# Solving Normal-Form Games

Brian Zhang

# Recap: Normal-Form Games

|  | 👊 | ✋ | ✌️ |
|---|---|---|---|
| **0.2** 👊 | 0 | -1 | +1 |
| **0.5** ✋ | +1 | 0 | -1 |
| **0.3** ✌️ | -1 | +1 | 0 |

✳ SIMULTANEOUS

(No turns)

✳ Strategy for a player is just a probability distribution over actions

# Two-Player Zero-Sum Normal-Form Games

- NE doesn't have problems as in general-sum or multiplayer games

- In a sense, NE is optimal in that no opponent can exploit you

  - If I were to play any other strategy than ⅓, ⅓, ⅓ in rock paper scissors, you could exploit me

- NE can leave utility on the table against imperfect opponents

  - If you always play Rock, NE will still just play ⅓, ⅓, ⅓

- But this is a price usually worth paying when playing experts or other AI programs

|   | R | P | S |
|---|---|---|---|
| R | 0 | -1 | 1 |
| P | 1 | 0 | -1 |
| S | -1 | 1 | 0 |

# Computing NE in Two-Player Zero-Sum Normal-Form Games (This Lecture)

1. LP for small games
2. Iterative Approaches
   - Best-Response Dynamics (doesn't converge)
   - Fictitious Play aka Follow the Leader (FTL)
3. No-Regret Algorithms
   - Follow the Regularized Leader (FTRL)
   - Regret Matching
4. Optimistic regret minimization

*Running example:*
Weighted RPS

|   | R | P | S |
|---|---|---|---|
| R | 0 | -2 | 1 |
| P | 2 | 0 | -1 |
| S | -1 | 1 | 0 |

# LP Approach

$$\max_{x \in \Delta^m} \min_{y \in \Delta^n} x^\top A y$$

|     |     |     | $y$ |     |     |
| --- | --- | --- | --- | --- | --- |
|     |     |     | 1/4 | 1/4 | 1/2 |
|     |     |     | R   | P   | S   |
|     | 1/4 | R   | 0   | -2  | 1   |
| $x$ | 1/4 | P   | 2   | 0   | -1  |
|     | 1/2 | S   | -1  | 1   | 0   |

# LP Approach

$$\max_{x \in \mathbb{R}^m} \begin{cases} \min_{y \in \mathbb{R}^n} \quad x^\top A y \\ \text{s.t. } \mathbf{1}^\top y = \mathbf{1}, \\ \qquad y \geq \mathbf{0} \end{cases}$$

$$\text{s.t. } \mathbf{1}^\top x = \mathbf{1},$$
$$x \geq \mathbf{0}$$

LP duality

find the largest value $v$ s.t.

$$\max_{v \in \mathbb{R}} \quad v$$

every strategy of the opponent gives us expected value at least $v$

$$\text{s.t.} \quad A^\top x \geq \mathbf{1} v$$

| | | | $y$ | | |
|---|---|---|---|---|---|
| | | | 1/4 | 1/4 | 1/2 |
| | | | R | P | S |
| | 1/4 | R | 0 | -2 | 1 |
| $x$ | 1/4 | P | 2 | 0 | -1 |
| | 1/2 | S | -1 | 1 | 0 |

# LP Approach

$$\max_{\substack{x \in \mathbb{R}^m \\ v \in \mathbb{R}}} v$$

find the largest value $v$ s.t.
for some $x$

$$\text{s.t. } \mathbf{1}^\top x = 1,$$
$$x \geq 0$$

$x$ is a valid mixed strategy

$$A^\top x \geq \mathbf{1}v$$

every strategy of the opponent gives us expected value at least $v$

find the largest value $v$ s.t.

every strategy of the opponent gives us expected value at least $v$

$$\max_{v \in \mathbb{R}} v$$

$$\text{s.t. } A^\top x \geq \mathbf{1}v$$

|   |   | $y$ |   |   |
|---|---|---|---|---|
|   |   | 1/4 | 1/4 | 1/2 |
|   |   | R | P | S |
| $x$ | 1/4 | R | 0 | -2 | 1 |
|   | 1/4 | P | 2 | 0 | -1 |
|   | 1/2 | S | -1 | 1 | 0 |

# LP Approach

- Solving our game results in the following
- We maximize the value that the opponent can get against us
- Any deviation would allow the opponent to exploit us more

For **P2's** strategy: take dual values of constraint $A^\top x \geq 1v$, or solve
$$\min_{y \in \Delta^n} \max_{x \in \Delta^m} x^\top A y$$

|  |  | R | P | S |
|---|---|---|---|---|
| 1/4 | R | 0 | -2 | 1 |
| 1/4 | P | 2 | 0 | -1 |
| 1/2 | S | -1 | 1 | 0 |

|  |  | R | P | S |
|---|---|---|---|---|
|  | EV | 0 | 0 | 0 |

# Iterative Approaches

*We'll make this precise soon*

- Only relatively small games can be solved via LP
- For larger games we need iterative approaches
- Most iterative approaches *approach* a NE *on average*
  - Can be stopped any time
- What we'll cover
  - Best Response Dynamics (doesn't converge to NE)
  - Fictitious Play aka Follow the Leader (isn't no-regret)
  - Follow the Regularized Leader (*e.g.,* gradient descent, multiplicative weights)
  - Regret Matching
  - Regret Matching Plus
  - Optimistic regret minimization

# Best Reponse Dynamics

$$x_i^{t+1} = \arg \max_{x_i} \; u_i\left(x_i, x_{-i}^t\right)$$

*Best respond to the opponent's **last** strategy*

Question: Does

$$\frac{1}{T}\sum_{t=1}^{T} x_i^t \; \xrightarrow{T \to \infty} \; \text{NE?}$$

**No!**

$$\frac{1}{T}\sum_{t=1}^{T} x_i^t \to \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix} \neq \text{NE} = \begin{bmatrix} 1/4 \\ 1/4 \\ 1/2 \end{bmatrix}$$

|   | R | P | S |
|---|---|---|---|
| R | 0 | -2 | 1 |
| P | 2 | 0 | -1 |
| S | -1 | 1 | 0 |

# Fictitious Play (Follow the Leader)

$$x_i^{t+1} = \arg\max_{x_i} \; \frac{1}{t}\sum_{\tau=1}^{t} u_i(x_i, x_{-i}^{\tau})$$

*Best respond to the opponent's **average** strategy*

Question: Does

$$\frac{1}{T}\sum_{t=1}^{T} x_i^t \xrightarrow{T\to\infty} \text{NE?}$$

**Yes!** (for zero-sum games)
[Robinson 1951]

…but possibly with very slow rate $T^{-1/n}$
if the tiebreaking is done adversarially
[Daskalakis & Pan 2014]

***Open question*** ["Karlin's weak conjecture", Karlin 1959]:
Does FP *with non-adversarial tiebreaking* converge with
rate $O_n(T^{-1/2})$ in all zero-sum normal-form games?

# No-Regret Algorithms

- What if I'm playing a repeated game against someone who knows I am playing fictitious play?
- Then they would know exactly what my next move will be and could choose a best response every time
- Can we find iterative algorithms that will not be *too bad* even when the opponent knows the algorithm?
- No-regret algorithms do exactly this
    - And achieve faster convergence than FP as well!

# Regret Minimization

for $t = 1, \ldots, T$:

- Agent chooses an *action distribution* $x^t \in X := \Delta^n$
- Environment chooses a *utility vector* $u^t \in [0,1]^n$
- Agent observes $u^t$ and gets utility $\langle u^t, x^t \rangle$

$\Delta^n$ = set of distributions on $n$ things
$= \{x \in \mathbb{R}^n : x \geq 0, \sum x_i = 1\}$

Agent goal: Minimize *regret.*

"How well do we do against best, fixed strategy in hindsight?"

$$R^T := \boxed{\max_{\widehat{x} \in X} \left\{ \sum_{t=1}^{T} \langle u^t, \widehat{x} \rangle \right\}} - \boxed{\sum_{t=1}^{T} \langle u^t, x^t \rangle}$$

Maximum utility that was achievable by the **best fixed** action in hindsight

Utility that was actually accumulated

❇️ Goal: have $R^T$ grow sublinearly with respect to time $T$, e.g., $R^T = O_n(\sqrt{T})$

No assumption on utilities!
Must be able to handle adversarial environments

# What does regret minimization have to do with zero-sum games?

Nash equilibrium in a 2-player 0-sum normal-form game with payoff matrix $A$:

$$\max_{x \in \Delta^m} \min_{y \in \Delta^n} x^\top A y$$

✻ **IDEA: Self-play. Make two regret minimizers play each other**

for $t = 1, \dots, T$:
- $x^t \leftarrow$ request strategy from P1's regret minimizer
- $y^t \leftarrow$ request strategy from P2's regret minimizer
- Pass utility $A y^t$ to P1's regret minimizer
- Pass utility $-A^\top x^t$ to P2's regret minimizer

$$R_1^T := \max_{\hat{x} \in \Delta^m} \left\{ \sum_{t=1}^T \langle A y^t, \hat{x} \rangle \right\} - \sum_{t=1}^T \langle A y^t, x^t \rangle \leq O_m(\sqrt{T})$$

$$R_2^T := \max_{\hat{y} \in \Delta^n} \left\{ \sum_{t=1}^T \langle -A^\top x^t, \hat{y} \rangle \right\} - \sum_{t=1}^T \langle -A^\top x^t, y^t \rangle \leq O_n(\sqrt{T})$$

Add these two lines and divide by $T$ to get the average

✻ **TAKEAWAY**

**The average strategies converge to a Nash equilibrium!**

$$\max_{\hat{x} \in \Delta^m} \{ \hat{x}^\top A \bar{y} \} - \min_{\hat{y} \in \Delta^n} \{ \bar{x}^\top A \hat{y} \} \leq O_{m,n}\left( \frac{1}{\sqrt{T}} \right)$$

where $\bar{x} = \frac{1}{T} \sum_{t=1}^T x^t$ and $\bar{y} = \frac{1}{T} \sum_{t=1}^T y^t$

# Regret Minimization: Follow the Leader (Fictitious Play)

First attempt: Follow the leader. That is, play the best action in hindsight so far:

$$x^{t+1} = \arg \max_{x \in X} \sum_{\tau \leq t} \langle u^\tau, x \rangle$$

**This does not work!**

Counterexample: $n = 2$ actions,

$$u^t = \begin{cases} [1/2, 0] & t = 1 \\ [0, 1] & t > 1, \text{even} \\ [1, 0] & t > 1, \text{odd} \end{cases}$$

Best action in hindsight has utility $\approx T/2$
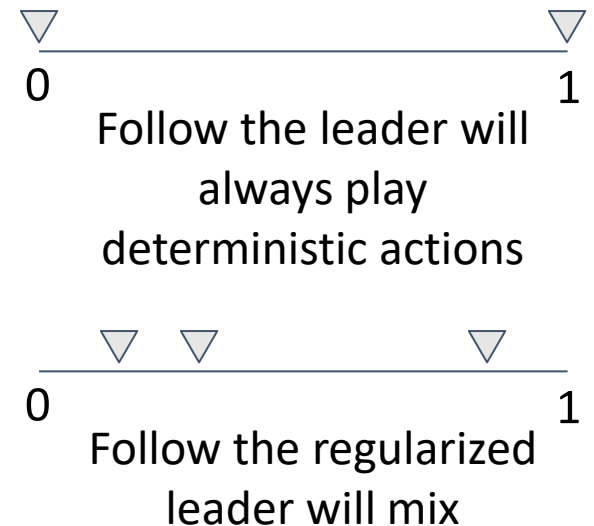Follow-the-leader always plays the wrong action and therefore gets utility $\approx 0$

*More generally: No algorithm outputting only pure actions can have no regret*

# Follow the *Regularized* Leader

**Idea**: Add a strictly convex *regularizer* $R : X \rightarrow \mathbb{R}$

$$x^{t+1} = \arg \max_{x \in X} \sum_{\tau \leq t} \langle u^\tau, x \rangle - \frac{1}{\eta} R(x)$$

- This prevents each iterate from being deterministic

- The resulting algorithm **is** no-regret (for $\eta \propto 1/\sqrt{T}$)

- Intuitively, **updates toward high-regret actions, but not too much**

Follow the leader will always play deterministic actions

Follow the regularized leader will mix

# Follow the *Regularized* Leader

**Idea**: Add a strictly convex *regularizer* $R : X \to \mathbb{R}$

$$x^{t+1} = \arg\max_{x \in X} \sum_{\tau \leq t} \langle u^\tau, x \rangle - \frac{1}{\eta} R(x)$$

**Example 1: *quadratic***

$$R(x) = \frac{1}{2} \|x\|_2^2$$

Closed-form optimization:    $\Pi_X$ = projection onto $X$

$$x^{t+1} = \Pi_X \left( \eta \cdot \sum_{\tau=1}^{t} u^\tau \right)$$

a.k.a. gradient descent!

Follow the leader will always play deterministic actions

Follow the regularized leader will mix

# Follow the *Regularized* Leader

**Idea**: Add a strictly convex *regularizer* $R : X \to \mathbb{R}$

$$x^{t+1} = \arg \max_{x \in X} \sum_{\tau \leq t} \langle u^\tau, x \rangle - \frac{1}{\eta} R(x)$$

**Example 2: *negative entropy***

$$R(x) = \sum_a x[a] \log x[a]$$

Closed-form optimization:

$$x^{t+1} \propto \exp\left( \eta \cdot \sum_{\tau=1}^{t} u^\tau \right)$$

a.k.a. multiplicative weights update (MWU), hedge, (discrete-time) replicator dynamics, randomized weighted majority, …

0 ▽ ——————————————— ▽ 1

Follow the leader will always play deterministic actions

0 ▽ ▽ ——————— ▽ — 1

Follow the regularized leader will mix

# A Common Template for Regret Minimizers

- Given utility vectors $\boldsymbol{u}^1, \ldots, \boldsymbol{u}^t$, we compute the empirical regrets up to time t of each action:

$$r^t[a] := \sum_{\tau=1}^{t} (u^\tau[a] - \langle \boldsymbol{u}^\tau, \boldsymbol{x}^\tau \rangle)$$

- Then, intuitively the next strategy $\boldsymbol{x}^{t+1}$ gives mass to actions in a manner related to how much regret they have accumulated

# A Common Template for Regret Minimizers

Empirical regret:

$$r^t[a] := \sum_{\tau=1}^{t} (u^\tau[a] - \langle \boldsymbol{u}^\tau, \boldsymbol{x}^\tau \rangle)$$

Hyperparameter
("learning rate")

| Algorithm | Rule |
|---|---|
| Gradient descent | $\boldsymbol{x}^{t+1} = \Pi_X(\eta \cdot \boldsymbol{r}^t)$ |
| Multiplicative weights update (MWU) <br> (aka Hedge, Randomized Weighted Majority, …) | $\boldsymbol{x}^{t+1} \propto \exp(\eta \cdot \boldsymbol{r}^t)$ |
| Regret matching (RM) <br> [Hart & Mas-Collel 2000] | $\boldsymbol{x}^{t+1} \propto \max\{0, \boldsymbol{r}^t\}$ |

No learning rate.
Scale-invariant!

# RM Regret Bound Proof

$$\mathbf{x}^{t+1} \propto [\boldsymbol{r}^t]^+ \quad \text{where}$$

$$\boldsymbol{r}^t := \boldsymbol{r}^{t-1} + \boldsymbol{g}^t$$
$$\boldsymbol{g}^t := \boldsymbol{u}^t - \langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \cdot \mathbf{1}$$

*Note:* $\langle \boldsymbol{g}^t, \boldsymbol{x}^t \rangle = 0$

$$\|[\boldsymbol{r}^{t+1}]^+\|_2^2 \leq \|[\boldsymbol{r}^t]^+ + \boldsymbol{g}^{t+1}\|_2^2 \qquad \text{using inequality } [x+y]^+ \leq |[x]^+ + y| \text{ for } x, y \in \mathbb{R}$$

$$= \|[\boldsymbol{r}^t]^+\|_2^2 + \|\boldsymbol{g}^{t+1}\|_2^2 + 2(\boldsymbol{g}^{t+1})^\top [\boldsymbol{r}^t]^+ \quad \mathbf{0}$$

induction

$$\|[\boldsymbol{r}^T]^+\|_2^2 \leq \sum_{t=1}^{T} \|\boldsymbol{g}^t\|_2^2 \leq nT \qquad \text{since } \boldsymbol{g}^t \in [-1,1]^n$$

$$R^T := \max_a r^T[a] \leq \|[\boldsymbol{r}^T]^+\|_2 \leq \sqrt{nT} \qquad \square$$

# A Common Template for Regret Minimizers

**Empirical regret:**   $r^t := r^{t-1} + g^t$

**Simple modification:**  $r_+^t := [r_+^{t-1} + g^t]^+$

(Floor regrets at 0 after every iteration)

| Algorithm | Rule |
|---|---|
| Gradient descent | $x^{t+1} = \Pi_X(\eta \cdot r^t)$ |
| Multiplicative weights update (MWU)<br>(aka Hedge, aka Randomized Weighted Majority) | $x^{t+1} \propto \exp\{\eta \cdot r^t\}$ |
| Regret matching (RM)<br>[Hart & Mas-Collel 2000] | $x^{t+1} \propto [r^t]^+$ |
| Regret matching plus (RM+)<br>[Tammelin 2014] | $x^{t+1} \propto [r_+^t]^+$ |

*(Regret bound proof is identical)*

# A Common Template for Regret Minimizers

All of these algorithms guarantee that after seeing any number T of utilities $\boldsymbol{u}^1, \ldots, \boldsymbol{u}^T$, the regret cumulated by the algorithm satisfies

$$R^T \leq C \sqrt{T}$$

Constant that depends on number of actions

MWU: $C = \sqrt{\log n}$

RM, RM+, GD: $C = \sqrt{n}$

**Remember**:
This holds without any assumption about the way the utilities are selected by the environment!

**Consequence**: when using these algorithms in self-play in 2-player 0-sum games, the **average strategy** converges to a Nash equilibrium at a rate of $C / \sqrt{T}$

Reminder: Self-play

for $t = 1, \ldots, T$:
- $\boldsymbol{x}^t \leftarrow$ request strategy from P1's regret minimizer
- $\boldsymbol{y}^t \leftarrow$ request strategy from P2's regret minimizer
- Pass utility $\boldsymbol{A}\boldsymbol{y}^t$ to P1's regret minimizer
- Pass utility $-\boldsymbol{A}^\top \boldsymbol{x}^t$ to P2's regret minimizer

# State-of-the-art variant in practice: Discounted RM (DRM)

- Linear RM (LRM)
  - Weight iteration t by t (in regrets and averaging)
  - RM+ floors regrets at 0. Can we combine this with linear RM? Theory: Yes. Practice: No! Does very poorly.

- But less-aggressive combinations do well: **Discounted RM**
  - On each iteration, multiply positive regrets by $t^\alpha / (t^\alpha+1)$
  - On each iteration, multiply negative regrets by $t^\beta / (t^\beta+1)$
  - Weight contributions toward average strategy on iteration $t$ by $t^\gamma$
  - Worst-case convergence bound only a small constant worse than that of RM
  - RM: $\alpha = \beta = +\infty$
  - RM+: $\alpha = +\infty, \ \beta = -\infty$
  - For $\alpha = 1.5, \beta = 0, \gamma = 2$, consistently outperforms RM+ in practice

# What Regret Minimizers are Used in Practice?

| Follow the Regularized Leader (FTRL) (*e.g.*, gradient descent, multiplicative weights) | Regret Matching (RM) & Regret Matching+ (RM+) |
|---|---|
| ✔ Works for general convex sets | ✘ Only for **simplex** domains |
| ✔ Widely used & understood | ✘ Not as well studied theoretically |
| ✘ Slow in practice | ✔ Fast in practice |
| ✘ Has hyperparameters (stepsize) | ✔ No hyperparameters |

🍀 Modern variants of this, such as DCFR, are the standard in extensive-form game solving!

✔ Can incorporate optimism about future losses to converge faster in 2-player 0-sum games

**?** Unknown

...until recently ✔

# Optimistic (Predictive) Regret Minimizers

| Algorithm | Standard (non-optimistic) rule | Optimistitic (aka Predictive) rule |
|-----------|-------------------------------|-----------------------------------|
| GD | $\boldsymbol{x}^{t+1} = \Pi_X(\eta \cdot \boldsymbol{r}^t)$ | |
| MWU | $\boldsymbol{x}^{t+1} \propto \exp\{\eta \cdot \boldsymbol{r}^t\}$ | Replace $\boldsymbol{r}^t$ |
| RM | $\boldsymbol{x}^{t+1} \propto [\boldsymbol{r}^t]^+$ | with $\boldsymbol{r}^t + \boldsymbol{g}^t$ |
| RM+ | $\boldsymbol{x}^{t+1} \propto [\boldsymbol{r}^t_+]^+$ | |

Typically, one-line change in implementation

All of these algorithms guarantee that after seeing any number $T$ of utilities $\boldsymbol{u}^1, \ldots, \boldsymbol{u}^T$, the regret cumulated by the algorithm satisfies
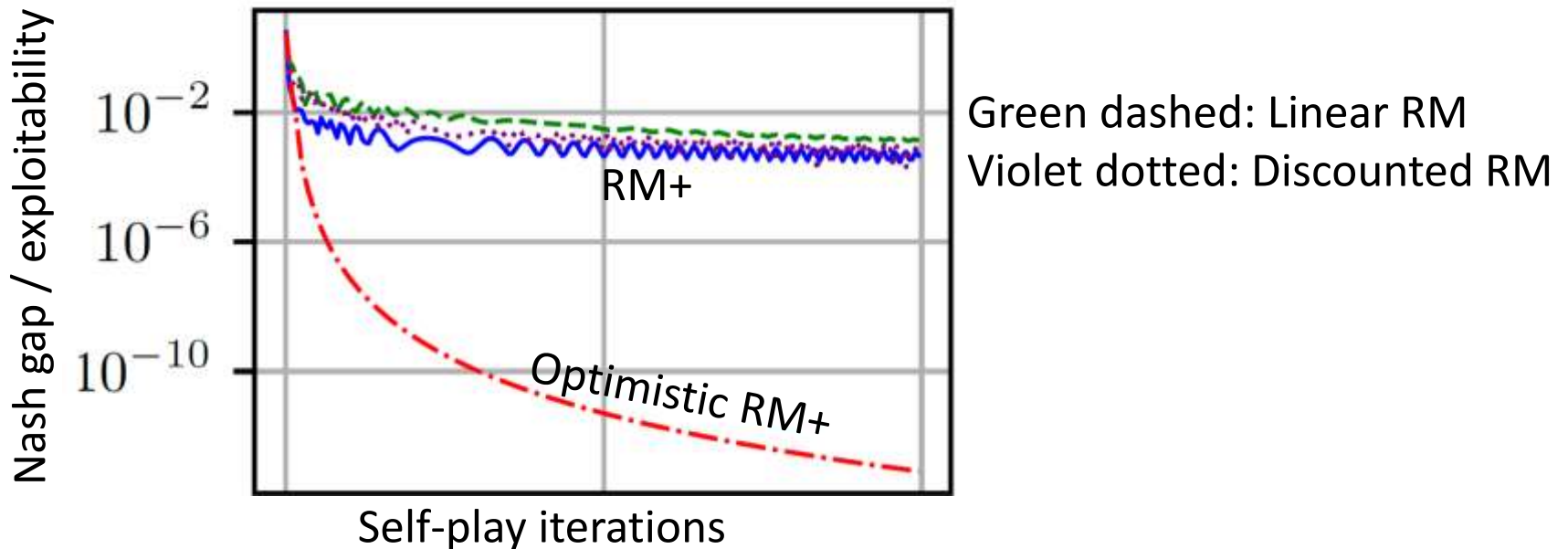
$$R^T \leq C \sqrt{\sum_{t=1}^{T} \|\boldsymbol{g}^t - \boldsymbol{g}^{t-1}\|_2^2} \qquad \text{(where } \boldsymbol{g}^0 := \boldsymbol{0}\text{)}$$

**Remember**: This holds without any assumption about the way the utilities are selected by the environment!

**Takeaway message:** still $\approx \sqrt{T}$ regret, but much smaller when there is little change to the utilities over time

# Empirical Performance



Green dashed: Linear RM
Violet dotted: Discounted RM

(RM was omitted as it is typically much slower than RM+)

# Practical State-of-the-Art

- In general, Discounted RM and Optimistic RM+ are the fastest in practice
  - For some games, like poker, Discounted RM is empirically consistently faster than Optimistic RM+
  - For many other games, Optimistic RM+ is significantly faster

[Farina, Kroer, and Sandholm; Faster Game Solving via Predictive Blackwell Approachability: Connecting Regret Matching and Mirror Descent, AAAI 2021]

# Beyond Zero-Sum Games

Correlated strategy profile:

*Note: not* $\Delta(A_1) \times \cdots \times \Delta(A_n)$

$$\mu^T := \frac{1}{T} \sum_{t=1}^{T} (x_1^t \otimes x_2^t \otimes \cdots x_n^t) \in \Delta(A_1 \times \cdots \times A_n)$$

*the product distribution in $\Delta(A_1) \times \cdots \times \Delta(A_n)$ whose marginal on $A_i$ is $x_i^t \in \Delta(A_i)$*

Regret guarantee: for all players $i$:

$$\max_{x_i^*} \frac{1}{T} \sum_{t=1}^{T} \left[ u_i(x_i^*, x_{-i}^t) - u_i(x_i^t, x_{-i}^t) \right] \leq O_n\left(\frac{1}{\sqrt{T}}\right)$$

$$= \max_{x_i^*} \underset{x \sim \mu^T}{\mathbb{E}} \left[ u_i(x_i^*, x_{-i}) - u_i(x_i, x_{-i}) \right]$$

$\mu^T$ is an $\epsilon$-"*coarse-correlated equilibrium*" (CCE) where $\epsilon = O_n(1/\sqrt{T})$

Note: A CCE that happens to be a product distribution ($\mu^T \in \Delta(A_1) \times \cdots \times \Delta(A_n)$) is a Nash equilibrium

# References

**Fictitious play:**

- J Robinson (*Ann. Math.* 1951), "An iterative method of solving a game"
- C Daskalakis, Q Pan (*FOCS* 2014), "A Counter-Example to Karlin's Strong Conjecture for Fictitious Play"
- S Karlin (1959), *Mathematical Methods and Theory in Games, Programming, and Economics*

**Blackwell Approachability (used in the original correctness proof of RM/RM+):**

- D Blackwell (*Pacific J. of Math*. 1956), "An analog of the minmax theorem for vector payoffs"

**Regret Matching and Regret Matching Plus:**

- S Hart, A Mas-Colell (*Econometrica* 2000), "A Simple Adaptive Procedure Leading to Correlated Equilibrium"
- O Tammelin (*arXiv* 2014), "Solving large imperfect information games using CFR+"
- N Brown, T Sandholm (*AAAI* 2019), "Solving Imperfect-Information Games via Discounted Regret Minimization"
- **Simple proof of correctness presented in this lecture due to** G Farina (2023), https://www.mit.edu/~gfarina/2023/6S890f23_L05_learning_algorithms/L05.pdf

**Predictivity:**

- CK Chiang et al. (*COLT* 2012), "Online optimization with gradual variations"
- A Rakhlin, K Sridharan (*COLT* 2013), "Online Learning with Predictable Sequences"
- G Farina, C Kroer, T Sandholm (*AAAI* 2021), "Faster Game Solving via Predictive Blackwell Approachability: Connecting Regret Matching and Mirror Descent"