

Hannes Strass

Faculty of Computer Science, Institute of Artificial Intelligence, Computational Logic Group

# Counterfactual Regret Minimisation

Lecture 8, 10th Jun 2024 // Algorithmic Game Theory, SS 2024

# Previously ...

- A **behaviour strategy** assigns move probabilities to information sets.
- A **belief system** assigns probabilities to histories in information sets.
- An **assessment** is a pair (behaviour strategy profile, belief system).
- A **sequentially rational** assessment plays best responses “everywhere”.
- An assessment satisfies **consistency of beliefs** whenever the belief system’s probabilities match what is expected from everyone playing according to the behaviour strategy profile.
- An assessment is a **weak sequential equilibrium** iff it is both sequentially rational and satisfies consistency of beliefs.
- Mixed Nash equilibria for normal-form games and subgame perfect equilibria for sequential perfect-information games are special cases of weak sequential equilibria for extensive-form games.

# Motivation

## Main Question

- How to algorithmically solve imperfect-information games ...
- ...or at least devise good strategies or play them well in practice?

## Transformation to Normal Form?

Incurs an **exponential blowup**:

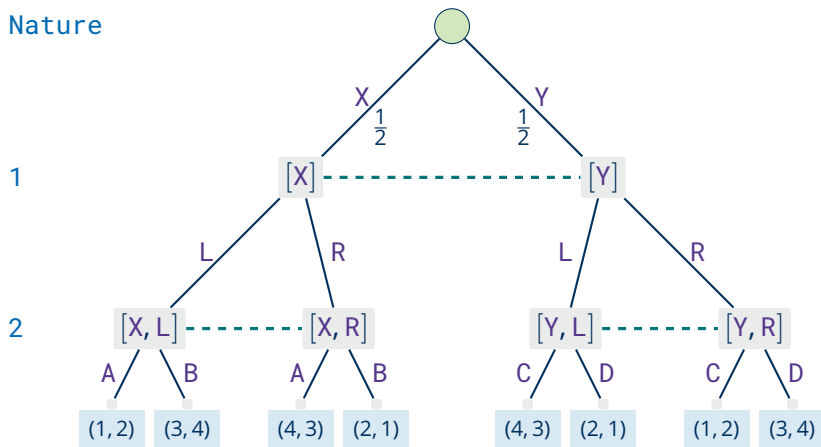
For every player  $i \in P$ , there are up to  $|M_i| |\{ \mathcal{J}_j \in \mathcal{J} \mid p^{(\mathcal{J}_j)} = i \}|$  many behaviour strategies (pure strategies in the normal-form game).

## Algorithms for sequential (perfect-information) games?

- Player  $i$ 's best move in  $\mathcal{J}_j \in \mathcal{J}$  depends on the player's beliefs  $\beta_i: \mathcal{J}_j \rightarrow [0, 1]$ .
- Consistent beliefs about  $\mathcal{J}_j$  in turn depend (in general) on probabilities of moves in other information sets (even on other paths of play).

# Motivation: Example

Nature



The best move for 2 in  $\{[X, L], [X, R]\}$  depends on what 2 does in  $\{[Y, L], [Y, R]\}$ :  
If 2 prefers C, then 1 will prefer L and thus 2 should prefer B. (Same for D and A.)

# Motivation: Regret Matching

Before minimising regret in imperfect-information extensive-form games, we start with the simpler case of normal-form games ...

## Recall

Let  $(P, \mathbf{S}, \mathbf{u})$  be a noncooperative game in normal form,  $i \in P$ , and  $s_j \in S_j$ . The **regret** of  $i$  playing  $s_j$  w.r.t. opponent profile  $\boldsymbol{\pi}_{-i}$  is

$$r_{\boldsymbol{\pi}_{-i}, s_j} := \left( \max_{\pi_k \in \Pi_i} U_i(\boldsymbol{\pi}_{-i}, \pi_k) \right) - U_i(\boldsymbol{\pi}_{-i}, s_j)$$

Difference between what player  $i$  could have had optimally vs. what they got.  
Regret is zero iff a best response is played.

↪ Minimise regret over time in order to approach playing best responses.

# Overview

Correlated Equilibria

Regret Matching

Counterfactual Regret Minimisation

# Correlated Equilibria

# Correlated Equilibria: Motivation

## Traffic Lights

Two cars both want to cross an intersection. If a car stops, it does not get to the other side. If only one car goes, it gets to the other side. If both cars go, there is an accident.

(Car1, Car2)	Stop	Go
Stop	(0,0)	(0,1)
Go	(1,0)	(-100,-100)

- The pure Nash equilibria are (Stop, Go) and (Go, Stop):  
In both equilibria, one car never gets to move.
- Another mixed Nash equilibrium is  $\left( \left( \frac{100}{101}, \frac{1}{101} \right), \left( \frac{100}{101}, \frac{1}{101} \right) \right)$ :  
Both cars mostly stop and there is a positive probability of accidents.
- A more desirable outcome would be:  $\left\{ (\text{Stop, Go}) \mapsto \frac{1}{2}, (\text{Go, Stop}) \mapsto \frac{1}{2} \right\}$ :  
However, mixed Nash equilibria cannot realise this. **Traffic lights can!**



# Correlated Equilibrium: Intuition

- An external device chooses a strategy profile  $\mathbf{s} \in \mathcal{S}$  randomly.
- The distribution  $\psi: \mathcal{S} \rightarrow [0, 1]$  for this is fixed and known to all players.
- For a chosen  $(s_1, \dots, s_n) \in \mathcal{S}$ , each player  $i \in P$  gets private advice  $s_i \in S_i$ .
- Knowing  $\{\mathbf{s} \in \mathcal{S} \mid \psi(\mathbf{s}) > 0\}$ , player  $i$  may be able to infer advice of others.
- Correlated equilibrium now means:  
No player has an incentive to deviate from the signal's advice.

## Example

In the traffic lights game, assume  $\psi = \left\{ (\text{Stop}, \text{Go}) \mapsto \frac{1}{2}, (\text{Go}, \text{Stop}) \mapsto \frac{1}{2}, \dots \right\}$ :

- If Car1 receives signal Stop, then it knows Car2 must have received Go.
- Thus its best choice is to Stop.
- Symmetrically for Car1 receiving signal Go, and Car2.

# Correlated Equilibrium: Definition

Definition [Aumann, 1974]

Let  $(P, \mathbf{S}, \mathbf{u})$  be a game in normal form with  $P = \{1, \dots, n\}$ .

A probability distribution  $\psi$  on  $\mathcal{S} = S_1 \times \dots \times S_n$  is a **correlated equilibrium** iff for every  $i \in P$ ,  $s_j \in S_j$ , and  $s_k \in S_i$ , we have

$$\sum_{\substack{\mathbf{s} \in \mathcal{S}, \\ \mathbf{s}_j = s_j}} \left( \psi(\mathbf{s}) \cdot \left( u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s}) \right) \right) \leq 0$$

**Roughly:** Following the signal's advice incurs no (positive) regret.

Observation

Every (mixed) Nash equilibrium  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_n)$  induces a correlated equilibrium  $\psi_{\boldsymbol{\pi}} := \{(s_1, \dots, s_n) \mapsto \pi_1(s_1) \cdot \dots \cdot \pi_n(s_n) \mid (s_1, \dots, s_n) \in \mathcal{S}\}$ .

**Correlated:** Players no longer mix their strategies independently.

# Correlated Equilibrium: Example (1)

## Battle of the Partners

Two partners, **Cat** and **Dee**, think about how to spend the evening. Each has their personal preference what to do, but overall they want to spend the evening together.

(Cat, Dee)	Cinema	Dancing
Cinema	(10,7)	(2,2)
Dancing	(0,0)	(7,10)

For the mixed Nash equilibrium  $\pi = (\pi_{\text{Cat}}, \pi_{\text{Dee}}) = \left( \left( \frac{2}{3}, \frac{1}{3} \right), \left( \frac{1}{3}, \frac{2}{3} \right) \right)$ , we get

$$\psi_{\pi} = \left\{ \begin{array}{l} (\text{Cinema}, \text{Cinema}) \mapsto \frac{2}{9}, (\text{Cinema}, \text{Dancing}) \mapsto \frac{4}{9}, \\ (\text{Dancing}, \text{Cinema}) \mapsto \frac{1}{9}, (\text{Dancing}, \text{Dancing}) \mapsto \frac{2}{9} \end{array} \right\}$$

with  $U_{\text{Cat}}(\psi_{\pi}) = U_{\text{Dee}}(\psi_{\pi}) = 4\frac{2}{3}$ .

## Correlated Equilibrium: Example (2)

To verify that  $\psi_\pi$  is a correlated equilibrium, we have the following cases:

- $i = \text{Cat}, s_j = \text{Cinema}, s_k = \text{Dancing}$ :

$$\begin{aligned} & \psi(\text{Cinema}, \text{Cinema}) \cdot (u_{\text{Cat}}(\text{Dancing}, \text{Cinema}) - u_{\text{Cat}}(\text{Cinema}, \text{Cinema})) + \\ & \psi(\text{Cinema}, \text{Dancing}) \cdot (u_{\text{Cat}}(\text{Dancing}, \text{Dancing}) - u_{\text{Cat}}(\text{Cinema}, \text{Dancing})) = \\ & \frac{2}{9} \cdot (0 - 10) + \frac{4}{9} \cdot (7 - 2) = -\frac{20}{9} + \frac{20}{9} \leq 0 \end{aligned}$$

- $i = \text{Cat}, s_j = \text{Dancing}, s_k = \text{Cinema}$ :

$$\begin{aligned} & \psi(\text{Dancing}, \text{Cinema}) \cdot (u_{\text{Cat}}(\text{Cinema}, \text{Cinema}) - u_{\text{Cat}}(\text{Dancing}, \text{Cinema})) + \\ & \psi(\text{Dancing}, \text{Dancing}) \cdot (u_{\text{Cat}}(\text{Cinema}, \text{Dancing}) - u_{\text{Cat}}(\text{Dancing}, \text{Dancing})) = \\ & \frac{1}{9} \cdot (10 - 0) + \frac{2}{9} \cdot (2 - 7) = \frac{10}{9} + \left(-\frac{10}{9}\right) \leq 0 \end{aligned}$$

Due to  $u_{\text{Dee}}(s_1, s_2) = u_{\text{Cat}}(s_2, s_1)$ , this also covers the cases for  $i = \text{Dee}$ .

# Correlated Equilibria: Example (3)

Assume that both **Cat** and **Dee** have access to the result of one fair coin toss:

- If the coin shows heads, both go to the concert;
- if the coin shows tails, both go to the cinema.

This leads to the following (additional) correlated equilibrium:

$$\psi = \left\{ (\text{Cinema}, \text{Cinema}) \mapsto \frac{1}{2}, (\text{Dancing}, \text{Dancing}) \mapsto \frac{1}{2}, \dots \right\}$$

with associated payoffs  $U_{\text{Cat}}(\psi) = U_{\text{Dee}}(\psi) = \frac{1}{2} \cdot 10 + \frac{1}{2} \cdot 7 = 8\frac{1}{2}$ .

To verify that  $\psi$  is a correlated equilibrium, we (essentially) verify that:

$$\begin{aligned} \psi(\text{Cinema}, \text{Cinema}) \cdot (u_{\text{Cat}}(\text{Dancing}, \text{Cinema}) - u_{\text{Cat}}(\text{Cinema}, \text{Cinema})) &\leq 0 \\ \psi(\text{Dancing}, \text{Dancing}) \cdot (u_{\text{Cat}}(\text{Cinema}, \text{Dancing}) - u_{\text{Cat}}(\text{Dancing}, \text{Dancing})) &\leq 0 \end{aligned}$$

which holds because  $\frac{1}{2} \cdot (0 - 10) = -5 \leq 0$  and  $\frac{1}{2} \cdot (2 - 7) = -2\frac{1}{2} \leq 0$ .

# Correlated Equilibria Form a Convex Set

## Theorem

Let  $G = (P, \mathbf{S}, \mathbf{u})$  be a strategic game in normal form.

For any two correlated equilibria  $\psi_1$  and  $\psi_2$ , and for any  $\alpha \in [0, 1]$ , we find that  $\psi_\alpha := \{\mathbf{s} \mapsto \alpha \cdot \psi_1(\mathbf{s}) + (1 - \alpha) \cdot \psi_2(\mathbf{s}) \mid \mathbf{s} \in \mathcal{S}\}$  is a correlated equilibrium.

## Proof.

Let  $\alpha \in [0, 1]$  and consider any  $i \in P, s_j, s_k \in S_i$ . We have

$$\begin{aligned} & \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} (\psi_\alpha(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s}))) \\ &= \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} ((\alpha \cdot \psi_1(\mathbf{s}) + (1 - \alpha) \cdot \psi_2(\mathbf{s})) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s}))) \\ &= \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} \left( (\alpha \cdot \psi_1(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s}))) + ((1 - \alpha) \cdot \psi_2(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s}))) \right) \\ &= \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} \left( \alpha \cdot \psi_1(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s})) \right) + \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} \left( (1 - \alpha) \cdot \psi_2(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s})) \right) \\ &= \alpha \cdot \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} \left( \psi_1(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s})) \right) + (1 - \alpha) \cdot \sum_{\mathbf{s} \in \mathcal{S}, s_j = s_j} \left( \psi_2(\mathbf{s}) \cdot (u_i(s_k, \mathbf{s}_{-i}) - u_i(\mathbf{s})) \right) \leq 0 \quad \square \end{aligned}$$

# Regret Matching

# Learning to Play

## Learning in Games: General Setting

- A (normal-form) game is played repeatedly for time points  $t = 1, 2, \dots$
- After the game at time point  $t$  has ended, player (say)  $i$  has access to all strategy profiles  $\mathbf{s}^1, \mathbf{s}^2, \dots, \mathbf{s}^t$  played previously, and their payoffs to  $i$ .
- Using this information, the player can devise a (mixed) strategy  $\pi_i^{t+1}$  to play at time point  $t + 1$ .

How can we evaluate whether a learner (player) is “doing well”?

## Hindsight Rationality

After playing the game for  $t \rightarrow \infty$  time points, the player “cannot think of” a function  $\Phi: \Pi_i \rightarrow \Pi_i$  that would strictly increase their payoff in hindsight.

Can **learning** (dynamic, local) lead to **equilibria** (static, global)?



# Regret Matching

In what follows, we assume a fixed normal-form game  $G = (P, \mathbf{S}, \mathbf{u})$  to be played at time points  $t = 1, 2, \dots, T$  and take the perspective of  $i \in P$ .

At each time step  $t \leq T$ ,  $i$ 's one-time regret of not having played  $s_k \in S_i$  is:

$$r_i^t(s_k) := u_i(s_k, \mathbf{s}_{-i}^t) - u_i(\mathbf{s}^t)$$

At time point  $T$ , the accumulated regret of a strategy  $s_k \in S_i$  is thus:

$$R_i^T(s_k) := \sum_{1 \leq t \leq T} r_i^t(s_k)$$

The probabilities at  $T + 1$  are then set to be proportional to **positive** regret:

$$\pi_i^{T+1}(s_j) := \begin{cases} \frac{[R_i^T(s_j)]^+}{R_i^{T,+}} & \text{if } R_i^{T,+} > 0, \quad \text{where } R_i^{T,+} := \sum_{s_k \in S_i} [R_i^T(s_k)]^+, \\ \frac{1}{|S_i|} & \text{otherwise.} \end{cases} \quad \text{for } s_j \in S_i$$

( $[x]^+ := \max\{x, 0\}$  for all  $x \in \mathbb{R}$ .)

# Regret Matching: Example

(Cat, Dee)	Cinema	Dancing
Cinema	(10,7)	(2,2)
Dancing	(0,0)	(7,10)

We denote  $\overline{\text{Cinema}} = \text{Dancing}$  and  $\overline{\text{Dancing}} = \text{Cinema}$ .

$T$	$\mathbf{s}^T = (s_{\text{Cat}}^T, s_{\text{Dee}}^T)$	$r_{\text{Cat}}^T(\overline{s_{\text{Cat}}^T})$	$R_{\text{Cat}}^T(\text{Cinema})$	$R_{\text{Cat}}^T(\text{Dancing})$	$\pi_{\text{Cat}}^{T+1}$
1	(Cinema, Dancing)	5	0	5	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}
2	(Dancing, Cinema)	10	10	5	{Cinema $\mapsto$ $\frac{2}{3}$ , Dancing $\mapsto$ $\frac{1}{3}$ }
3	(Cinema, Dancing)	5	10	10	{Cinema $\mapsto$ $\frac{1}{2}$ , Dancing $\mapsto$ $\frac{1}{2}$ }
4	(Cinema, Cinema)	-10	10	0	{Cinema $\mapsto$ 1, Dancing $\mapsto$ 0}
5	(Cinema, Dancing)	5	10	5	{Cinema $\mapsto$ $\frac{2}{3}$ , Dancing $\mapsto$ $\frac{1}{3}$ }
6	(Cinema, Cinema)	-10	10	-5	{Cinema $\mapsto$ 1, Dancing $\mapsto$ 0}

# Regret Matching: Correctness

For a given play sequence  $(\mathbf{s}^t)_{t=1}^T$ , and every  $\mathbf{s}' \in \mathcal{S}$ , define the **relative frequency** of  $\mathbf{s}'$  after  $T$  rounds via

$$\bar{\varphi}^T(\mathbf{s}') := \frac{1}{T} \cdot |\{1 \leq t \leq T \mid \mathbf{s}^t = \mathbf{s}'\}|$$

Theorem [Hart and Mas-Colell, 2000]

Let  $G = (P, \mathbf{S}, \mathbf{u})$  be a noncooperative game in normal form.

If every player plays according to regret matching, then  $(\bar{\varphi}^t)_{t=1}^T$  converges to the set of correlated equilibria of  $G$  as  $T \rightarrow \infty$ .

**More precisely:** For any  $\varepsilon > 0$ , there is a  $T_0 \geq 0$  such that for all  $T > T_0$ , there is a correlated equilibrium  $\psi_T$  of  $G$  whose distance from  $\bar{\varphi}^T$  is at most  $\varepsilon$ .

**Note:** The result does not say that relative frequencies converge to a *point*.

↪ Since all players must use regret matching, it will be used in **self-play**.

# Regret Matching in Self-Play: Example

(Cat, Dee)	Cinema	Dancing
Cinema	(10,7)	(2,2)
Dancing	(0,0)	(7,10)

We denote  $R_i^T = (R_i^T(\text{Cinema}), R_i^T(\text{Dancing}))$  for  $i \in \{\text{Cat}, \text{Dee}\}$ .

$T$	$\mathbf{s}^T = (s_{\text{Cat}}^T, s_{\text{Dee}}^T)$	$R_{\text{Cat}}^T$	$R_{\text{Dee}}^T$	$\pi_{\text{Cat}}^{T+1}$	$\pi_{\text{Dee}}^{T+1}$
1	(Cinema, Dancing)	(0,5)	(5,0)	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}	{Cinema $\mapsto$ 1, Dancing $\mapsto$ 0}
2	(Dancing, Cinema)	(10,5)	(5,10)	{Cinema $\mapsto$ $\frac{2}{3}$ , Dancing $\mapsto$ $\frac{1}{3}$ }	{Cinema $\mapsto$ $\frac{1}{3}$ , Dancing $\mapsto$ $\frac{2}{3}$ }
3	(Cinema, Dancing)	(10,10)	(10,10)	{Cinema $\mapsto$ $\frac{1}{2}$ , Dancing $\mapsto$ $\frac{1}{2}$ }	{Cinema $\mapsto$ $\frac{1}{2}$ , Dancing $\mapsto$ $\frac{1}{2}$ }
4	(Dancing, Dancing)	(5,10)	(0,10)	{Cinema $\mapsto$ $\frac{1}{3}$ , Dancing $\mapsto$ $\frac{2}{3}$ }	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}
5	(Cinema, Dancing)	(5,15)	(5,10)	{Cinema $\mapsto$ $\frac{1}{4}$ , Dancing $\mapsto$ $\frac{3}{4}$ }	{Cinema $\mapsto$ $\frac{1}{3}$ , Dancing $\mapsto$ $\frac{2}{3}$ }
6	(Dancing, Dancing)	(0,15)	(-5,10)	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}
7	(Dancing, Dancing)	(-5,15)	(-15,10)	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}	{Cinema $\mapsto$ 0, Dancing $\mapsto$ 1}

# Rate of Convergence

For a given sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$  of mixed-strategy profiles, define the (external) **overall regret** of player  $i \in P$  after  $T$  rounds via

$$R_i^T := \max_{\hat{\boldsymbol{\pi}} \in \Pi_i} \left\{ \sum_{t=1}^T (U_i(\hat{\boldsymbol{\pi}}, \boldsymbol{\pi}_{-i}^t) - U_i(\boldsymbol{\pi}_i^t, \boldsymbol{\pi}_{-i}^t)) \right\}$$

## Theorem

Let  $G = (P, \mathbf{S}, \mathbf{u})$  be a normal-form game and let player  $i \in P$  use regret matching in the sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$  of mixed-strategy profiles.

Then  $R_i^T \leq \omega \cdot \sqrt{T}$ , where the constant  $\omega \in \mathbb{R}$  depends only on  $\mathbf{u}$ .

The **average overall regret** is then  $\bar{R}_i^T := \frac{1}{T} \cdot R_i^T$ .

## Proposition

$\bar{R}_i^T$  tends to zero as  $T \rightarrow \infty$  iff  $\bar{\boldsymbol{\varphi}}^T$  tends to the set of correlated equilibria.

# The Case of Two-Player Zero-Sum Games

For a given sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$  of mixed-strategy profiles, define the **average mixed strategy**  $\bar{\pi}_i^T$  of player  $i \in P$  after  $T$  rounds via

$$\bar{\pi}_i^T(s_j) := \frac{1}{T} \cdot \sum_{t=1}^T \pi_i^t(s_j) \quad \text{for } s_j \in S_i$$

## Theorem

Let  $G = (P, \mathbf{S}, \mathbf{u})$  be a **two-player, zero-sum** normal-form game, i.e.  $P = \{1, 2\}$ , and let  $(\boldsymbol{\pi}^t)_{t=1}^T$  be obtained from both players using regret matching. Then as  $T \rightarrow \infty$ , the pair  $(\bar{\pi}_1^T, \bar{\pi}_2^T)$  converges to the set of **Nash equilibria** of  $G$ .

# Regret Matching<sup>+</sup>

The computation of the accumulated (possibly negative) regret of a strategy  $s_k \in S_i$  can be rewritten as:

$$R_i^T(s_k) := R_i^{T-1}(s_k) + r_i^T(s_k) \quad \text{with } R_i^0(s_k) := 0$$

Tammelin [2014] observed a better convergence when this is replaced by

$$R_i^{T,+}(s_k) := \left[ R_i^{T-1}(s_k) \right]^+ + r_i^T(s_k)$$

The probabilities at  $T + 1$  are again set to be proportional to positive regret:

$$\pi_i^{T+1}(s_j) := \begin{cases} \frac{R_i^{T,+}(s_j)}{R_i^{T,+}} & \text{if } R_i^{T,+} > 0, \quad \text{where } R_i^{T,+} := \sum_{s_k \in S_i} R_i^{T,+}(s_k), \\ \frac{1}{|S_i|} & \text{otherwise.} \end{cases} \quad \text{for } s_j \in S_i$$

RM<sup>+</sup> reacts more quickly when a previously poor action improves over time.

# Counterfactual Regret Minimisation



# From Normal Form to Extensive Form

## Solving Imperfect-Information Games: Main Ideas

- Traverse the game tree in a backward induction-like fashion.
- Apply regret matching at each decision point (information set).

## Problem

Optimal moves depend on probabilities of moves in other information sets.

## Solution of Zinkevich, Johanson, Bowling, and Piccione [2007]

- Define new notion of **counterfactual regret**:  
Assume the player played to deliberately reach a certain information set.
- Then for **games with perfect recall**:
  - Regret matching can be applied to each information set independently.
  - Counterfactual regret is an upper bound for actual regret (main theorem).
  - Thus minimising counterfactual regret minimises actual regret.

# Remember, Remember

## Recall

$P(h' | h, \boldsymbol{\pi})$  is the probability that  $h'$  is reached when playing  $\boldsymbol{\pi}$  from  $h$  on:

- $P(h | h, \boldsymbol{\pi}) = 1$  for all  $h \in H$ ,
- $P(\square | h, \boldsymbol{\pi}) = 0$  for all  $h \neq \square$ , and
- $P([h'; m] | h, \boldsymbol{\pi}) = \pi_{p(\mathcal{J}_{h'})}(m | \mathcal{J}_{h'}) \cdot P(h' | h, \boldsymbol{\pi})$ .

## Recall

The probability of reaching information set  $\mathcal{J}_j$  when playing  $\boldsymbol{\pi}$  is thus

$$P(\mathcal{J}_j | \boldsymbol{\pi}) := \sum_{h \in \mathcal{J}_j} P(h | \boldsymbol{\pi}) \quad \text{where } P(h | \boldsymbol{\pi}) \text{ denotes } P(h | \square, \boldsymbol{\pi})$$

## Recall

Player  $i$ 's expected utility of playing  $\boldsymbol{\pi}$  when history  $h$  has been reached is

$$U_i(\boldsymbol{\pi} | h) := \sum_{z \in Z} P(z | h, \boldsymbol{\pi}) \cdot u_i(z)$$

# Towards Counterfactual Regret

## Definition

Consider an extensive-form game with player  $i \in P$  and information sets  $\mathcal{J}$ .

1. The **counterfactual probability** of **playing to reach**  $h \in H$  is given by

$$P([\ ] | \boldsymbol{\pi}_{-i}) = 1 \text{ and } P([h'; m] | \boldsymbol{\pi}_{-i}) := \begin{cases} \pi_k(m | h') \cdot P(h' | \boldsymbol{\pi}_{-i}) & \text{if } p(h') = k \neq i, \\ P(h' | \boldsymbol{\pi}_{-i}) & \text{otherwise.} \end{cases}$$

2. The **counterfactual probability** of **playing to reach**  $\mathcal{J}_j \in \mathcal{J}$  is

$$P(\mathcal{J}_j | \boldsymbol{\pi}_{-i}) := \sum_{h \in \mathcal{J}_j} P(h | \boldsymbol{\pi}_{-i})$$

3. The **counterfactual utility** of **playing to reach**  $\mathcal{J}_j$  and then playing  $\boldsymbol{\pi}$  is

$$U_i(\boldsymbol{\pi} | \mathcal{J}_j) = \frac{\sum_{h \in \mathcal{J}_j} P(h | \boldsymbol{\pi}_{-i}) \cdot U_i(\boldsymbol{\pi} | h)}{P(\mathcal{J}_j | \boldsymbol{\pi}_{-i})} = \frac{\sum_{h \in \mathcal{J}_j} P(h | \boldsymbol{\pi}_{-i}) \cdot \sum_{z \in Z} P(z | h, \boldsymbol{\pi}) \cdot u_i(z)}{P(\mathcal{J}_j | \boldsymbol{\pi}_{-i})}$$

We **counterfactually** assume that  $i$  **intentionally** played to reach  $\mathcal{J}_j$ .

# Counterfactual Regret

## Definition

Consider  $i \in P$  and  $\mathcal{J}_j \in \mathcal{J}$  with  $p(\mathcal{J}_j) = i$ .

1. Denote the set of legal moves of  $i$  in  $\mathcal{J}_j$  by

$$M_i(\mathcal{J}_j) := \{m \in M_i \mid [h; m] \in H \text{ for some } h \in \mathcal{J}_j\}$$

2. For behaviour strategy profile  $\boldsymbol{\pi}$  and move  $m \in M_i(\mathcal{J}_j)$ , define modified profile  $\langle \boldsymbol{\pi} \rangle_m^{\mathcal{J}_j}$  to be just like  $\boldsymbol{\pi}$ , except that in  $\mathcal{J}_j$  it always chooses  $m$ .

3. The **immediate counterfactual regret** at time  $T$  is then defined by

$$r_i^T(\mathcal{J}_j) := \max_{m^* \in M_i(\mathcal{J}_j)} \sum_{t=1}^T P(\mathcal{J}_j \mid \boldsymbol{\pi}_{-i}^t) \cdot \left( U_i \left( \langle \boldsymbol{\pi}^t \rangle_{m^*}^{\mathcal{J}_j} \mid \mathcal{J}_j \right) - U_i(\boldsymbol{\pi}^t \mid \mathcal{J}_j) \right)$$

for any sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$  of behaviour strategy profiles.

**Key Feature:**  $r_i^T$  can be minimised by controlling only  $\pi_i(\mathcal{J}_j): M_i(\mathcal{J}_j) \rightarrow [0, 1]$ .

# Overall Regret $\leq$ Immediate Regret

Given sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$ , the (external) **overall regret** of player  $i$  at time  $T$  is:

$$R_i^T = \max_{\boldsymbol{\pi}_i^* \in \Pi_i} \sum_{t=1}^T (U_i(\boldsymbol{\pi}_i^*, \boldsymbol{\pi}_{-i}^t) - U_i(\boldsymbol{\pi}^t))$$

where  $U_i(\boldsymbol{\pi})$  denotes  $U_i(\boldsymbol{\pi} \mid \square)$ .

Theorem [Zinkevich, Johanson, Bowling, and Piccione, 2007]

In any extensive-form game with perfect recall, for any player  $i \in P$  and any sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$  of behaviour strategy profiles:

$$R_i^T \leq \sum_{\substack{J_j \in \mathcal{J}, \\ \rho(J_j)=i}} \left[ r_i^T(J_j) \right]^+$$

**Thus:** Minimising immediate regret in each  $J_j$  minimises overall regret.

# Regret Matching at Information Sets

## Definition

Consider the sequence  $(\boldsymbol{\pi}^t)_{t=1}^T$  of behaviour strategy profiles of past play.

1. Let  $\mathcal{J}_j \in \mathcal{J}$  with  $p(\mathcal{J}_j) = i$  and  $m \in M_i(\mathcal{J}_j)$ . The **accumulated regret** of  $m$  is

$$R_i^T(\mathcal{J}_j, m) := \sum_{t=1}^T P(\mathcal{J}_j | \boldsymbol{\pi}_{-i}^t) \cdot \left( U_i \left( \langle \boldsymbol{\pi}^t \rangle_m^{\mathcal{J}_j} \mid \mathcal{J}_j \right) - U_i(\boldsymbol{\pi}^t \mid \mathcal{J}_j) \right)$$

2. The probability of playing  $m$  at  $\mathcal{J}_j$  at time  $T + 1$  is set to

$$\pi_i^{T+1}(\mathcal{J}_j)(m) := \begin{cases} \frac{[R_i^T(\mathcal{J}_j, m)]^+}{R_i^{T,+}} & \text{if } R_i^{T,+} > 0, \quad \text{where } R_i^{T,+} := \sum_{m \in M_i(\mathcal{J}_j)} [R_i^T(\mathcal{J}_j, m)]^+ \\ \frac{1}{|M_i(\mathcal{J}_j)|} & \text{otherwise.} \end{cases}$$

# CFR: Algorithm (1)

Initialisation of global variables:

```
function init() {  
    foreach  $i \in \{1, 2\}$  do {  
        foreach  $\mathcal{J}_j \in \mathcal{J}$  with  $p(\mathcal{J}_j) = i$  do {  
            foreach  $m \in M_i(\mathcal{J}_j)$  do {  
                 $regret[j][m] := 0$  // accumulated regret table  
                 $strategy[j][m] := 0$  // accumulated strategy table  
                 $profile[1][j][m] := 1/|M_i(\mathcal{J}_j)|$  // move distribution for  $\mathcal{J}_j$  at  $t = 1$   
            }  
        }  
    }  
}
```

Main Loop:

```
function solve( $T$ ) {  
    foreach  $t \in \{1, 2, \dots, T\}$  do {  
        foreach  $i \in \{1, 2\}$  do {  
            cfr( $\square, i, t, 1, 1$ )  
        }  
    }  
}
```

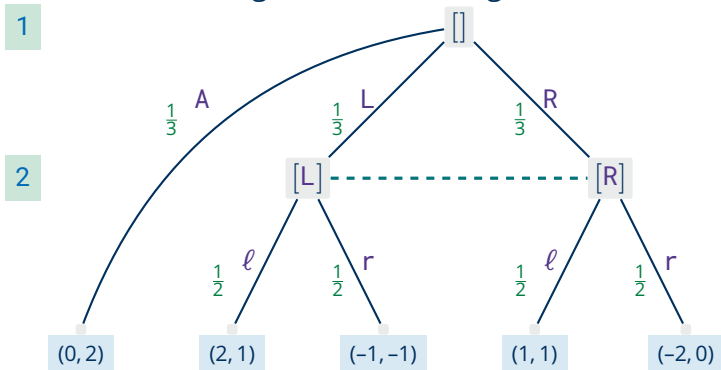
# CFR: Algorithm (2)

```
function cfr( $h, i, t, p_1, p_2$ ) { // history, player, time point, reach probabilities
  if IS-TERMINAL( $h$ ) then return UTILITY $_i(s)$ 
   $v_h := 0$  // initialise expected payoff at  $h \in \mathcal{J}_j$ 
  foreach  $m \in M_{p(\mathcal{J}_j)}(\mathcal{J}_j)$  do {  $v'_h[j][m] := 0$  } // initialise payoffs of single moves
  foreach  $m \in M_{p(\mathcal{J}_j)}(\mathcal{J}_j)$  do {
    if TURN( $h$ ) = 1 then {  $v'_h[j][m] := \mathbf{cfr}([h; m], i, t, \text{profile}[t][j][m] \cdot p_1, p_2)$  }
    else {  $v'_h[j][m] := \mathbf{cfr}([h; m], i, t, p_1, \text{profile}[t][j][m] \cdot p_2)$  }
     $v_h := v_h + \text{profile}[t][j][m] \cdot v'_h[j][m]$  // accumulate currently expected payoff
  }
  if TURN( $h$ ) =  $i$  then { // players minimise immediate regret of own moves
     $r^+ := 0$  // initialise sum of positive regrets
    for  $m \in M_i(\mathcal{J}_j)$  do { // update values needed for regret matching
       $\text{regret}[j][m] := \text{regret}[j][m] + p_{3-i} \cdot (v'_h[m] - v_h)$  // update accumulated cf regret
       $\text{strategy}[j][m] := \text{strategy}[j][m] + p_i \cdot \text{profile}[t][j][m]$  // update "frequency" of move
       $r^+ := r^+ + [\text{regret}[j][m]]^+$  // accumulate positive regret sum for normalisation
    }
    if  $r^+ > 0$  then { foreach  $m \in M_i(\mathcal{J}_j)$  do { // apply regret matching at  $\mathcal{J}_j$ 
       $\text{profile}[t+1][j][m] := [\text{regret}[j][m]]^+ / r^+$  } }
    else { foreach  $m \in M_i(\mathcal{J}_j)$  do {
       $\text{profile}[t+1][j][m] := 1 / |M_i(\mathcal{J}_j)|$  } } }
  }
  return  $v_h$  }
```



# CFR: Example

Recall the following extensive-form game  $G_4$ :



- (1) Initialise move probabilities by uniform distributions
- (2) Traverse game tree for  $T = 1, i = 1$
- (3) Traverse game tree for  $T = 1, i = 2$
- (4) Update move probabilities according to regret matching

# CFR: Convergence and Correctness

Theorem [Zinkevich, Johanson, Bowling, and Piccione, 2007]

For any extensive-form game with perfect recall, if player  $i$  selects actions according to regret matching at information sets, then

$$r_i^T(\mathcal{J}_j) \leq \omega \cdot \sqrt{|M'_i|} \cdot \sqrt{T} \quad \text{whence} \quad R_i^T \leq \omega \cdot |\{\mathcal{J}_j \in \mathcal{J} \mid p(\mathcal{J}_j) = i\}| \cdot \sqrt{|M'_i|} \cdot \sqrt{T}$$

where  $\omega \in \mathbb{R}$  only depends on  $\mathbf{u}$ , and  $|M'_i| := \max_{\mathcal{J}_j \in \mathcal{J}_{p(\mathcal{J}_j)=i}} |M_i(\mathcal{J}_j)|$ .

- The bound on overall regret is **linear** in the number of information sets.
- The overall regret grows **sublinearly** in  $T$ , so the **average overall regret**  $\bar{R}_i^T := \frac{1}{T} \cdot R_i$  tends to zero as  $T \rightarrow \infty$ .

Theorem

In any two-player, zero-sum extensive-form game with perfect recall, if both players select actions according to regret matching at information sets, then the average strategy profiles tend to the set of Nash equilibria as  $T \rightarrow \infty$ .

# CFR Algorithm: Remarks

- Histories/information sets of **Nature** can be treated in the algorithm via sampling a move from  $M_{\text{Nature}}(J_j)$  with the specified distribution.
- At each time step  $t = 1, 2, \dots, T$  (and for each  $i \in P$ ), the call to **cfr**( $[], i, t, 1, 1$ ) leads to a full traversal of the game tree.
- After **solve**( $T$ ), the final values of  $strategy[j][m]$  can be normalised to obtain the behaviour strategies tending towards Nash equilibria.
- Additional techniques, e.g. game abstraction, are used in practice to reduce the number of information sets (per player) to a manageable size.
- By using regret matching<sup>+</sup> in place of regret matching, we obtain CFR<sup>+</sup>.
- CFR<sup>+</sup> also uses linear weighting to compute average strategies:

$$\bar{\pi}_i^{T,+}(s_j) := \frac{2}{T^2+T} \cdot \sum_{t=1}^T (t \cdot \pi^t(s_j))$$

- Bowling et al. [2015] used CFR<sup>+</sup> (with additional optimisations) to “essentially weakly solve” heads-up limit hold'em poker.

# Conclusion

## Summary

- The **regret** is the difference between a player's best possible strategy and their actual strategy.
- A **correlated equilibrium** can be seen as providing players with private signals they can use to best-respond to each other's strategies.
- The **regret matching** algorithm uses self-play to steer play towards the set of correlated equilibria.
- In the case of two-player zero-sum games, regret matching tends towards the set of (mixed) Nash equilibria.
- The **counterfactual regret minimisation** algorithm applies regret matching to every information set of an (imperfect-information) extensive-form game (with perfect recall).

**Action Point:** Implement  $\text{CFR}^{(+)}$  and use it to solve Simplified Poker.