# Solving Games

## Combinatorial Algs for Matrix Games

Lemke - Howson    (Last Week)

PNS    (Porter - Nudelman - Shoham)

## Optimization

General Formulation  $\rightarrow$  MIP

Zero - Sum Games  .

## Learning / Evolution

Fictitious Play  $\leftarrow$ Alpha Star

Regret Matching

$\quad\mathrel{\rule[-1.5ex]{0.4pt}{2ex}\!\!\rule{1ex}{0.4pt}}$ Counterfactual Regret Minimization  $\leftarrow$ Deep Stack

---

## Optimization

Alg 4 DM

$\text{minimize}_{\pi, U} \quad \sum_i \left( U^i - U^i(\pi) \right)$  $\leftarrow$ 0 when NE

$\text{subject to} \quad U^i \geq U^i(a^i, \pi^{-i})$  $\leftarrow$ Best Response  $\quad \forall i, a^i$

$\left( \begin{array}{l} \sum \pi^i(a^i) = 1 \\ \pi^i(a^i) \geq 0 \end{array} \right)$  "  "

---

## Mixed - Integer Program

$\text{maximize}_{\pi, u_a^i, u^i, b_a^i} \quad 1 \quad \overset{\leftarrow}{\phantom{.}} \quad \sum_i u_i$  social optimum

$\text{subject to} \quad \sum_a \pi^i(a) = 1 \quad \forall i, a$  $\Big]$ $\pi^i$ is prob

$\qquad\qquad \pi^i(a) \geq 0$

linear equations that encode indifference  $\longrightarrow \boxed{u_a^i = \sum_{a^{-i}} \pi^{-i}(a^{-i}) U^i(a, a^{-i})}$

regret  $\longrightarrow$  $u^i \geq u_a^i$

$\qquad\qquad r_a^i = u^i - u_a^i$  $\Big\}$ defining regret

$\qquad\qquad \pi^i(a) \leq 1 - b_a^i$  $\leftarrow$ if inactive $\pi^i(a) = 0$

$\qquad\qquad r_a^i \leq U_i \, b_a^i$  $\leftarrow$ if active, regret is zero  $\Big\}$ defining which actions are active in NE

$U_i \equiv$ max difference between utilities for player $i$

$\qquad\qquad b_a^i \in \{0, 1\}$

$\qquad\qquad\qquad \overset{\uparrow}{\text{active}} \, \overset{\uparrow}{\text{inactive}}$

Only NE are feasible solutions

Can change objective to find NE we want, e.g. social optimum, fairest

Outperforms PNS when supports are large

$$\underset{x,y,u,v}{\text{minimize}} \quad u - \cancel{x^T A y} + v - \cancel{x^T B y}$$

$$\text{subject to} \quad u \geq \hat{x}^T A y \qquad \forall \text{ one-hot vectors } \hat{x}$$
$$v \geq x^T B \hat{y} \qquad \forall \qquad '' \qquad '' \qquad \hat{y}$$
$$\sum_i x_i = 1$$
$$x \geq 0$$
$$\sum_i y_i = 1$$
$$y \geq 0$$

$$\begin{array}{cc} x & y \\ A & B \end{array}$$

<span style="color:red">Zero sum $B = -A$</span>

<span style="color:blue">$$\underset{x,v}{\text{minimize}} \quad v$$
$$\text{such that} \quad v \geq x^T B \hat{y}$$
$$\sum_i x_i = 1$$
$$x \geq 0$$</span>

<span style="color:red">Linear Program</span>

---

# Learning Methods

## Fictitious Play

Initialize $N^i(a) \leftarrow 0 \qquad \forall i, a$

for $t \in 1..T$
$\quad \pi_t^i \leftarrow \text{normalize}(N^i) \qquad \forall i$
$\quad a^i \leftarrow \text{best-response}(\pi^{-i})$
$\quad N^i(a^i) \leftarrow N^i(a^i) + 1$

return $\pi_T$

Does not always converge
⑥ Guaranteed to converge in 2 player games if
− Constant sum
− nondegenerate $2 \times n$ — 2005
− potential game
− solvable by iterated elimination of strictly-dominated strategies

| 0,0 | 2,1 | 1,2 |
|-----|-----|-----|
| 1,2 | 0,0 | 2,1 |
| 2,1 | 1,2 | 0,0 |



opponent model     policy

agent 1 $P(\text{action})$

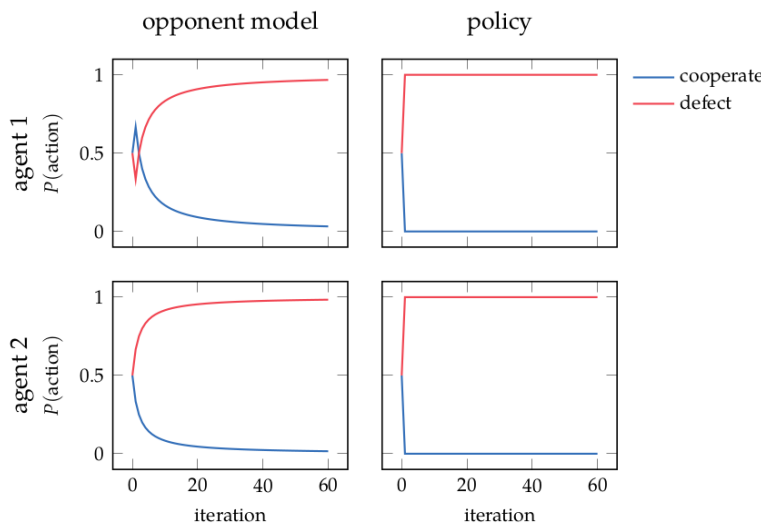agent 2 $P(\text{action})$

cooperate
defect

iteration

Figure 24.2. Two fictitious play agents learning and adapting to one another in a prisoner's dilemma game. The first row illustrates agent 1's learned model of 2 (left) and agent 1's policy (right) over iteration. The second row follows the same pattern, but for agent 2. To illustrate variation in learning behavior, the initial counts for each agent's model over the other agent's action were assigned to a random number between 1 and 10.



opponent model     policy

agent 1 $P(\text{action})$

agent 2 $P(\text{action})$
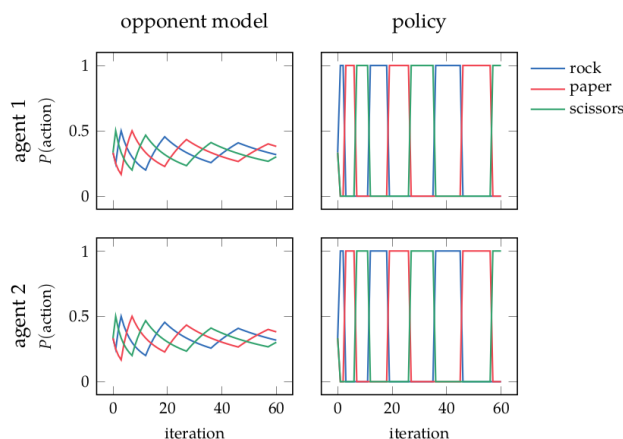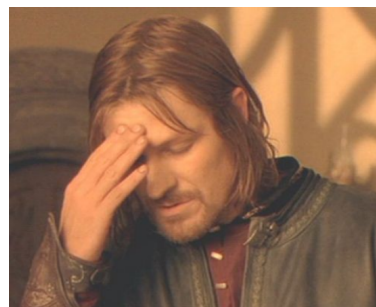
rock
paper
scissors

iteration

Figure 24.3. A visualization of two fictitious play agents learning and adapting to one another in a rock-paper-scissors game. The first row illustrates agent 1's learned model of 2 (left) and agent 1's policy (right) over time. The second row follows the same pattern, but for agent 2. To illustrate variation in learning behavior, the initial counts for each agent's model over the other agent's action were assigned to a random number between 1 and 10. In this zero-sum game, fictitious play agents approach convergence to their stochastic policy Nash equilibrium.

# Regret

How much better one could have done by taking another action

$$R^i(a^i, \pi) = U^i(a^i, \pi^{-i}) - U^i(\pi)$$

Odd

| Even | H | T |
|---|---|---|
| H | 1,-1 | -1,1 |
| T | -1,1 | 1,-1 |

$\pi^1 = [0.5, 0.5]$
$\pi^2 = [0.5, 0.5]$

$R^1(H, \pi) = (0.5 \cdot 1 + 0.5 \cdot -1) - (0.25 \cdot 1 + 0.25 \cdot -1 + 0.25 \cdot 1 + 0.25 \cdot -1) = 0$
$R^1(T, \pi) = 0$
$R^2(H, \pi) = 0$
$R^2(T, \pi) = 0$

$\pi^1 = [1, 0]$
$\pi^2 = [1, 0]$

$R^1(H, \pi) = (1) - (1) = 0$
$R^1(T, \pi) = (-1) - (1) = -2$
$R^2(H, \pi) = 0$
$R^2(T, \pi) = (1) - (-1) = 2$

## Regret Matching    a.k.a. Blackwell's algorithm

strat. prof. sum          cumulative regret

Initialize   $\bar{\pi}_0^i \leftarrow [0, 0 0 ..]$      $\bar{R}_0^i \leftarrow [0, 0, 0 ....]$      $\forall i$

for $t \in 1..T$

does not converge to Nash →

$\pi_t^i \leftarrow \text{normalize}(\bar{R}_t^i)$    $\forall i$

$\bar{\pi}_t \leftarrow \bar{\pi}_{t-1} + \pi_t$    ← if argument is 0, return uniform

$R_t^i(a) \leftarrow \max(U^i(a, \pi_t^{-i}) - U^i(\pi_t), 0)$      $\forall i, a$

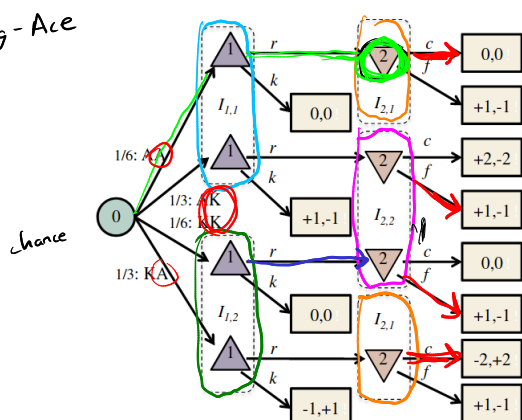$\bar{R}_t^i \leftarrow \bar{R}_{t-1}^i + R_t^i$      $\forall i$

return $\text{normalize}(\bar{\pi}_t)$

← converges for some classes    e.g. zero sum

## Extensive - Form Games

King-Ace



chance

$h \equiv$ history : sequence of actions
node in tree

$I \equiv$ information set
histories in the same $I$ are indistinguishable to that player

policy : mapping from each information set to a distribution over actions

Evaluating strategies

$P_\pi(h) \equiv$ prob. of reaching $h$ under $\pi$

$U(\pi) = \sum_{h \in Z} U(h) P_\pi(h)$    ← terminal

|      | cc   | cf    | ff  | fc  |
|------|------|-------|-----|-----|
| r r  | 0    | $-\frac{1}{6}$ | 1   | $\frac{7}{6}$ |
| b r  | $-\frac{1}{3}$ | $-\frac{1}{6}$ | $\frac{5}{6}$ | $\frac{2}{3}$ |
| r b  | $\frac{1}{3}$ | 0 | $\frac{1}{6}$ | $\frac{1}{2}$ |
| b b  | 0    | 0     | 0   | 0   |

## Counterfactual Regret Minimization (CFR)

Key Idea: break overall regret into terms for each $I$ that can be added together to bound overall regret

### Counterfactual Utility

$$U^i(\pi, I) = \frac{\sum_{h \in I, h' \in Z} P_\pi^{-i}(h) P_\pi(h, h') U^i(h')}{P_\pi^i(I)}$$

$P_\pi(h, h') \equiv$ probability of going from $h$ to $h'$ under $\pi$

$P_\pi^{-i}(h) \equiv$ probability of reaching $h$ if $i$ deliberately tries to reach $h$, everyone else plays according to $\pi$

### Counterfactual Regret

$$R_\pi^i(I, a) = P_\pi^i(I) \left( U^i(\pi|_{I \to a}, I) - U^i(\pi, I) \right)$$

$\pi|_{I \to a} \equiv$ play w/ policy $\pi$ except at $I$, take $a$

$$\bar{R}_{t, imm}^i(I) = \frac{1}{t} \max_{a \in A(I)} \sum_{\tau=1}^{t} R_{\pi_\tau}^i(I, a)$$

$$\bar{R}_{t, imm}^{i+} = \max\left( \bar{R}_{t, imm}^i(I), 0 \right)$$

**Theorem** $\quad \bar{R}_t^i \leq \sum_I \bar{R}_{t, imm}^{i+}(I) \qquad \left[ \text{Zinkevich et al. 07} \right]$

CFR algorithm: Apply regret matching at each $I$ with regret $R_\pi^i(I, a)$

In 2007, abstracted limit Texas Hold em

$10^{18}$ game states $\longrightarrow 10^{12}$ in abstraction

Poker-specific optimization: 18.5k reachable states
6.5k reachable info sets

750 iterations / sec on single core