Guarantees for Self-Play via Polymatrix Decomposability Revan MacQueen

Thesis Presentation. August 28, 2023











Self-Play

Success of Self-Play











Background

Game Theory Basics

Normal-Form Games





Utility Functions



Nash Equilibrium



Definition. A Nash equilibrium is a **strategy profile** such that no players wish to deviate to another strategy.

Approximate Nash Equilibrium



Definition. A strategy profile is an ϵ -Nash equilibrium if no player can gain more than ϵ utility by deviating.

Coarse Correlated Equilibrium (CCE)





Coarse Correlated Equilibrium (CCE)



Each player gets expected utility of $0.5\,$

No player wishes to **deviate** from their recommendations

What If Players Lose the Ability to Correlate?



What If Players Lose the Ability to Correlate?



 S^{μ} $\mu \in \Delta(P)$ Marginal CCF strategy profile



Self-Play & No-Regret Learning









no-regret algorithm

Goal: Average regret $\rightarrow 0$





Online Interaction



We characterize the strategy produced by self-play as the marginal strategy of a CCE

Self-Play Wishlist





$$u_i(s_i^{\mu}, s_{-i}^{\mu}) - \min_{s_{-i}} u_i(s_i^{\mu}, s_{-i}) \le \epsilon$$





Approximate Nash equilibrium
















When Does Self-Play Have Guarantees?

Two Player Constant-Sum Games







2 player

 Self-play will produce an approximate Nash equilibrium Self-Play Wishlist
☑ Low vulnerability
☑ Similar values
☑ Nearly Exchangeable

What About Constant-Sum Multi-player Games?



What About Constant-Sum Multi-player Games?



<u>Self-Play Wishlist</u>
Low vulnerability
Similar values
Nearly Exchangeable









There are multi-player games that are structurally similar to two-player constant-sum games

Constant-Sum Polymatrix (CSP) Games Bregman & Fokin, 1987; Cai et al., 2011





Overall utility = **sum** of subgame utilities

- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!

G

- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!



G

- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!



- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!



- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!



- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!



Polynomial time!

- Not every game has an **exact** CSP decomposition
- But every game can be **approximated** by a CSP game!



Counterexample

Offense-Defense



Counterexample

Offense-Defense



Counterexample

Offense-Defense



Subgame Stability



Approximate Subgame Stability



(ϵ, γ) -Subgame Stability



Main Result

Guarantees For Self-Play in Multi-Player Games





Experiments

Conjecture. No limit Texas hold 'em is approximately **constant-sum polymatrix** and **subgame stable**.







Conjecture. No limit Texas hold 'em is approximately **constant-sum polymatrix** and **subgame stable**.

Conjecture. No limit Texas hold 'em is approximately **constant-sum polymatrix** and **subgame stable...**

in parts of the game space that are actually played.

Kuhn Poker

Kuhn, 1950; Abou Risk & Szafron, 2010

- 3 players
- 4 cards, 1 round of betting with fixed bet size
- CFR was previously shown to converge to Nash equilibria in Kuhn poker

Hypothesis. Kuhn poker is approximately **constantsum polymatrix** and **subgame stable...** in parts of the game space that are learnable by selfplay.



Experiment Design

 Use a self-play algorithm to generate a set of strategies for each player

• Compute a set of match-ups between all these strategies.

 Check if this set of strategy profiles is approximately subgamestable and constant-sum polymatrix.

Hypothesis. Kuhn poker is approximately constantsum polymatrix and subgame stable...

in parts of the game space that are learnable by selfplay.



Hypothesis. Kuhn poker is approximately **constantsum polymatrix** and **subgame stable...** in parts of the game space that are learnable by selfplay.






Results

Hypothesis. Kuhn poker is approximately constantsum polymatrix and subgame stable...

in parts of the game space that are learnable by selfplay.

•
$$\max_{s \in S'} |u_i(s) - \check{u}_i(s)| \le 0.0022$$

• CFR strategies were 2.9e - 5-Nash of each subgame



Other Experiments

Leduc Poker & Tiny Hanabi

 We repeated these experiments on Leduc poker and a toy Hanabi game.

• We had similar results to Kuhn poker for Leduc poker.

We found that toy Hanabi was not well-approximated by a subgame stable CSP game

Generalizations

Strategic Equivalence



Definition. Strategic Equivalence (Moulin & Vial, 1978)

For any s_{-i} ,

 s'_i is preferred to s_i in G

 \leftarrow

 s'_i is preferred to s_i in \check{G}

Strategic Equivalence to CSP Games

• We generalize the algebraic characterization of strategic equivalence to multiplayer games



Strategic Equivalence to CSP Games



Strategic equivalence extends some (but not all) guarantees.



Low vulnerability

Similar values

Mearly Exchangeable

Strategically subgame stable CSP games generalize in a meaningful way strategically constant-sum games from Moulin & Vial.

Conclusion





Main Result

Guarantees for Multi-player Games



Theorem. Given G, if there exists a CSP game \check{G} such that:

$$1.\max_{s\in S} \left| u_i(s) - \check{u}_i(s) \right| \le \delta$$

2. \check{G} is $(2n\delta, \gamma)$ -subgame stable

Then...

1. $\operatorname{Vul}_i(s^{\mu}, S_{-i}) \leq |E_i|\gamma + 2\delta$

 $|2. \exists v_i : |v_i - u_i(s^{\mu})| \le |E_i|\gamma + \delta|$

3. If *s* is a profile where each strategy is a marginal strategy from a *different* CCE...



CFR Computes Approximate Nash in Leduc Poker



Tiny Hanabi



Strategic Equivalence

Definition 6.1.1 (Strategic Equivalence). G is strategically equivalent (SE) to G' for player i if

$$u_i(s'_i, s_{-i}) \ge u_i(s_i, s_{-i}) \iff u'_i(s'_i, s_{-i}) \ge u'_i(s_i, s_{-i}) \quad \forall s_i, s'_i \in S_i, s_{-i} \in S_{-i}.$$

$$u_i(\rho) = \lambda_i u'_i(\rho) + \hat{u}_i(\rho_{-i}) \quad \forall \rho \in \mathbf{P},$$

Loss Functions (1)

$\mathcal{L}^{\delta}(\pi; \check{u}, u) \doteq \sum_{i \in N} |\check{u}_i(\pi) - u_i(\pi)|$

Loss Functions (2)

$$\mathcal{L}_{ij}^{\gamma}(\pi_{ij}, \pi_{ij}^{*}; \check{u}) \doteq \max\left(\check{u}_{ij}(\pi_{i}^{*}, \pi_{j}) - \check{u}_{ij}(\pi_{ij}), 0\right) \\ + \max\left(\check{u}_{ji}(\pi_{i}, \pi_{j}^{*}) - \check{u}_{ji}(\pi_{ij}), 0\right).$$

$$\mathcal{L}^{\gamma}(\pi, \pi^*; \check{u}) \doteq \sum_{(i,j)\in E} \mathcal{L}^{\gamma}_{ij}(\pi_{ij}, \pi^*_{ij}; \check{u})$$

Loss Functions (3)

$$\mathcal{L}(\Pi^{\delta}, \Pi^{\gamma}, \Pi^{*}; \check{u}, u) \doteq \frac{\lambda}{B_{1}} \sum_{\pi \in \Pi^{\delta}} \mathcal{L}^{\delta}(\pi; \check{u}, u) + \frac{(1-\lambda)}{B_{2}} \sum_{\pi \in \Pi^{\gamma}} \sum_{\pi^{*} \in \Pi^{*}} \mathcal{L}^{\gamma}(\pi, \pi^{*}; \check{u})$$

$$\check{u} \leftarrow \check{u} - \eta \cdot \nabla_{\check{u}} \mathcal{L}(\Pi^{\delta}, \Pi^{\gamma}, \Pi^*; \check{u}, u).$$

Sample-Based Algorithm

Algorithm 2 Compute G**Input:** $G, \Pi', \eta, T_{in}, T_{out}, \lambda, B_1, B_2$ Initialize \check{u} to all 0 $\Pi^{\times} \leftarrow X_{i \in N} \hat{\Pi}_i$ for $t_{out} \in 1...T_{out}$ do $\Pi^* \leftarrow \text{getBRs}(\check{G}, \Pi')$ for $t_{in} \in 1...T_{in}$ do $\Pi^{\delta} \leftarrow \text{sample batch of size } B_1 \text{ u.a.r. from } \Pi^{\times}$ $\Pi^{\gamma} \leftarrow \text{sample batch of size } B_2 \text{ u.a.r. from } \Pi'$ $g \leftarrow \nabla_{\check{u}} \mathcal{L}(\Pi^{\delta}, \Pi^{\gamma}, \Pi^*; \check{u}, u)$ $\check{u} \leftarrow \check{u} - \eta \cdot \frac{g}{\|g\|_2}$ end for end for {Lastly, output δ and γ } $\delta \leftarrow \max_{\pi \in \Pi^{\times}} |u_i(\pi) - \check{u}_i(\pi)|$ $\gamma \leftarrow \max_{\pi \in \Pi'} \max_{i \neq j \in N \times N} \left(\check{u}_{ij}(BR_{ij}(\pi_j), \pi_j) - \check{u}_{ij}(\pi_i, \pi_j) \right)$ return $\check{u}, \gamma, \delta$

Compute Subgame Stability (1)

```
Algorithm 1 Compute \gamma
```

```
Input: G = (N, E, P, u), a polymatrix game
\gamma \leftarrow -\infty
for (i, j) \in E do
   for \rho'_i \in P_i do
       if LP1(i, j, \rho'_i) not infeasible then
           \gamma_{ij}^{\rho_i'} \leftarrow \text{LP1}(i, j, \rho_i')
           \gamma \leftarrow \max(\gamma, \gamma_{ii}^{\rho_i'})
       end if
   end for
   for \rho'_i \in P_j do
       if LP1(j, i, \rho'_i) not infeasible then
           \gamma_{ji}^{\rho'_j} \leftarrow \text{LP1}(j, i, \rho'_j)
           \gamma \leftarrow \max(\gamma, \gamma_{ii}^{\rho'_j})
       end if
   end for
end for
return \gamma
```

Computing CSP Decompositions

$$\begin{array}{ll} \min & \underline{\delta} \\ \text{s.t.} & u_i(\rho) - \sum_{j \in -i} \check{u}_{ij}(\rho_i, \rho_j) \leqslant \underline{\delta} \quad \forall i \in N, \rho \in \mathbf{P} \\ & u_i(\rho) - \sum_{j \in -i} \check{u}_{ij}(\rho_i, \rho_j) \geqslant -\underline{\delta} \quad \forall i \in N, \rho \in \mathbf{P} \\ & \check{u}_{ij}(\rho_i, \rho_j) + \check{u}_{ji}(\rho_i, \rho_j) = c_{ij} \quad \forall i \neq j \in N, (\rho_i, \rho_j) \in \mathbf{P}_{ij}, \end{array}$$