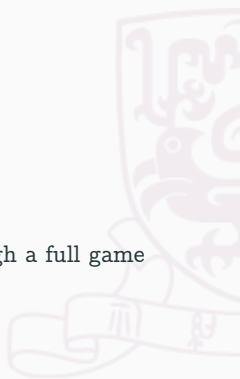
$$r(a_1 \rightarrow \text{scissors}) = u_1(\text{scissors}, a_{2,t}) - u_1(a_{1,t}, a_{2,t}) = 1 - (-1) = 2$$

- Regret matching idea

$$\sigma(a) = \frac{r(a_1 \rightarrow a)}{\sum_{b \in A_1} r(a_1 \rightarrow b)}$$

A learning problem

- For  $t = 1, \dots, n$ 
  - Player 1 chooses  $\sigma_{1,t}$
  - Player 1 executes actions prescribed by  $\sigma_{1,t}$  through a full game
  - Player 1 receives payoff  $u_{1,t}$



# Counterfactual Regret

- Counterfactual value of a history

$$v_i(\sigma, h) = \sum_{z \in Z, h \subseteq z} \pi_{-i}^\sigma(h) \pi^\sigma(h, z) u_i(z)$$

- Counterfactual regret of not taking  $a$  in  $h$

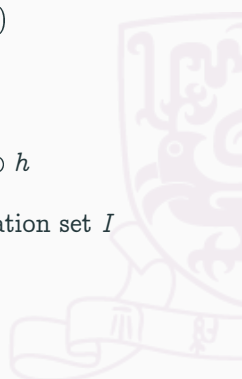
$$r_i^\sigma(h, a) = v_i(\sigma_{I \rightarrow a}, h) - v_i(\sigma, h), \quad I \supset h$$

- Counterfactual regret of not taking  $a$  in an information set  $I$

$$r_i^\sigma(I, a) = \sum_{h \in I} r_i^\sigma(h, a)$$

- Cumulative counterfactual regret

$$R_{i,t}(I, a) = \sum_{s=1}^t r_i^{\sigma^s}(I, a)$$





# Learning the Nash Equilibrium

## Version 1: Counterfactual Regret Minimization (CFR)

- For  $t = 1, \dots, n$ 
  - Player 1 chooses strategy

$$\sigma_{1,t}(l, a) = \begin{cases} \frac{R_{1,t}^+(l, a)}{\sum_{b \in A_1} R_{1,t}^+(l, b)} & \text{if } \sum_{b \in A_1} R_{1,t}^+(l, b) > 0 \\ \frac{1}{|A_1|} & \text{otherwise} \end{cases}$$

- Player 1 executes actions prescribed by  $\sigma_{1,t}$  through a full game
- Player 1 receives payoff  $u_{1,t}$
- Player 1 computes instantaneous regret  $r_i^{\sigma^t}$  over information sets observed over the game

$$R^+ = \max\{0, R\}$$

# Learning the Nash Equilibrium

## Theorem

If CFR is run over  $n$  steps, then the regret is bounded as

$$R_n = \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_{2,t}) - \frac{1}{n} \sum_{t=1}^n u_1(\sigma_{1,t}, \sigma_{2,t}) \leq |\mathcal{I}_i| \sqrt{\frac{A_1}{n}}$$

## Remarks

- $\lim_{n \rightarrow \infty} R_n \leq 0 \Rightarrow$  Hannan's consistency
- Rate of convergence  $O(1/\sqrt{n})$
- Player 1 receives payoff  $u_{1,t}$
- Linear dependence on the number of information sets
- In self-play EWF converges to the Nash equilibrium

# Learning the Nash Equilibrium

## Version 2: Counterfactual Regret Minimization+ (CFR+)

- For  $t = 1, \dots, n$ 
  - At  $t$  even player 1 chooses strategy

$$\sigma_{1,t}(l, a) = \begin{cases} \frac{Q_{1,t}(l, a)}{\sum_{b \in A_1} Q_{1,t}(l, b)} & \text{if } \sum_{b \in A_1} Q_{1,t}(l, b) > 0 \\ \frac{1}{|A_1|} & \text{otherwise} \end{cases}$$

- At  $t$  odd player 1 chooses strategy  $\sigma_{1,t} = \sigma_{1,t-1}$
  - Player 1 executes actions prescribed by  $\sigma_{1,t}$  through a full game
  - Player 1 receives payoff  $u_{1,t}$
  - Player 1 computes instantaneous regret  $r_i^{\sigma_t}$  over information sets observed over the game
- Return

$$\hat{\sigma}_{1,n} = \sum_{t=1}^n \frac{2t}{n^2 + n} \sigma_{1,t}$$

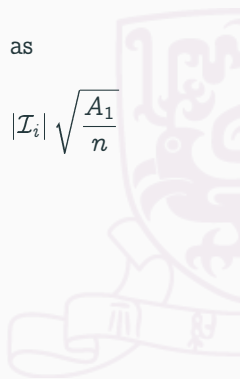
$$Q_{1,t} = (Q_{1,t-1} + r_i^{\sigma_{t-1}})^+ \quad \text{instead of } R_{1,t}^+ = \left( \sum_{s=1}^{t-1} r_i^{\sigma_s} \right)^+$$

If CFR+ is run over  $n$  steps, then the regret is bounded as

$$R_n = \max_{\sigma_1} \frac{1}{n} \sum_{t=1}^n u_1(\sigma_1, \sigma_{2,t}) - \frac{1}{n} \sum_{t=1}^n u_1(\sigma_{1,t}, \sigma_{2,t}) \leq |\mathcal{I}_i| \sqrt{\frac{A_1}{n}}$$

Remarks

- Same performance as CFR
- Empirically is more reactive
- Empirically  $\hat{\sigma}_{1,t}$  tends to converge



## The problem

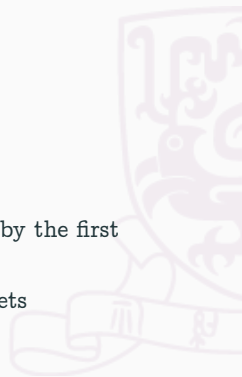
- Four rounds of cards, four rounds of betting, discrete bets
- About  $10^{18}$  states,  $3.2 \times 10^{14}$  information sets

## Abstraction: cluster together similar histories

- Symmetries (reducing to  $10^{13}$  information sets)
- Clustering
  - Buckets based on (roll-out) hand strength
  - Hierarchical buckets (e.g., second hand is indexed by the first bucket as well)
  - About  $1.65 \times 10^{12}$  states,  $5.73 \times 10^7$  information sets

## Engineering:

- Rounding:  $\sigma(a) = 0.0$  if smaller than threshold, fixed-point arithmetic
- Dynamic compression regret and strategy (from 262 TiB to 10.9 TiB)
- Distribute recursive computation of regret and strategy over rounds



# CFR in Large Problems: Heads-up Limit Texas Hold'em

