DeepMind

Game-Theoretic Approaches for Multiagent Reinforcement Learning in Partially Observable Environments ... and a Plea to "Go Wide"

Marc Lanctot

Joint work with many, many collaborators!

Joint Work with Many Collaborators!



DeepMind



• <u>Part 1</u> [30 min]: Game-Theoretic Approaches to Multiagent RL in Partially Observable Environments

• <u>Part 2</u> [20 min]: Beyond Zero-Sum Games and Beyond Domain-Specific Evaluation





• <u>Part 1 [30 min]</u>. Game-Theoretic Approaches to Multiagent RL in Partially Observable Environments

- <u>Part 2</u> [20 min]: Beyond Zero-Sum Games and Beyond Domain-Specific Evaluation
 - Part 1: 45 min
 - Part 2:
 - Questions: 5 min

10 min

DeepMind

Game-Theoretic Approaches to MARL in Partially Observable Games



Inspired by two-player zero-sum games

Two main "streams"

- 1. Fictitious play / best response stream
- 2. No-regret stream





• Fictitious Play:



• Start with an arbitrary policy per player (π_0^1, π_0^2) ,





• Fictitious Play:



- Start with an arbitrary policy per player (π_0^1, π_0^2) ,
 - Then, play best response

against a uniform distribution over the past policies of the opponent ($BR^{-i}_{1,..,n-1}$).





• Fictitious Play:



- Start with an arbitrary policy per player (π_0^1, π_0^2) ,
 - Then, play best response
 - against a uniform distribution over the past policies of the opponent ($BR^{-i}_{1 - n-1}$).





• Fictitious Play:

• Start with (R, P, S)= (1, 0, 0), (1, 0, 0)







• Fictitious Play:

	R	Ρ	
R	0	-1	
Ρ	1	0	

- Start with (R, P, S)= (1, 0, 0), (1, 0, 0)
- Iteration 1:
 - $\circ BR_{1}^{1},BR_{1}^{2} = P, P$
 - (½, ½, O), (½, ½, O)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)
 (½)





• Fictitious Play:

	R	Ρ	Ρ	
R	0	-1	-1	
Ρ	1	0	0	
Р	1	0	0	

- Start with (R, P, S)= (1, 0, 0), (1, 0, 0)
- Iteration 1:
 - $\circ BR_{1}^{1},BR_{1}^{2} = P, P$
 - (½, ½, O), (½, ½, O)
- Iteration 2:
 - $\circ \quad \mathsf{BR}^{1}_{2}, \mathsf{BR}^{2}_{2} = \mathsf{P}, \mathsf{P}$
 - $\circ \quad (\frac{1}{3}, \frac{2}{3}, 0), (\frac{1}{3}, \frac{2}{3}, 0)$





• Fictitious Play:

	R	Ρ	Ρ	S	
R	0	-1	-1	1	
Ρ	1	0	0	-1	
Ρ	1	0	0	-1	
S	-1	1	1	0	

- Start with (R, P, S)= (1, 0, 0), (1, 0, 0)
- Iteration 1:
 - $\circ \quad \mathsf{BR}^{1}_{1}, \mathsf{BR}^{2}_{1} = \mathsf{P}, \mathsf{P}$
 - (½, ½, O), (½, ½, O)
- Iteration 2:
 - $\circ BR_{2}^{1},BR_{2}^{2} = P, P$
 - $\circ \quad (\frac{1}{3}, \frac{2}{3}, 0), (\frac{1}{3}, \frac{2}{3}, 0)$
- Iteration 3:
 - $\circ \quad \mathsf{BR}^{1}_{3}, \mathsf{BR}^{2}_{3} = \mathsf{S}, \mathsf{S}$
 - $\circ \quad (1/_4, 1/_2, 1/_4), \ (1/_4, 1/_2, 1/_4)$



• Fictitious Play:

	R	Ρ	Ρ	S	S	
R	0	-1	-1	1	1	
Ρ	1	0	0	-1	-1	
Ρ	1	0	0	-1	-1	
S	-1	1	1	0	0	
S	-1	1	1	0	0	

- Start with (R, P, S)= (1, 0, 0), (1, 0, 0)
- Iteration 1:
 - $\circ \quad \mathsf{BR}^{1}_{1}, \mathsf{BR}^{2}_{1} = \mathsf{P}, \mathsf{P}$
 - (½, ½, O), (½, ½, O)
- Iteration 2:
 - \circ BR¹₂, BR²₂ = P, P
 - $\circ \quad (\frac{1}{3}, \frac{2}{3}, 0), (\frac{1}{3}, \frac{2}{3}, 0)$
- Iteration 3:
 - $\circ BR_{3}^{1},BR_{3}^{2} = S, S$
 - $\circ \quad (1/_4, 1/_2, 1/_4), (1/_4, 1/_2, 1/_4)$



• double oracle [HB McMahan 2003]:



DeepMind

- Start with an arbitrary policy per player $(\pi^{1}_{0}, \pi^{2}_{0})$,
 - Compute (pⁿ,qⁿ) by solving the game at iteration n
 - Then, best response against
 (pⁿ,qⁿ) and get a new best
 response (BR¹_n,BR¹_n).



• double oracle:











• Start with (R, P, S)= (1, 0, 0), (1, 0, 0)

• double oracle:

	R	Ρ	
R	0	-1	
Ρ	1	0	

- Start with (R, P, S)= (1, 0, 0), (1, 0, 0)
- Iteration 1:
 - $\circ \quad \mathsf{BR}^{1}_{1}, \mathsf{BR}^{2}_{1} = \mathsf{P}, \mathsf{P}$
 - \circ Solve the game : (0, 1, 0), (0, 1,

0)





• double oracle:

	R	Ρ	S	
R	0	-1	1	
Ρ	1	0	-1	
S	-1	1	0	

- Start with (R, P, S)= (1, 0, 0), (1, 0, 0)
- Iteration 1:
 - $\circ \quad \mathsf{BR}^{1}_{1}, \mathsf{BR}^{2}_{1} = \mathsf{P}, \mathsf{P}$
 - Solve the game : (0, 1, 0), (0, 1,
 0)
- Iteration 2:
 - $\circ BR_{2}^{1}, BR_{2}^{2} = S, S$
 - $\circ \quad \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right), \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right)$





• Regret-matching (Hart & Mas-Colell '00):







• Regret-matching (Hart & Mas-Colell '00):

0

Ρ

0

S

• For row player 1, column player fixed



(showing row player's utility)



Regrets

Cumulative

0

R

- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player fixed

$$\circ \quad t=0, \ \pi^{1}_{0} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}), \ R^{1} = 0$$









- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player **fixed**

0

S

- $\circ \quad t=0, \ \pi^{1}_{0} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}), \ R^{1} = 0$
- Reg. *not* playing rock:

0

R

= (-0.4 + 0.4) - 0 = 0

0

Ρ



(showing row player's utility)



Regrets

Cumulative



- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player **fixed**
 - $\circ \quad t=0, \ \pi^{1}_{0} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}), \ R^{1} = 0$
 - Reg. *not* playing paper:
 - **•** = (0.2 0.4) 0 = -0.2





(showing row player's utility)



O DeepMind

Multi-Agent and Al

- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player fixed
 - $\circ \quad t=0, \ \pi^{1}_{0} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}), \ R^{1} = 0$
 - Reg. *not* playing scissors:
 - **•** = (-0.2 + 0.4) 0 = +0.2





(showing row player's utility)



O DeepMind

- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player fixed
 - $t=1, \pi_1^1 = (0, 0, 1)$

Normalize positive cumulative regrets (set others to 0)









- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player fixed

o $t=1, \pi_1^1 = (0, 0, 1), R^1 = 0.2$









- Regret-matching (Hart & Mas-Colell '00):
 - For row player 1, column player fixed

 $\circ \quad t=2, \ \pi_{2}^{1} = (0, \ 0, \ 1), \ R^{1} = 0.2$









Fictitious Self-Play (FSP) [Heinrich, Lanctot, & Silver '15]

- Idea: Fictitious play + reinforcement learning in one online agent
- Update rule in sequential setting *equivalent* to standard fictitious play (matrix game)

- 1. Best response (BR):
 - Estimate a best response
 - Trained via RL (e.g. Q-learning)
 - Circular buffer of (s, a, s', r) tuples

2. Average policy (AVG):

- Estimate the time-average policy
- Trained via supervised learning
- Reservoir buffer of (s, a) pairs



Neural Fictitious Self-Play (NFSP) [Heinrich & Silver '16]

• Approximate NE via two neural networks:



• Competitive with strong computer poker programs when it was released





Policy-Space Response Oracles (Lanctot et al. '17)





Best Response Policy Iteration and Diplomacy [Anthony et al. '20]

- Classic board game
- 7-player game
- Simultaneous moves
- 10²¹ 10⁶⁴ legal actions *per turn*
- Mixed-motives:
 - Winning requires alliances
 - Players negotiate for territory

Current focus on *no-press variant*.





Best-Response Policy Iteration [Anthony et al. '20]



Best-Response Policy Iteration [Anthony et al. '20]

- Input:
 - Base policy π_b
 - Candidate policy π_c
 - Environment dynamics T(s, a) -> s'
 - Value function V(s')
- Algorithm:
 - At each turn, given state s:
 - Sample a few base profiles a_{i} from $\pi_{b}(s)$ for all players but *i*
 - Sample several candidate actions a_i from $\pi_c(s)$
 - Plug the sampled actions into $T(s, a'_i, a_j) \rightarrow s'$
 - Get V(s') (for player i)
 - Play the candidate action with the best average value against the base profiles → sampled best response (SBR)





BRPI Policy Improvement

But what best response should we be imitating?

• Iterated Best Responses



BRPI Policy Improvement

But what best response should we be imitating?

- Iterated Best Responses
- Fictitious Play (1) -- "à la NFSP"



BRPI Policy Improvement

But what best response should we be imitating?

- Iterated Best Responses
- Fictitious Play (1) -- "à la NFSP"
- Fictitious Play (2)


BRPI in Diplomacy: Results

	SL [90]	A2C [90]	SL (ours)	FPPI-1	IBR	FPPI-2
SL [90]	14.2%	8.3%	16.3%	2.3%	1.8%	0.8%
A2C [90]	15.1%	14.2%	15.3%	2.3%	1.7%	0.9%
SL (ours)	12.6%	7.7%	14.1%	3.0%	1.9%	1.1%
FPPI-1	26.4%	28.0%	25.9%	14.4%	7.4%	4.5%
IBR	20.7%	30.5%	25.8%	20.3%	12.9%	10.9%
FPPI-2	19.4%	32.5%	20.8%	22.4%	13.8%	12.7%

Table 1: Average scores for 1 row player vs 6 column players. BRPI methods give an improvement over A2C or supervised learning. All numbers accurate to a 95% confidence interval of $\pm 0.5\%$. Bold numbers are the best value for single agents against a given set of 6 agents, italics are for the best result for a set of 6-agents against each single agent.

Human-Level No-Press Diplomacy [Gray et al. '20]

- Human data \rightarrow DipNet
- DipNet provides policy π and value net v
- Use regret matching in stage game
- Get payoffs from sims / search
 - Use policy for rollouts + selection
 - Value net after some horizon
- Human-level play on webdiplomacy



Counterfactual regret minimization (CFR) (Zinkevich et al. '08):

Basis of success in Poker AI for two-player zero-sum games:



Initial policies iteration, t = 0



Counterfactual regret minimization (CFR) (Zinkevich et al. '08):

Basis of success in Poker AI for two-player zero-sum games:

Player 1
$$\pi_1^0 \rightarrow \pi_1^1$$

Player 2 $\pi_2^0 \rightarrow \pi_2^1$

Counterfactual regret minimization (CFR) (Zinkevich et al. '08):

Basis of success in Poker AI for two-player zero-sum games:





Counterfactual regret minimization (CFR) (Zinkevich et al. '08):

Basis of success in Poker AI for two-player zero-sum games:



Counterfactual Regret Minimization (CFR) [Zinkevich et. al' 08]

- Tabular method (policy iteration)
- Each information state, s:
 - Compute **counterfactual regrets** r(s,a)
 - Accumulate: R(s,a) += r(s,a)
 - Use regret-matching for new $\pi(s)$



e.g.



Advantage vs. Regrets

A key notion in CFR is an **immediate regret**:

$$r(s, a) = q_{\pi,i}^{c}(s, a) - v_{\pi,i}^{c}(s)$$
counterfactual q-value
joint policy
return to player *i*
(player to play at *S*)

 \rightarrow This is just a (counterfactual) advantage!

RL values vs. Counterfactual values

$$q_{\pi,i}(s,a) = \frac{1}{\beta_{-i}(s)} q_{\pi,i}^c(s,a)$$
"RL-style" q-value (conditioned on reaching s)
$$(s,a) = \frac{1}{\beta_{-i}(s)} q_{\pi,i}^c(s,a)$$
(weighted sum over histories)

Probability that s is reached given opponents' policies



Bayes Normalizer



$$\beta_{-i}(s,\pi) = \sum_{h \in s} \Pr(h)$$



Q-based Policy Gradient

A.K.A. "all-actions" policy gradient A.K.A. Mean Actor-Critic (Allen et al. '17)

$$\nabla_{\boldsymbol{\theta}}^{\text{QPG}}(s) = \sum_{a} [\nabla_{\theta} \pi(s, a; \boldsymbol{\theta})] \left(q(s, a; \mathbf{w}) - \sum_{b} \pi(s, b; \boldsymbol{\theta}) q(s, b, \mathbf{w}) \right)$$



Regret-based Policy Gradient [Srinivasan et al. '18]

Instead of maximizing objective, **minimize regret**:

$$\nabla_{\boldsymbol{\theta}}^{\text{RPG}}(s) = -\sum_{a} \nabla_{\boldsymbol{\theta}} \left(q(s, a; \mathbf{w}) - \sum_{b} \pi(s, b; \boldsymbol{\theta}) q(s, b; \mathbf{w}) \right)^{+}$$

where $(x)^{+} = \max(0, x)$

→ Gradient **descent** (instead of ascent)



Replicator Dynamics (Taylor & Jonker '78)

- Population state x evolves inspired by biologically inspired operators
- Proportion of member i, x_i , grows according to their fitness f

$$\dot{x}_i = oldsymbol{x}_i \left(f(oldsymbol{x})_i - oldsymbol{ar{f}}(oldsymbol{x})
ight)$$

 $ar{f}(oldsymbol{x}) = \sum_j x_j f(oldsymbol{x})_j$

in Reinforcement Learning terms:

$$\dot{\pi}(a) = \frac{\pi(a)}{[q(a) - v]}$$



Policy Gradient vs. Replicator Dynamics

Policy Gradient (Advantage Actor-Critic)

Replicator Dynamics

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\pi}} \left[\nabla_{\boldsymbol{\theta}} \log \pi(a_t \mid s_t; \boldsymbol{\theta}) A(s_t, a_t; \boldsymbol{w}, \boldsymbol{\theta}) \right]$$

$$\dot{\pi}(a) = \pi(a)A(a)$$



Neural Replicator Dynamics [Omidshafiei, Hennes, Morrill et al. '19]



$$\boldsymbol{\theta}_{t} = \boldsymbol{\theta}_{t+1} + \eta \sum_{s,a} \nabla_{\boldsymbol{\theta}} y_{t-1}(s_{t}, a_{t}; \boldsymbol{\theta}) A(s_{t}, a_{t}; \boldsymbol{\theta}, \boldsymbol{w})$$

$$\textbf{Logits, where policy is}_{\boldsymbol{\pi} = softmax(\boldsymbol{y})}$$

$$\textbf{Advantage q(s,a)-v(s)}$$

Neural Replicator Dynamics: Results



Advantage Regret-Matching Actor-Critic [Gruslys et al. '20]

<u>Goal:</u> sample-based model-free CFR with function approximation.

Monte Carlo CFR [Lanctot et al. '09]: Sample trajectories in interesting portions of the game tree.







Limitations of (Outcome Sampling) Monte Carlo CFR [Lanctot et al '09]

Memory hungry tabulates all information states

No generalization

Huge variance

We want to solve all those problems by using neural networks and RL-style trajectory sampling





Problem 1: Cumulative Regrets

Problem: neural networks can not accumulate regrets

Solution: reformulate CFR in terms of mean regrets

Instead of cumulative regrets $\ R^T(s,a)$

Learn an estimate
$$\hat{R}^T(s, a) \rightarrow \frac{R^T(s, a)}{T}$$

Inspired by Regression CFR (RCFR) [Waugh et al. '15]





Problem 2: Sampling and Variance

Corresponds to T CFR iterations

Problem: MCCFR estimator of $\hat{R}^T(s, a)$ can have huge variance **Solution:**

- Maintain a reservoir of past joint policies over epochs { 1, 2, ..., T }
- 2. At current **epoch** t, generate data:
 - a. Sample past checkpoint j ~ { 1... t-1 } uniformly
 - b. Sample our actions with exploratory behavior policy $u_i^t(s)$
 - c. Sample opponent actions using π^j
 - d. Tabulate sampled regrets $\hat{r}^{j}(s, a)$ with **history-based critics**
- 3. Train regressor to predict mean regrets



Problem 2: Sampling and Variance

Problem: MCCFR estimator of $\hat{R}^T(s, a)$ can have huge variance **Solution:**

- 1. Maintain a **reservoir of <u>past joint policies</u>** over **epochs** { 1, 2, ..., T }
- 2. At current **epoch** t, generate data:
 - a. Sample past checkpoint j ~ { 1... t-1 } uniformly
 - b. Sample our actions with exploratory behavior policy $u_i^t(s)$
 - c. Sample opponent actions using π^j
 - d. Tabulate sampled regrets $\hat{r}^{j}(s, a)$ with **history-based critics**
- 3. Train regressor to predict mean regrets

As a result, we learn
$$\hat{W}(s,a) \rightarrow W(s,a) = R(s,a) k(s)$$



Problem 3: History / world state critics

Problem:

We need to evaluate advantages
$$q_{\pi^j}(h,a) - \sum_b q_{\pi^j}(h,b)$$

But in epoch T we produce trajectories with an exploratory behavior policy $u_i^t(s)$

Solution:

Train value functions using off-policy-RL.





ARMAC Results: benchmark domains

Leduc Poker

 $10^{0} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{3}} \underbrace{10^{0}}_{10^{4}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{3}} \underbrace{10^{0}}_{10^{4}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{3}} \underbrace{10^{0}}_{10^{4}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{3}} \underbrace{10^{0}}_{10^{4}} \underbrace{10^{0}}_{10^{4}} \underbrace{10^{0}}_{10^{2}} \underbrace{10^{0}}_{10^{3}} \underbrace{10^{0}}_{10^{4}} \underbrace{10$

Stable learning curves.





Liars dice

ARMAC Results: action-abstracted no-limit Hold'em



Local Best Response (LBR) bound on exploitability: y-axis = exploitation value, x-axis = epochs Roughly matches the performance of state of art Poker bots from 2016



Multi-Agent and Al

ARMAC Results: Atari



Hypothesis: regret matching has nice exploratory policies





Learning with Regularization: Friction FoReL [Perolat et al. '21] Follow The Regularized Leader: Adding a

$$y_t^i(a^i) = \int_0^{\cdot} Q_{\pi_s}^i(a^i) ds \quad \text{and} \quad \pi_t^i = \operatorname{argmax}_{p \in \Delta A} \Lambda^i(p, y_t^i)$$

With : $\Lambda^i(p, y) = \langle y, p \rangle - \phi_i(p)$ and $\phi_i(p)$ is a regularisation for the policy projection.

In zero-sum two-player games, the following quantity is preserved and the learning trajectory is recurrent:

$$J(y) = \sum_{i=1}^{2} \left[\phi_i^*(y_i) - \langle y_i, \pi_i^* \rangle \right]$$



Adding a policy dependent term:

$$r_{\pi}^{i}(a) = r^{i}(a^{i}, a^{-i}) - \eta \log \frac{\pi^{i}(a^{i})}{\mu^{i}(a^{i})} + \eta \log \frac{\pi^{-i}(a^{-i})}{\mu^{-i}(a^{-i})}$$

This policy dependent term transforms a recurrent learning dynamic to a convergent one:

$$\frac{d}{dt}J(y) = \sum_{i=1}^{2} \underbrace{[V^{i}_{\pi^{i}_{t},\pi^{*-i}} - V^{i}_{\pi^{*}}]}_{\leq 0 \text{ because } \pi^{*} \text{ is a Nash}} -\eta \sum_{i=1}^{2} KL(\pi^{*i},\pi^{i}_{t})$$





Learning with Regularization Friction FoReL [Perolat et al. '21]

Video 50x slower



Increasing bias to the solution

We want to learn at high regularization and low bias!



Learning with Regularization: Friction FoReL [Perolat et al. '21]

Regularization centered around [1/3, 1/3, 1/3].

 π_0



Solution : [0.38, 0.48, 0.12] π_1

Regularization centered around [0.38, 0.48, 0.12]. π_1



Solution : [0.29, 0.62, 0.07]



Learning with Regularization: Friction FoReL [Perolat et al. '21]

Regularization centered around $[1_3, 1_3, 1_3]$.

 π_0



Solution : [0.38, 0.48, 0.12] π_1

Regularization centered around [0.38, 0.48, 0.12]. π_1



 $r_{\pi}^{i}(a) = r^{i}(a^{i}, a^{-i}) - \eta \log \frac{\pi^{i}(a^{i})}{\pi_{0}^{i}(a^{i})} + \eta \log \frac{\pi^{-i}(a^{-i})}{\pi_{0}^{-i}(a^{-i})}$ Converges to π_{1}

$$r_{\pi}^{i}(a) = r^{i}(a^{i}, a^{-i}) - \eta \log \frac{\pi^{i}(a^{i})}{\pi_{1}^{i}(a^{i})} + \eta \log \frac{\pi^{-i}(a^{-i})}{\pi_{1}^{-i}(a^{-i})}$$

Converges to π_{2}

Recentering around the previous fixed point will decrease strictly the distance to the Nash of the original game:

$$\Xi(\pi^*, \pi_k) - \Xi(\pi^*, \pi_{k-1}) = \underbrace{-\Xi(\pi_k, \pi_{k-1}) + \frac{1}{\eta} \sum_{i=1}^N (m_k^i + \delta_k^i + \kappa_k^i)}_{<0}$$

Learning with Regularization

Regularization centered around [1/3, 1/3, 1/3].

0.6-0.6-0.2-5 0,0 0,0 0,0 0,0 R

Solution : [0.38, 0.48, 0.12]

Regularization centered around [0.38, 0.48, 0.12].



Regularization centered around [0.29, 0.62, 0.07].



Solution : [0.19, 0.72, 0.07]

Regularization centered around [0.19, 0.72, 0.07].



Solution : [0.13 0.76 0.09]



Friction FoReL: Experiments in Sequential

Games

Convergence in Sequential Imperfect Information Games (Kuhn Tabular):

lr: 0.001, rf: 0.001 Ir: 0.05, rf: 0.05 fixed mu and decaying eta (to a fixed eta): 100 policy policy 10^{-1} avg policy eta : 1.0 1.2 eta 0.2 10-3 10-1 eta 0.05 NashConv 0.02 VashConv : 0.01 ota 10-5 eta : 0.0 10-2 10-7 NashConv for g 10^{-9} 10-3 1e+06 1e+06 0 0.0 500000 1000000 1500000 2000000 2500000 3000000 3500000 4000000 iterations iterations iterations Ir: 0.001, rf: 0.002 Ir: 0.05, rf: 0.1 fixed eta and refreshed mu 100 10 policy policy 4.0 avg policy eta · 10 eta : 0.5 3.5 leduc poker 10eta : 0.2 0.05 NashConv 10-1 3.0 eta : 0.02 NashConv 10-5 eta : 0.0 2.5 game : 2.0 10-3 for 1.5 10-2 NashConv 1.0 10^{-4} 05 1e+06 1e+06 0 0.0 iterations iterations 500000 1000000 1500000 2000000 2500000 3000000 3500000 4000000 iterations

Convergence in Sequential Imperfect Information Games (Leduc with Neural Network and a NeuRD loss):

Hidden Information Game Competition (HIGC)

Hidden Information Games Competition (2021)

About Contact Docs Games Rules Schedule Registration Talk to us at Discord G

Hidden Information Games Competition (HIGC) tests AI bots on large two-player zero-sum games with imperfect information, such as Poker or DarkChess. The players do not have perfect knowledge about everything that goes on in the game: they receive only partial observations about the world's real state. The goal of the AI is to play the games as well as possible against any opponent.





Join the community at the **Discord server**

- Reconnaissance Blind Chess
- Gin Rummy
- One hidden game!

higcompetition.info/









DeepMind

A Plea to "Go Wide": Beyond Zero-Sum and Domain-Specific Evaluation



Games, RL, and AI



Arthur Samuel - Checkers



IBM - Chess (DeepBlue)



Poker: DeepStack & Libratus





DeepMind - Go (AlphaGo)



Gerald Tesauro - Backgammon

Minimax



<u>Max-min</u>: P1 looks for a π_1 such that

 $v_1 = \max_{\pi_1} \min_{\pi_2} u_1(\pi_1, \pi_2)$

<u>Min-max</u>: P1 looks for a π_1 such that

$$v_1 = \min_{\pi_2} \max_{\pi_1} u_1(\pi_1, \pi_2)$$

In two-player, zero-sum these are the same!

John von Neumann 1928

---> The Minimax Theorem





Consequences of Minimax

The optima
$$\ \pi^*=(\pi_1^*,\pi_2^*)$$

- These exist! (They sometimes might be stochastic.)
- Called a minimax-optimal joint policy. Also, a Nash equilibrium.
- They are interchangeable:

•
$$\forall \pi^*, \pi^{*\prime} \Rightarrow (\pi_1^*, \pi_2^{*\prime}), (\pi_1^{*\prime}, \pi_2^*)$$
 also minimax-optimal

• Each policy is a best response to the other.




Minimax (Outside Two-Player Zero-Sum Games)

<u>Max-min</u>: Player i looks for a π_i such that:

$$v_i^{maxmin} = \max_{\pi_i} \min_{\pi_{-i}} u_i(\pi_i, \pi_{-i})$$
 (Paranoid)





Minimax (Outside Two-Player Zero-Sum Games)

<u>**Max-min**</u>: Player i looks for a π_i such that:

$$v_i^{maxmin} = \max_{\pi_i} \min_{\pi_{-i}} u_i(\pi_i, \pi_{-i})$$
 (Paranoid)

<u>Min-max</u>: Player j looks for a π_j (against π_i , i.e. in π_{-i}) such that

$$v_i^{minmax} = \min_{\pi_{-i}} \max_{\pi_i} u_i(\pi_i, \pi_{-i})$$
 (Optimistic)





A joint policy used by all n players:
$$\pi = (\pi_1, \pi_2, \cdots, \pi_n)$$

such that
$$\forall i, u_i(\pi) \ge \max_{\pi'_i} u_i(\pi'_i, \pi_{-i})$$





Suppose each player computes an equilibrium: π^A, π^B, π^C





$$\pi^A, \pi^B, \pi^C$$

•
$$(\pi_1^A,\pi_2^B,\pi_3^C)$$
 is generally not an equilibrium





$$\pi^A, \pi^B, \pi^C$$

- $(\pi_1^A,\pi_2^B,\pi_3^C)$ is generally not an equilibrium
- Each equilibrium might have different values for all players





$$\pi^A, \pi^B, \pi^C$$

- $(\pi_1^A,\pi_2^B,\pi_3^C)$ is generally not an equilibrium
- Each equilibrium might have different values for all players
- Which equilibrium should you "choose"?
 - Equilibrium selection problem





$$\pi^A, \pi^B, \pi^C$$

- $(\pi_1^A, \pi_2^B, \pi_3^C)$ is generally not an equilibrium
- Each equilibrium might have different values for all players
- Which equilibrium should you "choose"?
 - Equilibrium selection problem
- Also PPAD-Hard (Daskalakis, Goldberg, Papadimitriou + Chen & Deng)





• How do you tell if chess program is super-human?



- How do you tell if chess program is super-human?
- What is "super-human" Iterated Prisoner's Dilemma?



- How do you tell if chess program is super-human?
- What is "super-human" Iterated Prisoner's Dilemma?
- If minimax is optimal, why doesn't it win RoShamBo competitions?



- How do you tell if chess program is super-human?
- What is "super-human" Iterated Prisoner's Dilemma?
- If minimax is optimal, why doesn't it win RoShamBo competitions?
- What should general agents learn in a multiagent setting?



Back to the Essentials: What's the Goal?







 $\cdots \leftarrow \pi_n^{t-2} \leftarrow \pi_n^{t-1} \leftarrow \pi_n^t \quad \bigcirc_{\text{Agent N}}$

















Agent i perceives:

$$\cdots, (o_i^{t-2}, r_i^{t-2}), (o_i^{t-1}, r_i^{t-1}), (o_i^t, r_i^t)$$

• Set of **deviations** considered: Φ



Agent i perceives:

$$\cdots, (o_i^{t-2}, r_i^{t-2}), (o_i^{t-1}, r_i^{t-1}), (o_i^t, r_i^t)$$

- Set of deviations considered: Φ
- π_i^t chosen in a way that minimizes regret w.r.t. Φ



Agent i perceives:

$$\cdots, (o_i^{t-2}, r_i^{t-2}), (o_i^{t-1}, r_i^{t-1}), (o_i^t, r_i^t)$$

- Set of **deviations** considered: Φ
- π_i^t chosen in a way that minimizes regret w.r.t. Φ
- ... based on history of correlated play!

Agent i perceives:

$$\cdots, (o_i^{t-2}, r_i^{t-2}), (o_i^{t-1}, r_i^{t-1}), (o_i^t, r_i^t)$$

- Set of **deviations** considered: Φ
- π_i^t chosen in a way that minimizes regret w.r.t. Φ
- ... based on history of correlated play!

Classes of (extensive-form) correlated equilibria



Deviation sets $\, \Phi \,$

- Functions of the entire sequence of past plays
- Sequential deviations too (not just entire policy deviations)
- New counterfactual deviations and associated equilibria in self-play

arXiv.org > cs > arXiv:2012.05874

Computer Science > Computer Science and Game Theory

[Submitted on 10 Dec 2020 (v1), last revised 17 Dec 2020 (this version, v2)]

Hindsight and Sequential Rationality of Correlated Play

Dustin Morrill, Ryan D'Orazio, Reca Sarfati, Marc Lanctot, James R. Wright, Amy Greenwald, Michael Bowling

Efficient Deviation Types and Learning for Hindsight Rationality in Extensive-Form Games

Dustin Morrill, Ryan D'Orazio, Marc Lanctot, James R. Wright, Michael Bowling, Amy Greenwald

→ See also first COMARL seminar by Michael Bowling



Correlated Equilibrium (CE)

- CE arise as a *result* of learning
- (More) compatible with Bayesian perspective
- Compatible with a prescriptive agenda (via n
- Tractable!! (to find one)





Robert Aumann



	Rock	Paper	Scissors
Rock	0,0	0,1	1,0
Paper	1,0	0,0	0, 1
Scissors	0, 1	1,0	0,0

Figure 7.6: Shapley's Almost-Rock-Paper-Scissors game.



	Rock	Paper	Scissors
Rock	0,0	0,1	1,0
Paper	1,0	0,0	0, 1
Scissors	0, 1	1,0	0,0

Figure 7.6: Shapley's Almost-Rock-Paper-Scissors game.



	Rock	Paper	Scissors
Rock	0,0	0, 1	1,0
Paper	1,0	0,0	0, 1
Scissors	0, 1	1,0	0,0

9. Then recommend each player their

r recommendation?

Figure 7.6: Shapley's Almost-Rock-Paper-Scissors game.



	Rock	Paper	Scissors
Rock	0,0	0, 1	1,0
Paper	1,0	0,0	0, 1
Scissors	0, 1	1,0	0,0

9. Then recommend each player their

r recommendation? ty ½. Is this a CE? If so, which one

Figure 7.6: Shapley's Almost-Rock-Paper-Scissors game.



Calibrated Learning [Foster & Vohra '97]

The following process converges to a CE:

Repeat over many trials $t \rightarrow T$:

- 1. Compute a "calibrated forecast" of the opponents' policies (i.e. asymptotically consistent as $t \rightarrow +inf$)
- 2. Best respond to the forecast

Claim: basis for a "best response stream" outside zero-sum



Bayesian Perspectives and Meta-Learning

Meta-learning of Sequential Strategies

Pedro A. Ortega, Jane X. Wang, Mark Rowland, Tim Genewein, Zeb Kurth-Nelson, Razvan Pascanu, Nicolas Heess, Joel Veness, Alex Pritzel, Pablo Sprechmann, Siddhant M. Jayakumar, Tom McGrath, Kevin Miller, Mohammad Azar, Ian Osband, Neil Rabinowitz, András György, Silvia Chiappa, Simon Osindero, Yee Whye Teh, Hado van Hasselt, Nando de Freitas, Matthew Botvinick, and Shane Legg DeepMind

Meta-trained agents implement Bayes-optimal agents

Vladimir Mikulik, Grégoire Delétang, Tom McGrath, Tim Genewein, Miljan Martic, Shane Legg, Pedro A. Ortega[†] DeepMind London, UK



A Generalized Training Approach for MARL [Muller et al. '19]

- Game-theoretic training:
 - (Meta-)solve empirical game
 - Find best response oracles
- Outside two-player zero-sum:
 - Use α-Rank as a tractable solution concept





Joint Policy-Space Response Oracles (JPSRO) [Marris et al. '21]

- Game-theoretic training:
 - (Meta-)solve empirical game
 - Find best response oracles
- Finds joint distribution





- New CE meta-solvers driving by Gini impurity (eq. selection)
- Converges to CE & CCE

 $\rightarrow @ICML$

- Works for n-player general-su
- Stochastic meta-solver poss



2-player Trade Comm (simplified trading game)





• Progress in Deep RL benefited from the generality of Atari





• Progress in Deep RL benefited from the **generality** of Atari



- Progress in Deep RL benefited from the **generality** of Atari
- Games provide the formalism, have been the heart of multiagent RL





Figure 2: An initial board (left) and a situation requiring a probabilistic choice for A (right).









Figure 1: A 10 by 10 grid world.











- Progress in Deep RL benefited from the generality of Atari
- Games provide the formalism, have been the heart of multiagent RL
- The MARL community is split across many different types:
 - Competitive vs. Cooperative vs. Mixed / Social Dilemmas





Figure 2: An initial board (left) and a situation requiring a probabilistic choice for A (right).




A Plea to "Go Wide": Beyond Domain-Specific Evaluation

- Progress in Deep RL benefited from the **generality** of Atari
- Games provide the formalism, have been the heart of multiagent RL
- The MARL community is split across many different types:
 - Competitive vs. Cooperative vs. Mixed / Social Dilemmas
- A generally intelligent agent is canable across many environments!

$$\Upsilon(\pi) := \sum_{\mu \in E} 2^{-K(\mu)} V_{\mu}^{\pi}.$$
 Value achieved Measure of Intelligence

Sum over environments

"Intelligence measures an agent's ability to achieve goals in a wide range of environments"



-- Shane Legg & Marcus Hutter '07, "Universal Intelligence"

A Plea to "Go Wide": Beyond Domain-Specific Evaluation

- Progress in Deep RL benefited from the **generality** of Atari
- Games provide the formalism, have been the heart of multiagent RL
- The MARL community is split across many different types:
 - Competitive vs. Cooperative vs. Mixed / Social Dilemmas
- A generally intelligent agent is capable **across many environments**!
- An artificial general intelligence (AGI) can handle all cases



A Plea to "Go Wide": Beyond Domain-Specific Evaluation

- Progress in Deep RL benefited from the **generality** of Atari
- Games provide the formalism, have been the heart of multiagent RL
- The MARL community is split across many different types:
 - Competitive vs. Cooperative vs. Mixed / Social Dilemmas
- A generally intelligent agent is capable **across many environments**!
- An artificial general intelligence (AGI) can handle all cases
- Games:
 - Are externally defined or inspired by human problems
 - Have been traditional benchmarks of rationality for thousands of years
 - Have formally-defined interactions
 - Can be easily simulated and run many times
 - \circ $\;$ Are enjoyable to play and easy to demonstrate with humans $\;$



Why Generality in MARL?

Independent learners with their training partners:



Independent learners with similarly-trained





Self-Play / Independent Agents Do Not Generalize











How Do We Get There?

Some ideas:

- 1. More domains in empirical evaluations
- 2. More openly available code ... for MARL specifically
- 3. Increased focus on *ad-hoc setting* for evaluation
 - a. Especially with human(s) in the loop, where possible
- 4. Wider and more dynamic evaluation regimes

Importance of Open-Source and Widely Usable Libraries

theano







OpenSpiel: A Framework for RL in Games

Supports:

- n-player games
- Zero-sum, coop, general-sum
- Perfect / imperfect info
- Simultaneous-move games

github.com/deepmind/open_spiel/



- One main API for all games
- Atari Learning Env. of multiagent / games
- Release 1.0 coming in August (!)
 - with new mean-field game API



MAVA: A New Open-Source Framework for MARL

- Scalable multiagent training framework: • Scalable multiagent training framework: • Scalable multiagent training framework:
 - Built on open-source technologies:
 - Reverb
 - Launchpad
 - Acme
- Integration with many environments / libraries!
 - PyMARL
 - OpenSpiel
 - PettingZoo
 - Flatland
 - Robocup
 - Starcraft Multi-Agent Challenge (SMAC)



https://github.com/instadeepai/Mava

Scalable Evaluation of MARL with Melting Pot [Leibo, Duéñez-Guzmán, Vezhnevets, Agapiou, et al. '21]

- Compare algorithms/agents
- Focus on **generalization**
- Over 80 unique test scenarios!
- Eval in held out test scenarios
- Agnostic to training method



Potential General Evaluation Regimes

- A suite of general learning agents that can play all sides/roles
- A suite of games (domains) to evaluate on

Agent match-ups:

- Agent versus fixed reference set
- Agent versus agent

- Agent sampling distr: fixed versus adaptive
- Inter-match memory: none (blank slate) vs. learning (lifelong)



Potential General Evaluation Regimes: A vs. fixed ref set

... with

à la Melting Pot. Evaluating agent A: test-time adaptation For all games in selected games, G:

• Decide on a fixed reference set Ref(G)

• E.g. self-play RL benchmark, known strategies, etc.

- Evaluate A over/with other agents in Ref(G)
- Get overall return V(G)

Report V(G) for all selected games



Potential General Evaluation Regimes: A vs. fixed ref set

... with

à la Melting Pot. Evaluating agent A: test-time adaptation Agent A is always learning.

For all games in selected games, G:

- Ensure observations encode identification of G.
- Decide on a fixed reference set Ref(G)
 - E.g. self-play RL benchmark, known strategies, etc.
- Evaluate A over/with other agents in Ref(G)
- Get overall return V(G)

Report V(G) for all selected games



Potential General Evaluation Regimes: A vs. fixed ref set

... with

à la Melting Pot. Evaluating agent A: adaptive distr. For all games in selected games, G:

- Decide on a fixed reference set Ref(G)
 - E.g. self-play RL benchmark, known strategies, etc.
- Initialize OppDist(Ref(G)) to uniform
- For epochs t = 1 ... T:
 - Evaluate A over/with other agents in Ref(G) using OppDist
 - Get overall return V(G, t)
 - Adjust OppDist adversarially

Report V(G, t = 1.. T) for all selected games.



Potential General Evaluation Regimes: Agent vs. Agent

- Similarly to fixed reference set, except:
 - Reference set is not fixed
 - Other agents might be in reference sets

Goal: fully online, continuous, lifelong evaluation across many environments

Conclusions & Summary

- Game-theoretic approaches to partially observable games
 - Inspired by two-player zero-sum games
 - Shown to scale to very large domains
 - Two main streams: best response and no-regret
 - Many ideas can be generalized outside two-player zero-sum
 - Correlated equilibria as link between prescriptive/descriptive views



Conclusions & Summary

- Game-theoretic approaches to partially observable games
 - Inspired by two-player zero-sum games
 - Shown to scale to very large domains
 - Two main streams: best response and no-regret
 - Many ideas can be generalized outside two-player zero-sum
 - Correlated equilibria as link between prescriptive/descriptive views
- Let's "Go Wide": evaluate agents across many environments!
 - Pursuit of general agents requires general evaluation
 - Many efforts to help make this happen
 - "Never been a better time than right now" :)

Thank You! ... Any Questions?

Marc Lanctot

lanctot@deepmind.com

mlanctot.info/



