Iterative Empirical Game Solving via Single Policy Best Response

Max Olan Smith¹, Thomas Anthony², Michael P. Wellman¹ ¹Michigan, ²Google DeepMind







Problem:

Reduce amount of training required in multiagent systems.

Real Game





Problem:

Reduce amount of training required in multiagent systems.













Which policies to add to the empirical game?



1. Solution: Solve current empirical-game.



1. Solution: Solve current empirical-game.

2. Best Response:

Add best-responses for each Player to the current solution.





During training:

- Opponent sampled at start of episode.
- Agent does not observe their current opponent.
- Results in an increase in variance of state-outcomes.
- Making the learner require many more samples.



Epoch 3:







Contributions

Current Approach

- Replace mixed-strategy opponent with a single opponent policy.
- Focus training on single most salient opponent policy.
- Reducing variance during learning.





Contributions

Current Approach

- Replace mixed-strategy opponent with a single opponent.
- Focus training on single most salient opponent.
- Reducing variance during learning.
- Two new algorithms differing choice of new opponent:
 - Mixed-Oracles
 - Mixed-Opponents





Mixed-Oracles

Insight: each player adds one new policy each epoch.

Idea:

- Train best-response to new policy.
- Transfer best-responses from older policies.



Epoch 2:

Mixed-Oracles

Insight: each player adds one new policy each epoch.

Idea:

- Train best-response to new policy.
- Transfer best-responses from older policies.







Insight: each player adds one new policy each epoch.

Idea:

Train best-response to new policy. _

Mixed-Oracles

Transfer best-responses from older policies. _









Combine Responses

- Using our previous work: **Q-Mixing**.
 - Get approximate best-response to any mixed-strategy,
 - By averaging Q-values of pure-strategy best-responses.



Mixed-Oracles: Gathering-Small



Mixed-Opponents

- The opponent's mixed-strategy may not be the most salient opponent.
- Explore strategy-space more efficiently, reduce number of epochs.

Mixed-Opponents

- The opponent's new policy may not be the most salient opponent.
- Explore strategy-space more efficiently, reduce number of epochs.
- Idea: combine opponent policies by averaging their Q-values.
 - Considers each policy's value for each action respectively.
 - Giving rise to a unique greedy policy.



Gathering-Small 2-Player

Gathering-Open 3-Player



Conclusion

р

- Replace mixed-strategy objective with pure-strategy objective in PSRO.
- Offer reductions in training time and may converge to better solutions.



Mixed-Oracles: transfer knowledge about past opponents.

Mixed-Opponents:

aggregate opponents into single new policy.

Hyperparameter Ablation (Gathering-Small):

