

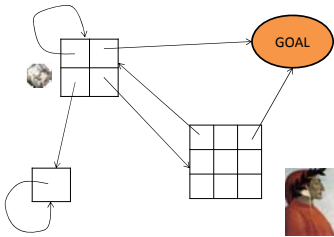
## Recent Results on Concurrent Reachability Games

Peter Bro Miltersen  
Aarhus University

**CTIC**  
交互计算  
ctic.au.dk

RP'15 1

### Concurrent Reachability Game (CRG)

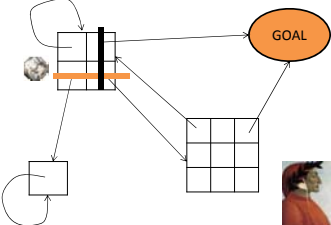


Row player wants pebble to reach GOAL

Column player wants to prevent pebble from reaching GOAL

RP'15 2

### Concurrent Reachability Game (CRG)

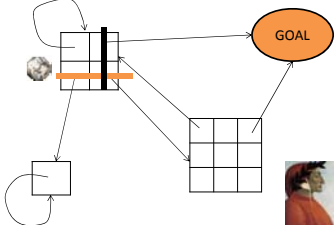


Row player wants pebble to reach GOAL

Column player wants to prevent pebble from reaching GOAL

RP'15 3

### Concurrent Reachability Game (CRG)

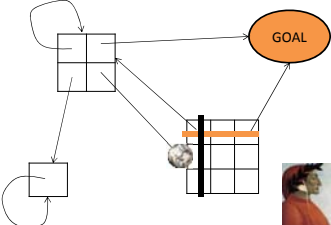


Row player wants pebble to reach GOAL

Column player wants to prevent pebble from reaching GOAL

RP'15 4

### Concurrent Reachability Game (CRG)

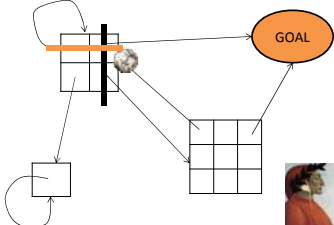


Row player wants pebble to reach GOAL

Column player wants to prevent pebble from reaching GOAL

RP'15 5

### Concurrent Reachability Game (CRG)



Row player wants pebble to reach GOAL

Column player wants to prevent pebble from reaching GOAL

RP'15 6

## Concurrent Reachability Game

- Arena:
  - Finite directed graph.
  - Distinguished terminal **GOAL** node.
  - Each non-goal node contains **a matrix** of outgoing arcs.
    - Let us also allow probabilistic transitions
- Play:
  - A pebble moves from position to position.
  - In each step, Row player chooses a row and Column player **simultaneously** chooses a column of the matrix.
  - The pebble moves along the appropriate arc.
  - If pebble reaches **GOAL**, Row player wins.
  - If this **never** happens, Column player wins.

RP'15

7

## Why study these games?

- Common generalization of
  - Turn-based reachability games.
  - Matrix games
  - Parity games
  - Shapley's discounted stochastic games.
- de Alfaro, Henzinger, Kupferman ('98+'07):
  - *Concurrent games capture the interaction of a system with its environment: in many concurrency models, in each state, both the system and the environment can independently propose moves (input or output signals), and the parallel execution of the moves determines the next state. Concurrent games also provide a natural model for distributed systems in which the moves are not revealed until their combined effect (the state transition) is apparent.*
- Natural model for **Poker tournaments**.

RP'15

8

Miltersen & Sørensen, AAMAS'07.  
My most **downloaded** paper according to ResearchGate (NOT most cited).  
(Download rate exceeds combined rate of all other papers)

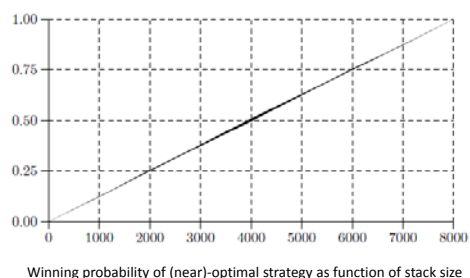
### A Near-Optimal Strategy for a Heads-Up No-Limit Texas Hold'em Poker Tournament

Peter Bro Miltersen  
University of Aarhus  
Åbogaards 34, Århus, Denmark  
bromille@daimi.au.dk

Troels Ejner Sørensen  
University of Aarhus  
Åbogaards 34, Århus, Denmark  
trold@daimi.au.dk

**ABSTRACT**  
We analyze a heads-up no-limit Texas Hold'em poker tournament with a fixed small blind of 300 chips, a fixed big blind of 600 chips and a total amount of 8000 chips on the table (until recently, these parameters defined the heads-up endgame of sit-n-go tournaments on the popular PartyPoker.com online poker site). Due to the size of this game, a computation of an optimal (i.e. minimax) strategy for the game is completely infeasible. However, combining an algorithm due to Koller, Megiddo and von Stengel with concepts of Everett and suggestions of Sklansky, we compute an optimal *push/fold* strategy, i.e. a strategy that would be optimal if our bet made by the player chosen by the strategy that the computed strategy. These exact parameters were chosen as they until recently defined the heads-up endgame of the \$5 through \$30 buy-in sit-n-go tournaments on the popular PartyPoker.com online poker site.  
We briefly review the rules of such a tournament (and of no-limit Texas Hold'em in general). The tournament is played between two players, Player 1 and Player 2. When the tournament starts, Player 1 receives  $s_1$  chips and Player 2 receives  $s_2$  chips where  $s_1 + s_2 = 8000$ . We want to keep  $s_1$  and  $s_2$  as parameters as the heads-up tournament we consider may be the endgame of a tournament with more people, so the two players should be able to enter the game with different stack sizes. The total number of chips on the

## PartyPoker heads-up Texas Hold'em tournament end-game



RP'15

10

## Values and Near-Optimal Strategies

- **Everett'57**: Each position  $i$  in a CRG has a **value**  $v_i$  so that

$$v_i = \min_{\text{stationary } y} \max_{\text{general } x} \mu_i(x, y) \\ = \sup_{\text{stationary } x} \min_{\text{general } y} \mu_i(x, y)$$

where  $\mu_i(x, y)$  is the probability of reaching GOAL when row player plays by strategy  $x$  and column player plays by strategy  $y$ .

RP'15

11

## Why sup instead of max?

- **Everett'57, Example 1**:
    - Column player hides a penny.
    - If Row player can guess if it is **heads** up or **tails** up, he gets the penny.
    - If he incorrectly guesses **heads**, play repeats.
    - If he incorrectly guesses **tails**, the game ends.
- 
- guarantees reaching GOAL with probability  $1 - \epsilon$ , so value of start position is 1.

RP'15

12

## The problem(s) we study in this talk

- **Compute the value of a given game**
  - Caveat: The value might be irrational
    - Exact decision problem: Compare value to given rational number
      - Interesting special case: Value-1 problem.
    - Compute (additive) approximation to the value
- **Synthesize good strategies**
  - Given a game and  $\epsilon$ , synthesize stationary strategy guaranteeing the value within  $\epsilon$ .

RP'15

13

## Classical results

- The value-1 and associated synthesis problem can be solved *in polynomial time*
  - [de Alfaro, Henzinger, Kupferman '98]
- The general problems can be solved *in polynomial space*
  - [Etessami & Yannakakis'05]
- A practically applicable *strategy improvement algorithm* for the row player
  - [de Alfaro, Chatterjee, Henzinger '06]
  - We used this algorithm to solve the poker tournament..
- **No "standard" hardness results!**
  - Non-standard hardness: SQRT-SUM hardness [EY'05], "hard for solving parity games".

RP'15

14

## PSPACE upper bound

- Given a game, we solve it as follows:
  1. Write down first order formulae defining the real numbers forming the solution.
  2. Use the best available decision procedure for the first order theory of the real numbers to turn the definitions into actual answers.



RP'15

15

## The recent results

- A worst case time complexity analysis of the strategy improvement algorithm of de Alfaro *et al.*
  - The complexity is **doubly exponential**.
- A **polynomial time** algorithm for the general problem *when the number of positions is constant*.
- A **polynomial hierarchy** algorithm for the approximation version of the general problem.
- A structural theorem on near-optimal strategies: "Monomial" strategies suffice.



RP'15

16

## References

- Kristoffer Arnsfelt Hansen, Rasmus Ibsen-Jensen, and Peter Bro Miltersen. **The complexity of solving reachability games using value and strategy iteration.** In *Proceedings of CSR'11*, volume 6651 of LNCS, pages 77–90.
- Kristoffer Arnsfelt Hansen, Michal Koucký, Niels Lauritzen, Peter Bro Miltersen, and Elias P. Tsigaridas. **Exact algorithms for solving stochastic games:** extended abstract. In *Proceedings of STOC'11*, pages 205–214.
- Søren Kristoffer Stiil Frederiksen and Peter Bro Miltersen. **Monomial strategies for concurrent reachability games and other stochastic games.** In *Proceedings of RP'13*, volume 8169 of LNCS, pages 122–134.
- Søren Kristoffer Stiil Frederiksen and Peter Bro Miltersen. **Approximating the value of a concurrent reachability game in the polynomial time hierarchy.** In *Proceedings of ISAAC'13*, Lecture Notes in Computer Science, volume 8283, pages 457–467, 2013.
- Kristoffer Arnsfelt Hansen, Søren Kristoffer Stiil Frederiksen, Tongyang Li, Peter Bro Miltersen. **Meticulous algebraic sampling theorems with applications to concurrent reachability games,** in preparation.

RP'15

17

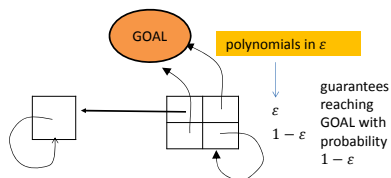
## The recent results

- A worst case time complexity analysis of the strategy improvement algorithm of de Alfaro *et al.*
  - The complexity is **doubly exponential**.
- A **polynomial time** algorithm for the general problem *when the number of positions is constant*.
- A **polynomial hierarchy** algorithm for the approximation version of the general problem.
- A structural theorem on near-optimal strategies: "Monomial" strategies suffice.

RP'15

18

## Why sup instead of max



The value of this game is 1.

RP'15

19

## A variant




- A variant:
  - Bob hides a die.
  - If Alice can guess its outcome  $x$ , and  $x$  is in 3,4,5,6, she gets a cookie.
  - If Alice can guess its outcome  $x$ , and  $x$  is in 1,2, she gets a cookie with probability  $x/2$ .
  - If she incorrectly guesses 3,4,5 or 6, play repeats.
  - If she incorrectly guesses 1 or 2, play ends.
- The value for Alice is  $\frac{2}{3}$  cookie.
- An  $\epsilon$ -optimal strategy for Alice:
  - Guess 3,4,5 each with probability  $\frac{1}{4}$
  - Guess 1 with probability  $\frac{2}{3}\epsilon$
  - Guess 2 with probability  $\frac{1}{3}\epsilon$

polynomials in  $\epsilon$

RP'15

20

## Two-level Purgatory

- Two-level Purgatory
  - Bob repeatedly hides a penny. 
  - If Alice can guess correctly if it is **heads** up or **tails** up **twice in a row**, she gets the penny.
  - If she ever incorrectly guesses **tails**, the game ends.
- The value for Alice is 1 penny.
- An  $\epsilon$ -optimal strategy for Alice:
  - Guess tails with probability  $\epsilon^2$  in initial state.
  - Guess tails with probability  $\epsilon^4$  after having guessed correctly once.

polynomials in  $\epsilon$

21

## Do "simple" strategies always suffice?

- In the preceding examples, the  $\epsilon$ -optimal strategies had probabilities expressed as (simple) polynomials in  $\epsilon$ .
- **Is this the case for all concurrent reachability games?**
- Classical result: de Alfaro, Henzinger, Kupferman '98: YES for games of value 1.
- Recent result [FM'13]: YES for all games.

RP'15

22

## Monomial strategies

**Definition 1.** A family of stationary strategies  $(x_k)_{0 \leq k \leq s_0}$  for Player I in a stochastic game is called *monomial* if for all states  $k$ , and all actions  $j$  available to Player I in state  $k$  except possibly one action, we have that  $x_{k,j}$  is given by a monomial in  $\epsilon$ , i.e., an expression of the form  $c_j^k \epsilon^{d_j^k}$ , where  $d_j^k$  is a non-negative integer and  $c_j^k$  is a non-negative real number.

**Theorem 1.** For any game  $G$ , there is an  $\epsilon_0 > 0$  and a monomial family of stationary strategies  $(x_k)_{0 \leq k \leq s_0}$  for Player I, so that for each  $\epsilon \in (0, \epsilon_0]$ , we have that  $x_k$  is  $\epsilon$ -optimal among stationary strategies.

Motivation: Opens up possibility of *symbolic* algorithm for finding near-optimal strategies for recursive games.

RP'15

23

## Proof



- Let  $v$  be the supremum of payoffs to Alice guaranteed by her stationary strategies.
- By definition of supremum, for every  $\delta > 0$  there is a stationary strategy guaranteeing at least  $v - \delta$  (\*)
- (\*) can be formalized in the first order theory of the reals.
- By the **Tarski Transfer Principle**, (\*) is also true if  $\mathbb{R}$  is replaced with  $\mathbb{R}(\langle \epsilon \rangle)$ : The **real closed field** of Puiseux series that converge for sufficiently small  $\epsilon > 0$ .
  - A Puiseux series:  $\sum_i a_i \epsilon^{i/K}$ .
- In particular, (\*) is true for  $\delta = \epsilon$ . Let the corresponding strategy be  $x$ .
- The probabilities of  $x$  are Puiseux series in  $\epsilon$ . For sufficiently small  $\epsilon$  they converge and describe an  $\epsilon$ -optimal strategy.
- Using perturbation theory for Markov chains (Solan'03), the Puiseux series can be "rounded" to their most significant term while preserving near-optimality.
- If the exponents in the leading terms is of the form  $i/K$ , we immediately have a monomial  $\epsilon^k$ -optimal strategy.
- This strategy is also  $\epsilon$ -optimal.

RP'15

24

### Open Problem

- Is there an elementary and constructive (in some sense) proof of this theorem?
  - Motivation: leads to algorithm?

RP'15

25

### The recent results

- A worst case time complexity analysis of the strategy improvement algorithm of de Alfaro *et al.*
  - The complexity is **doubly exponential**.
- A **polynomial time** algorithm for the general problem *when the number of positions is constant*.
- A **polynomial hierarchy** algorithm for the approximation version of the general problem.
- A structural theorem on near-optimal strategies: "Monomial" strategies suffice.

RP'15

26

### Strategy iteration for CRGs

Chatterjee, de Alfaro, Henzinger '06

```

1:  $t := 1$ 
2:  $x^t :=$  the uniform distribution at each position
3: while true do
4:    $y^t :=$  an optimal best reply to  $x^t$ ;
5:   for  $i \in \{0, 1, 2, \dots, N, N + 1\}$  do
6:      $v_i^t := \mu_i(x^t, y^t)$ 
7:   end for
8:    $t := t + 1$ 
9:   for  $i \in \{1, 2, \dots, N\}$  do
10:    if  $\text{val}(A_i(v^{t-1})) > v_i^{t-1}$  then
11:       $x_i^t := \text{maximin}(A_i(v^{t-1}))$ 
12:    else
13:       $x_i^t := x_i^{t-1}$ 
14:    end if
15:  end for
16: end while
    
```

Solve Markov Decision Process

Solve matrix game

RP'15

27

### Properties

- The valuations  $v_i^t$  converge to the values  $v_i$  (from below).
- The strategies  $x^t$  guarantee the valuations  $v_i^t$  for row player.
- **What is the number of iterations required to guarantee a good approximation?**

RP'15

28

### [HKLMT'11, HIM'11]

For all games with  $N$  positions and  $m$  actions for each player in each position,  $(1/\epsilon)^{m^{O(N)}}$  iterations is sufficient to arrive at  $\epsilon$ -optimal strategy.

- Proof relies heavily on R.A.G.
- Reliance on R.A.G. means that I'm not quite sure what the constant in the big-O is.
- For some games ("Purgatory games"),  $(1/\epsilon)^{m^{N-O(N)}}$  iterations are necessary to get  $\epsilon$ -optimal strategy.

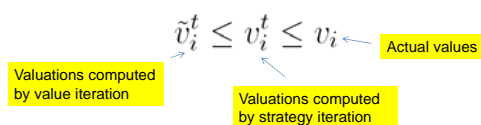


RP'15

29

### Step 1: Reduction to analysis of value iteration

- We can relate the valuations computed by strategy iteration to the valuations computed by **value iteration**.



RP'15

30

### Value iteration (dynamic programming)

```

1: t := 0
2: v0 := (0, 0, ..., 1) {the vector v0 is indexed 0, 1, ..., N, N + 1}
3: while true do
4:   t := t + 1
5:   v0t := 0
6:   vN+1t := 1
7:   for i ∈ {1, 2, ..., N} do
8:     vit := val(Ai(vt-1))
9:   end for
10: end while
    
```

Value iteration computes the value of a time bounded game, for larger and larger values of the time bound t, by **backward induction**.

### Step 2: Reduction to bounding patience

- We need to upper bound the difference in value between **time bounded** and **infinite** versions of the game.
- The difference in value between a time bounded and the infinite version of a concurrent reachability game is captured by the **patience** of its stationary near-optimal strategies.
  - Patience = 1/smallest non-zero probability used

**Lemma:** If the game has an  $\epsilon$ -optimal strategy with patience L, then for  $T = kNL^N$ , the value of the game with time bound T differs from the value of the original game by at most  $\epsilon + e^{-k}$ .

### Step 3: Bounding patience using R.A.G.

- Everett's characterization (1957) of value and near-optimal strategies:

Given valuations  $v_1, \dots, v_N$  for the positions and a given position k we define  $A^k(v)$  to be the  $m_k \times n_k$  matrix game where entry  $(i, j)$  is  $A_{ij}^k + \sum_{l=1}^{m_k} v_l p_{ij}^k$ . The value mapping operator  $M: \mathbb{R}^N \rightarrow \mathbb{R}^N$  is then defined by  $M(v) = (\text{val}(A^k(v)), \dots, \text{val}(A^N(v)))$ . Define relations  $\succ$  and  $\prec$  on  $\mathbb{R}^N$  as follows:

$$\begin{aligned}
 u \succ v & \text{ if and only if } \begin{cases} u_i > v_i & \text{if } v_i > 0 \\ u_i \geq v_i & \text{if } v_i \leq 0 \end{cases} \text{ for all } i. \\
 u \prec v & \text{ if and only if } \begin{cases} u_i < v_i & \text{if } v_i < 0 \\ u_i \leq v_i & \text{if } v_i \geq 0 \end{cases} \text{ for all } i.
 \end{aligned}$$

Next, we define the regions  $C_1(\Gamma)$  and  $C_2(\Gamma)$  as follows:

$$\begin{aligned}
 C_1(\Gamma) &= \{v \in \mathbb{R}^N \mid M(v) \succ v\}. \\
 C_2(\Gamma) &= \{v \in \mathbb{R}^N \mid M(v) \prec v\}.
 \end{aligned}$$

A **critical vector** of the game is a vector  $v$  such that  $v \in \overline{C_1(\Gamma) \cap C_2(\Gamma)}$ . That is, for every  $\epsilon > 0$  there exists vectors  $v_1 \in C_1(\Gamma)$  and  $v_2 \in C_2(\Gamma)$  such that  $\|v - v_1\|_1 \leq \epsilon$  and  $\|v - v_2\|_1 \leq \epsilon$ .

The following theorem of Everett characterizes the value of an Everett game and exhibits near-optimal strategies.

**Theorem 5 (Everett).** There exists a unique critical vector  $v$  for the value mapping  $M$ , and this is the value vector of  $\Gamma$ . Furthermore,  $v$  is a fixed point of the value mapping, and if  $v_1 \in C_1(\Gamma)$  and  $v_2 \in C_2(\Gamma)$  then  $v_1 \leq v \leq v_2$ . Let  $x$  be the stationary strategy for player I, where in position k an optimal strategy in the matrix game  $A^k(v_1)$  is played. Then for any k, starting play in position k, the strategy  $x$  guarantees expected payoff at least  $v_{1,k}$  for player I. The analogous statement holds for  $v_2 \in C_2(\Gamma)$  and Player II.

### Step 3: Bounding patience using R.A.G.

- Applying the fundamental theorem of linear programming and Cramer's rule:

Now we can rewrite the predicate  $\text{val}(A^k(v_1)) > v_{1k}$  to the following expression:  $\forall_{Bk} ((v_1 \in F_{Bk}^{A^k} + \Lambda \det((M_{Bk}^{A^k}(v_1))_{m_k+1})) > v_{1k} \det(M_{Bk}^{A^k}(v_1))) \vee ((v_1 \in F_{Bk}^{A^k} + \Lambda \det((M_{Bk}^{A^k}(v_1))_{m_k+1})) < v_{1k} \det(M_{Bk}^{A^k}(v_1)))$ , where the disjunction is over all potential basis sets, and each of the expressions  $v_1 \in F_{Bk}^{A^k}$  and  $v_1 \in F_{Bk}^{A^k}$  are shorthands for the conjunction of the  $m_k + 1$  polynomial inequalities describing the corresponding sets.

**Lemma 40.** There is a quantifier free formula with  $2N$  free variables  $v_1$  and  $v_2$  that expresses  $v_1 \in C_1(\Gamma)$ ,  $v_2 \in C_2(\Gamma)$ , and  $\|v_1 - v_2\|_1 \leq 2^{-m}$ . The formula uses at most  $(2N + 1) + 2(m + 2) \sum_{k=1}^N (m_k + m_k)$  different polynomials, each of degree at most  $m + 2$  and having coefficients of bitsize at most  $\max(\sigma, 2(N + 1)(m + 2))$ .

### Step 3: Bounding patience using R.A.G.

- Sampling theorem



**Theorem 13.10.** Let  $\mathcal{P}$  be a set of polynomials each of degree at most  $d$  in  $k$  variables with coefficients in a real closed field  $\mathbb{R}$ . Let  $D$  be the ring generated by the coefficients of  $\mathcal{P}$ . There is an algorithm that computes a set of  $2^{\sum_{j=1}^k \binom{d}{j}}$  points meeting every semi-algebraically connected component of the realization of every suitable sign condition on  $\mathcal{P}$  in  $\mathbb{R}(c, d)^k$ , described by univariate representations of degree bounded by  $(2d + 6)(2d + 5)^{k-1}$ .

The algorithm has complexity  $\sum_{j=1}^k \binom{d}{j} 4^j d^{O(k)}$  in  $D$ . There is also an algorithm computing the signs of all the polynomials in  $\mathcal{P}$  at each of these points with complexity  $\sum_{j=1}^k \binom{d}{j} 4^j d^{O(k)}$  in  $D$ .

If the polynomials in  $\mathcal{P}$  have coefficients in  $\mathbb{Z}$  of bitsize at most  $\tau$ , the bitsize of the coefficients of these univariate representations is bounded by  $\tau d^{O(k)}$ .

### Step 3: Bounding patience using R.A.G.

**Lemma 40.** There is a quantifier free formula with  $2N$  free variables  $v_1$  and  $v_2$  that expresses  $v_1 \in C_1(\Gamma)$ ,  $v_2 \in C_2(\Gamma)$ , and  $\|v_1 - v_2\|_1 \leq 2^{-m}$ . The formula uses at most  $(2N + 1) + 2(m + 2) \sum_{k=1}^N (m_k + m_k)$  different polynomials, each of degree at most  $m + 2$  and having coefficients of bitsize at most  $\max(\sigma, 2(N + 1)(m + 2))$ .

**Theorem 13.10.** Let  $\mathcal{P}$  be a set of polynomials each of degree at most  $d$  in  $k$  variables with coefficients in a real closed field  $\mathbb{R}$ . Let  $D$  be the ring generated by the coefficients of  $\mathcal{P}$ . There is an algorithm that computes a set of  $2^{\sum_{j=1}^k \binom{d}{j}}$  points meeting every semi-algebraically connected component of the realization of every suitable sign condition on  $\mathcal{P}$  in  $\mathbb{R}(c, d)^k$ , described by univariate representations of degree bounded by  $(2d + 6)(2d + 5)^{k-1}$ .

The algorithm has complexity  $\sum_{j=1}^k \binom{d}{j} 4^j d^{O(k)}$  in  $D$ . There is also an algorithm computing the signs of all the polynomials in  $\mathcal{P}$  at each of these points with complexity  $\sum_{j=1}^k \binom{d}{j} 4^j d^{O(k)}$  in  $D$ . If the polynomials in  $\mathcal{P}$  have coefficients in  $\mathbb{Z}$  of bitsize at most  $\tau$ , the bitsize of the coefficients of these univariate representations is bounded by  $\tau d^{O(k)}$ .

Limitation of consumer perspective!

+ separation bounds for roots of univariate polynomials

An  $\epsilon$ -optimal strategy with all probabilities either 0 or bounded from below by  $\epsilon^{m \cdot O(N)}$

[HKLMT'11, HIM'11]

For all games with  $N$  positions and  $m$  actions for each player in each position,  $(1/\epsilon)^{m^{O(N)}}$  iterations is sufficient to arrive at  $\epsilon$ -optimal strategy.

- Proof relies heavily on R.A.G.
- Reliance on R.A.G. means that *I'm not quite sure what the constant in the big-O is.*
- For some games ("Purgatory games"),  $(1/\epsilon)^{m^{N-o(N)}}$  iterations are necessary to get  $\epsilon$ -optimal strategy.

RP'15

37

The tight example

Generalized Purgatory  $P(N,m)$ :

- Column player repeatedly hides a number in  $\{1, \dots, m\}$ .
- Row player must try to guess the number.
- If he guesses correctly  $N$  times in a row, he wins the game.
- If he ever guesses incorrectly **overshooting** hidden number, he loses the game.
- These games all have value 1(!)
- Strategy iteration needs  $(1/\epsilon)^{m^{N-o(N)}}$  to get  $\epsilon$ -optimal strategy.

RP'15

38

Howard's algorithm is slow on  $P(7,2)$

#iterations:	Valuation of start position:
1	0.01347
10	0.03542
100	0.06879
1000	0.10207
10000	0.13396
100000	0.16461
1000000	0.19415
10000000	0.22263
100000000	0.24828
$> 2 \cdot 10^{65}$	0.9
$> 10^{128}$	0.99

RP'15

39

[HKLMT'11, HIM'11]

For all games with  $N$  positions and  $m$  actions for each player in each position,  $(1/\epsilon)^{m^{O(N)}}$  iterations is sufficient to arrive at  $\epsilon$ -optimal strategy.

- Proof relies heavily on R.A.G.
- ~~Reliance on R.A.G. means that *I'm not quite sure what the constant in the big-O is.*~~
- For some games ("Purgatory games"),  $(1/\epsilon)^{m^{N-o(N)}}$  iterations are necessary to get  $\epsilon$ -optimal strategy.

RP'15

40

[FHLM'15]

For all games with  $N$  positions and  $m$  actions for each player in each position,  $(1/\epsilon)^{m^{(4+o(1))N}}$  iterations is sufficient to arrive at  $\epsilon$ -optimal strategy.

- Proof relies heavily on R.A.G.
- ~~Reliance on R.A.G. means that *I'm not quite sure what the constant in the big-O is.*~~
- For some games ("Purgatory games"),  $(1/\epsilon)^{m^{N-o(N)}}$  iterations are necessary to get  $\epsilon$ -optimal strategy.

RP'15

41

Why does the algorithm work well "in practice"?

- The algorithm converges fast when the infinite game is approximated well by a time limited game with a reasonable time bound.
- Games you want to solve in practice (e.g., poker tournaments) tend to have this property...

RP'15

42

### The recent results

- A worst case time complexity analysis of the strategy improvement algorithm of de Alfaro *et al.*
  - The complexity is **doubly exponential**.
- A **polynomial time** algorithm for the general problem *when the number of positions is constant*.
- A **polynomial hierarchy** algorithm for the approximation version of the general problem.
- A structural theorem on near-optimal strategies: "Monomial" strategies suffice.

RP'15

43

### Exact Algorithms for Solving Stochastic Games

(Hansen, Lauritzen, Koucky, Miltersen, Tsigaridas, STOC'11)

- **Good news:** Algorithms solving stochastic games **exactly** in polynomial time when the number of positions is constant.
- **Bad news:** Complexity is something like  $\lfloor \exp(O(N \log N)) \rfloor, \dots$

RP'15

44

### Slogan of approach

- **Doing numerical analysis/optimization in dangerous waters using real algebraic geometry.**
  - The waters are dangerous because small perturbations mean everything...
  - Why?

RP'15

45

### The tight example

Generalized Purgatory  $P(N,m)$ :

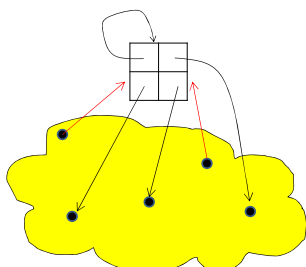
- Column player repeatedly hides a number in  $\{1, \dots, m\}$ .
- Row player must try to guess the number.
- If he guesses correctly  $N$  times in a row, he wins the game.
- If he ever guesses incorrectly **overshooting** hidden number, he loses the game.
- These games all have value 1(!)
- Strategy iteration needs  $(1/\epsilon)^{m^{N-O(N)}}$  to get  $\epsilon$ -optimal strategy.

In particular, this patience is necessary to be  $\epsilon$ -optimal!

RP'15

46

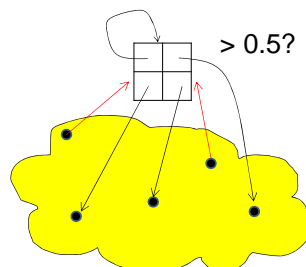
### Recursive Bisection Algorithm



RP'15

47

### Recursive Bisection Algorithm



RP'15

48



### Recursive Bisection Algorithm

1. Replace position with target value

0.5

RP'15 49

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game

0.5

0.8 0.2 0.9

RP'15 50

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinstatement position

0.5

0.8 0.2 0.9

RP'15 51

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinstatement position
4. Replace pointers

0.5

0.8 0.2 0.9

RP'15 52

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinstatement position
4. Replace pointers

0.5 0.9

0.8 0.2 0.9

RP'15 53

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinstatement position
4. Replace pointers

0.5 0.9

0.8 0.2 0.9

RP'15 54

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinststate position
4. Replace pointers

RP'15 55

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinststate position
4. Replace pointers
5. Solve matrix game

RP'15 56

### Recursive Bisection Algorithm

1. Replace position with target value
2. Recursively solve smaller game
3. Reinststate position
4. Replace pointers
5. Solve matrix game

RP'15 57

### What's the catch?

- We can compare the value of a position in an N-position game to a given rational number (and do binary search) if we recursively can solve an (N-1)-position game exactly!
- How to solve (N-1)-position game exactly using binary search?
  - 0.5 vs. 0.50000000000000000000000000000001
  - Will happen on simple examples such as Purgatory.
- To get algorithm, we must replace "exactly" with "approximately" or vice versa.

RP'15 58

### Real algebraic geometry to the rescue

- To reconcile approximately and exactly, we need **separation bounds**.
  - Separation bound: If games X and Y of certain parameters have different values, they differ by at least  $\epsilon$ .
- **Separation bounds for stochastic games using real algebraic geometry is the technical meat of the work.**

RP'15 59

### Isolated root theorem (HLKMT'11)

- Given a polynomial system:
  - $f_1(x_1, x_2, x_3, \dots, x_n) = \dots = f_m(x_1, x_2, x_3, \dots, x_n)$
- with each  $f_j$  in  $\mathbb{Q}[x_1, x_2, \dots, x_n]$ , of total degree  $d$  and with an **isolated** root  $x^*$  in  $\mathbb{R}^n$ .
- Then, the algebraic degree of each  $x^*_i$  is at most  $(2d+1)^n$ .
- Best(?) previously published bound:  $(O(d))^n$ .
- Open(?): Is  $d^n$  possible?

RP'15 60

### Open Problems

- Better exponent
- Polynomial time algorithm for finding monomial  $\epsilon$ -optimal strategy for constant number of positions and symbolic  $\epsilon$ ?

### The recent results

- A worst case time complexity analysis of the strategy improvement algorithm of de Alfaro *et al.*
  - The complexity is **doubly exponential**.
- A **polynomial time** algorithm for the general problem *when the number of positions is constant*.
- A **polynomial hierarchy** algorithm for the approximation version of the general problem.
- A structural theorem on near-optimal strategies: "Monomial" strategies suffice.

### Polynomial hierarchy?

- Not very impressive, is it?!
- Finite-state two-player zero-sum games can "usually" be solved in  $NP \cap coNP$ !
  - Parity games
  - Condon's simple stochastic games
- "Usual proof": Guess a pair of strategies and verify that they are in equilibrium.
- What is the catch?

### Values and Near-Optimal Strategies

- **Everett'57**: Each position  $i$  in a CRG has a *value*  $v_i$  so that

$$v_i = \min_{\text{stationary } y} \max_{\text{general } x} \mu_i(x, y) = \sup_{\text{stationary } x} \min_{\text{general } y} \mu_i(x, y)$$

First catch: No exact equilibrium...

No worries, instead of  $NP \cap coNP$ , we should still get TFNP upper bound for the approximation problem....

### The tight example

- Generalized Purgatory
- Column player chooses  $\alpha$ .
  - Row player chooses  $\beta$ .
  - If he guesses  $\beta$  correctly, he wins.
  - If he ever guesses wrong number, he loses.
- These games all have value  $\frac{1}{2}$ .
- Strategy iteration needs  $(1/\epsilon)^{m^{N \cdot o(N)}}$  to get  $\epsilon$ -optimal strategy.

Second catch: Exponentially many bits needed to express near-optimal strategy in fixed point notation!

In particular, this patience is necessary to be  $\epsilon$ -optimal!

### TFNP[NP] upper bound

- When fixed point notation of numbers will not do, what to go for?
- **Floating point notation!**
- Algorithm:
  - Guess stationary strategies in floating point notation.
    - Solan'03: Absorption probabilities of Markov chain change benignly when transition probabilities change with multiplicative error, so good guess is possible.
  - Guess best replies
  - Solve *rare event* absorbing Markov process using appropriate numerical algorithm: State reduction algorithm of Grassman, Takar and Heyman (1985).

### State reduction algorithm

Absorption probabilities can be approximated in polynomial time!

RP'15 67

### Open problem

- Can a rare event Markov *Decision Process* (i.e., with probabilities given in floating point) be approximately solved in polynomial time?
- This would lead to TFNP upper bound for the approximation problem, but also seems interesting in its own right..

RP'15 68

### The recent results

- A worst case time complexity of the strategy improvement algorithm.
- The complete complexity of the problem when the number of states is constant.
- A **polynomial** algorithm for the approximation version of the general problem.
- A **polynomial hierarchy** algorithm for the approximation version of the general problem.
- A structural theorem on near-optimal strategies: "Monomial" strategies suffice.

Open problem:  
Improve this!

RP'15 69

### Thank you!

RP'15 70