

Evolution of Societies via Reinforcement Learning

Yann Bouteiller, Karthik Soma, Giovanni Beltrame

Keywords: Multi-agent reinforcement learning, Evolutionary game theory, Simulation, Economy, Social dynamics, Opponent-learning awareness.

Summary

Diverse studies from computational neuroscience have found evidence of Reinforcement Learning (RL) driving learning in biological brains, amongst other learning protocols such as imitation. Furthermore, RL algorithms have now emerged as the core technique explicitly driving innovation in a growing number of industrial applications, including artificial language generation, drug discovery and finance. It is therefore natural to consider Multi-Agent Reinforcement Learning (MARL) as one of the revision protocols powering social dynamics. However, the mathematical complexity of this idea has prevented the development of a theory around it so far. In this paper, we hope to initiate this line of research by leveraging simulation. We develop computationally efficient implementations of two fundamental MARL revision protocols: "naive" Policy Gradient and Learning with Opponent-Learning Awareness (LOLA). These implementations enable us to conduct large-scale evolutionary simulations across classic interaction models from Evolutionary Game Theory. Our experiments yield various insights into how non-stationarity-aware learners affect the evolution of societies. In particular, we find that LOLA learners promote cooperation in the Stag Hunt model, delay cooperative outcomes in the Hawk-Dove model, and reduce strategy diversity in the Rock-Paper-Scissors model.

Contribution(s)

1. This paper introduces a methodology for analyzing the social dynamics that stem from MARL revision protocols in large populations.
Context: While Evolutionary Game Theory traditionally models biological evolution through fitness-based replication in pairwise interactions—an approach later extended to study imitation in social systems—our work explores how societies evolve when agents actively optimize their own fitness via continual learning. This framework is relevant for understanding the dynamics of competitive environments like economic markets.
2. We derive fast implementations of exact multi-agent Policy Gradient and exact Opponent-Learning Awareness targeting evolutionary simulations in pairwise matrix games.
Context: The scope of our paper is limited to non-repeated games with random agent pairing at each evolutionary step. This foundational work establishes a framework that future research can extend to more complex scenarios, including episodic MARL environments (such as the Iterated Prisoner's Dilemma) and structured assortment processes.
3. We conduct population simulations of 200,000 agents to analyze how naive learning and Opponent-Learning Awareness shape collective behavior in evolutionary settings.
Context: This work represents the first application of advanced MARL revision protocols in the context and scale of Evolutionary Game Theory simulations.

Abstract

The universe involves many independent co-learning agents as an ever-evolving part of our observed environment. Yet, in practice, Multi-Agent Reinforcement Learning (MARL) applications are typically constrained to small, homogeneous populations and remain computationally intensive. We propose a methodology that enables simulating populations of Reinforcement Learning agents at evolutionary scale. More specifically, we derive a fast, parallelizable implementation of Policy Gradient (PG) and Opponent-Learning Awareness (LOLA), tailored for evolutionary simulations where agents undergo random pairwise interactions in stateless normal-form games. We demonstrate our approach by simulating the evolution of very large populations made of heterogeneous co-learning agents, under both naive and advanced learning strategies. In our experiments, 200,000 PG or LOLA agents evolve in the classic games of Hawk-Dove, Stag-Hunt, and Rock-Paper-Scissors. Each game provides distinct insights into how populations evolve under both naive and advanced MARL rules, including compelling ways in which Opponent-Learning Awareness affects social evolution.

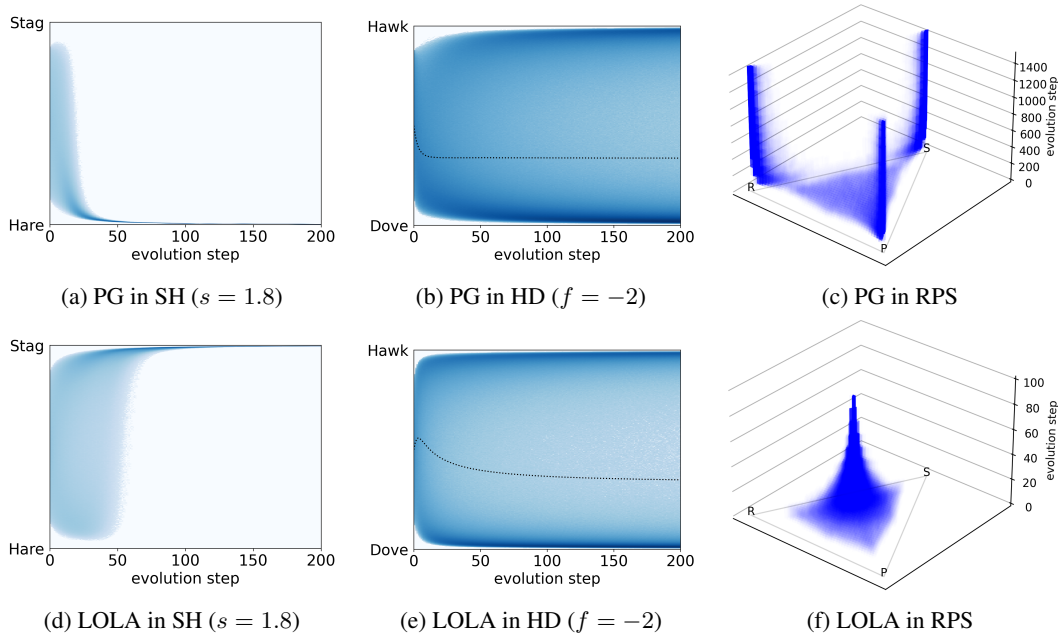


Figure 1: Populations of 200,000 RL agents evolving in the classic games of Stag Hunt, Hawk-Dove and Rock-Paper-Scissors (columns), via Policy Gradient and LOLA (rows). Each agent is a stochastic policy, represented as linear coordinates between pure strategies. Dark shades of blue indicate high concentrations of agents, and evolution steps correspond to one learning step performed per agent. In Hawk-Dove, black dots indicate the average policy over the entire population.

1 Introduction

Our universe is one of perpetual change, where countless agents co-exist, co-learn and co-evolve. From an individual agent’s perspective, other learning agents are a fundamentally non-stationary part of the environment, especially when incentives are in conflict (Papoudakis et al., 2019). In the realm of Reinforcement Learning (RL), the study of this type of possibly adversarial non-stationarity is a field known as Multi-Agent Reinforcement Learning (MARL). Some notable achievements do tackle MARL scenarios, but these are largely constrained by the high complexity of multi-agent

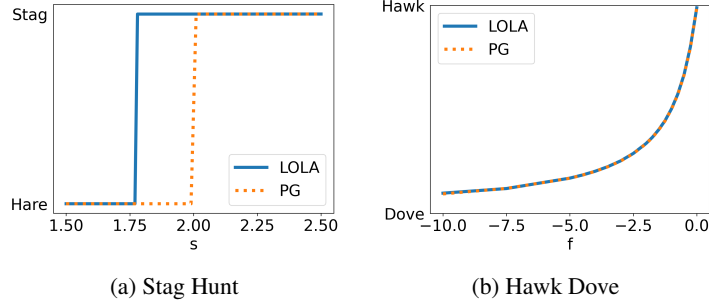


Figure 2: Final average policy over the population, depending on cost values.

training. In fact, these often resort to “self-play” (Silver et al., 2016; Berner et al., 2019), i.e., training a single neural network against one or few copies of itself, or to “centralized training” (Lowe et al., 2020; Yu et al., 2022), i.e., circumventing non-stationarity by using privileged global information.

We are interested in modeling plausible real-world social evolution processes that stem from RL revision protocols. We make the assumption that individuals continuously adapt their own strategies by following local RL rules, and we simulate how large societies evolve as a result. This setting is relevant to the study of social dynamics where individuals actively optimize their own fitness, such as economy and finance. More precisely, the scope of this paper is to simulate large populations of persistent RL agents, following an interaction model close to that of Evolutionary Game Theory (EGT). At each evolution step, agents are paired according to some assortment process (in our case, uniformly at random), and interact according to their respective policies. After each such interaction, agents adapt their policies according to their respective payoff and learning rule¹. Note that this type of social evolution is intrinsically driven by learning and individual preferences, in stark contrast to genetic evolution, which the literature usually considers as an extrinsic process happening among non-learning agents through their replication and death (Weibull, 1997). In our work, the concept of *fitness* does not necessarily refer to a measure of reproductive ability. Instead, it refers to a measure of expected completion of an arbitrary individual objective, i.e., an RL value.

We bring advanced intrinsic MARL protocols to the scale of evolutionary simulations, studying in particular the evolutionary effects of Opponent-Learning Awareness, an advanced MARL protocol able to take advantage of non-stationary dynamics (Foerster et al., 2017).

2 Preliminaries

This section reviews the concepts from Evolutionary Game Theory (EGT) and Multi-Agent Reinforcement Learning (MARL) used in this paper. A high-level overview of the literature contextualizing our work is provided in the Supplementary Material (Appendix A).

2.1 Evolutionary Game Theory (EGT)

EGT models evolution as a series of random pairwise interactions, where interactions are typically simple bi-matrix games (i.e., 2-player multi-armed bandits). Agents are sampled from a large population to be randomly paired and evaluated against their drawn opponent. The outcome of this interaction is a payoff for each opponent, whose expectation is called the agent’s *fitness* against the current population. An agent’s fitness depends on both its policy and the current configuration of the population. The agent’s policy is called its *type*, which is one of n possible types. In gene-inspired revision protocols, agents with a greater fitness replicate and thus tend to “invade” the population, whereas agents with a lower fitness tend to go “extinct”. In particular, EGT is interested in *evolu-*

¹Throughout this paper, we use vocabulary from RL and EGT interchangeably. In particular, “payoff” = “return”, “fitness” = “value”, “strategy” = “policy”, “pure strategy” = “action”, “revision protocol” = “learning rule”.

tionarily stable equilibria, which are configurations of the population where the different types are present in stable proportions under replication dynamics. In evolutionarily stable equilibria, the population configuration is robust to rare mutations, where few agents randomly switch from one type to another. In this paper, we will be studying more complex, learning-based evolutionary equilibria in three classic symmetric games: *Stag Hunt*, *Hawk-Dove*, and *Rock-Paper-Scissors*. These games can be described by bi-matrices:

	Stag	Hare		Hawk	Dove		Rock	Paper	Scissors
Stag	s, s	$0, 1$	Hawk	f, f	$2, 0$	Rock	$0, 0$	$-1, 1$	$1, -1$
Hare	$1, 0$	$1, 1$	Dove	$0, 2$	$1, 1$	Paper	$1, -1$	$0, 0$	$-1, 1$
						Scissors	$-1, 1$	$1, -1$	$0, 0$

(a) Stag Hunt (b) Hawk-Dove (c) Rock-Paper-Scissors

where rows represent the action chosen by the ego agent (bold) with corresponding payoffs in first position, and columns represent the action chosen by the other agent with corresponding payoffs in second position.

Stag hunt (SH) is a 2-action game illustrating the evolution of cooperation. In this game, agents need to hunt for food and choose to either go for a Stag, or go for a Hare. Hunting a Hare is easy: any agent choosing this option successfully receives a payoff of 1. Hunting a Stag is harder: both agents need to cooperate, otherwise the agent choosing to go for a Stag fails to catch anything and receives a payoff of 0. However, if both agents cooperate, they succeed and each receives a payoff of $s > 1$, which is better than going for Hares. The game of Stag hunt has two distinct pure strategy Nash equilibria²: (1) both agents always playing Stag, and (2) both agents always playing Hare.

Hawk-Dove (HD) is a 2-action game illustrating the evolution of conflict over shareable resources. Agents either choose to act as a “Hawk” or as a “Dove”. When a Dove encounters another Dove, they share the available food (each receives a payoff of 1). When a Dove encounters a Hawk, it yields and gets no food (payoff of 0) while the Hawk gets all of it (payoff of 2). But when a Hawk encounters another Hawk, they fight and both get injured (payoff of $f < 0$). The game of Hawk-Dove has two pure strategy Nash equilibria: (1) Hawk-Dove and (2) Dove-Hawk. But note that in these equilibria, Doves have a smaller payoff than Hawks. In the context of evolutionary genetics, this means that when Doves encounter almost only Hawks, they move toward extinction as Hawks invade. However, when Hawks encounter almost always Hawks, their expected payoff is even less than Doves encountering Hawks, and thus Hawks move toward extinction as Doves invade. In other words, there are population configurations in which it is not more relevant to be a Hawk than a Dove in terms of fitness, and replication dynamics naturally drive the population there.

Rock-Paper-Scissors (RPS) is a 3-action zero-sum³ game. It illustrates more complex situations where there are cycles in the preferences over actions. Scissors beats Paper, Paper beats Rock, and Rock beats Scissors. RPS has a single mixed-Nash equilibrium, where both players choose their actions uniformly at random. Similarly to the HD game, all populations whose average behavior is this equilibrium are neutrally stable under replication dynamics and drifting due to random mutations.

A population whose individuals are distributed amongst n types can be represented as a *population vector* $P \in \mathbb{R}^n$ whose components $0 \leq p_i \leq 1$ sum to 1 and represent the proportion of type i . Under the “imitation of the fittest” revision protocol (as well as several other bio-inspired revision protocols), large populations are known to follow a famous population dynamic over time (t), called the *Replicator Dynamic*:

$$\frac{dp_i}{dt} = p_i(v_i(P) - \bar{v}(P)) \quad (1)$$

where $v_i(P)$ is the fitness of type i in the population, and $\bar{v}(P)$ is the average fitness of all agents in the population. Denoting the vector of v_i ’s as Q , the vector of all ones as $\mathbf{1}$ and the Hadamard

²2-player equilibria where each agent always selects the same action.

³The sum of the two agents’ payoffs is always 0.

product⁴ as \odot , Equation 1 can be written in matrix form:

$$\frac{d}{dt}P = P \odot (Q - \mathbf{1}\bar{v}), \quad (2)$$

2.2 Multi-Agent Reinforcement Learning

While EGT traditionally focuses on populations where non-learning agents reproduce based on their fitness, our work examines how populations evolve when agents actively learn and adapt their strategies. Similar to EGT, we model social dynamics as a series of random pairwise interactions. Therefore, we are principally interested in 2-agent learning rules. In this paper, we will be specifically looking at two important such learning rules: *policy gradient* (PG), also referred to as “naive learning” in the MARL literature, and *learning with opponent-learning awareness* (LOLA) in its full form (i.e., using all terms from the first-order Taylor expansion).

Policy gradient (PG) is a fundamental learning rule from single-agent RL. It consists of following the first-order gradient of the value function with respect to the ego agent’s policy parameters. Let us consider a pair of agents. We denote the ego agent as agent 1, and the other agent as agent 2. Let us further denote their respective policies as π_1 and π_2 , parameterized by vectors θ_1 and θ_2 , of current values v_1 and v_2 . The naive policy gradient is:

$$\nabla_{\theta_1} v_1(\theta_1, \theta_2) \quad (3)$$

The reason why following this gradient is considered naive in MARL is that this does not take into account the non-stationarity introduced by the learning process of agent 2.

Learning with opponent learning awareness (LOLA) is an improved version of PG that takes into account the learning process of the other agent. More precisely, LOLA models the learning process of agent 2 as if agent 2 were a naive learner and differentiates through its learning step:

$$\nabla_{\theta_1} v_1(\theta_1, \theta_2 + \Delta\theta_2) \quad (4)$$

where $\Delta\theta_2 = \eta \nabla_{\theta_2} v_2(\theta_1, \theta_2)$ is the naive learning step of agent 2, η being its learning rate. LOLA approximates this gradient by the following first-order Taylor expansion:

$$\begin{aligned} \nabla_{\theta_1} v_1(\theta_1, \theta_2 + \Delta\theta_2) &\approx \nabla_{\theta_1} (v_1 + (\Delta\theta_2)^\top \nabla_{\theta_2} v_1) \\ &= \underbrace{\nabla_{\theta_1} v_1}_{\text{PG}} + \underbrace{\eta (\nabla_{\theta_1} \nabla_{\theta_2} v_1)^\top \nabla_{\theta_2} v_2 + \eta (\nabla_{\theta_1} \nabla_{\theta_2} v_2)^\top \nabla_{\theta_2} v_1}_{\text{opponent-learning awareness}} \end{aligned} \quad (5)$$

Differentiating through the learning step of the opponent has an important advantage in our discussion: it is a naturally plausible way of predicting non-stationarity (assuming we maintain an internal model of others) in order to adapt beforehand and actively steer it toward our own incentives.

2.3 Population-policy equivalence

In a population of agents playing only pure strategies, uniformly sampling agents is equivalent to sampling actions from the abstract stochastic policy defined by the probability vector P . Thus, Equation 2 can be viewed as a learning process, albeit at the population level, where the evolving population of non-learning agents P is itself a self-play learning agent (Bloembergen et al., 2015).

3 Methods

To model how learning affects societies, we adopt a philosophy similar to EGT. Namely, we consider large populations of independent learning agents, which are paired randomly at each evolution

⁴The matrix product takes precedence over the Hadamard product in all our notations.

iteration and interact in normal-form matrix games. Each agent has its own learning rule (i.e., either of the two presented in Section 2.2) that it applies to its own policy after each pairwise interaction. Whereas MARL usually thinks about these rules in the context of persistent interactions between fixed pairs of agents, in the evolutionary setting, they instead get applied after single interactions between random pairs. In other words, learning agents consider that the random opponent they interact with at each evolution step is a representative sample of the population and, for LOLA, of the current direction of its non-stationary dynamics.

3.1 Policy architecture

From an RL perspective, the normal-form matrix games presented in Section 2.1 are 2-agent multi-armed bandits. As this is a common assumption in multi-armed bandits (Sutton, 2018), we consider the policy architecture parameterized by the preference vector $\theta \in \mathbb{R}^n$, where n is the number of actions, projected to the probability simplex by a simple softmax function σ , which yields the probability vector $P \in \mathbb{R}^n$ of the policy selecting each of the n available actions:

$$P = \sigma(\theta) \quad (6)$$

This policy architecture has a useful property for our derivations: its gradient has a symmetric analytical form, which is

$$\nabla_{\theta} P = (\nabla_{\theta} P)^{\top} = \text{diag}(P) - PP^{\top} \quad (7)$$

3.2 Analytical Policy Gradient

Let us consider a symmetric normal-form game with n actions, played by a pair of agents denoted as agents 1 and 2. Since the game is symmetric, we can represent its bi-matrix as a single matrix $A \in \mathbb{R}^{n \times n}$, valid from the perspective of both agents⁵. Let us further assume that the policies of both agents are parameterized by $\theta_{1,2} \in \mathbb{R}^n$, with the simple policy architecture described in Equation 6:

$$P_{1,2} = \sigma(\theta_{1,2}) \quad (8)$$

The value functions of both agents are:

$$v_1 = P_1^{\top} A P_2 ; v_2 = P_2^{\top} A P_1 \quad (9)$$

Thus, the “naive” Policy Gradient of agent 1’s value with respect to agent 1’s parameters is:

$$\nabla_{\theta_1} v_1 = P_1 \odot (Q_1 - \mathbf{1}v_1)$$

where $Q_1 \in \mathbb{R}^n$ is the vector of action-values of the n available actions (derivation in Appendix B).

As a side note, this draws an interesting parallel with Equation 2: the PG update on the parameter vector θ_1 is the same as the Replicator update on the probability vector $P_1 = \sigma(\theta_1)$. In other words, each individual PG agent can itself be seen as an evolving population of abstract non-learning agents playing pure strategies, similarly to the equivalence noted in Bloembergen et al. (2015).

This formulation yields the following analytical PG formulation for symmetric normal-form games:

$$\nabla_{\theta_1} v_1 = P_1 \odot (I - \mathbf{1}P_1^{\top})AP_2 \quad (10)$$

$$\nabla_{\theta_2} v_2 = P_2 \odot (I - \mathbf{1}P_2^{\top})AP_1 \quad (11)$$

which only involves simple matrix operations, and therefore is a fast, easily parallelizable implementation.

⁵ A is formed by the first entries of the corresponding bi-matrix in Section 2.1.

3.3 Analytical LOLA

Similarly, we derive the following analytical formulation of the LOLA gradient in the symmetric normal-form game defined by matrix A (the derivation is provided in Appendices C and D):

$$\begin{aligned} \nabla_{\theta_1} v_1(\theta_1, \theta_2 + \Delta\theta_2) &\approx P_1 \odot X_1 A P_2 \\ &\quad + \eta(T^\top \odot X_1 A X_2^\top)(P_2 \odot X_2 A P_1) \\ &\quad + \eta(T^\top \odot X_1 A^\top X_2^\top)(P_2 \odot X_2 A^\top P_1) \end{aligned} \quad (12)$$

where $X_1 := I - \mathbf{1}P_1^\top$, $X_2 := I - \mathbf{1}P_2^\top$, and $T := P_2P_1^\top$. As for Policy Gradient, this formulation only involves simple matrix operations and is straightforward to parallelize.

3.4 Batched pairwise bandits

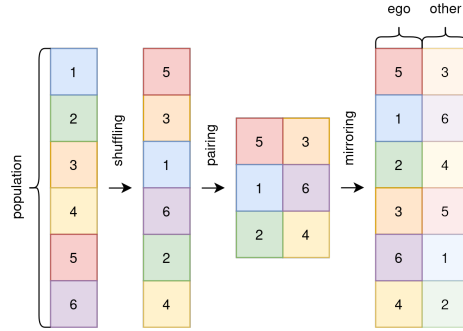


Figure 3: Pairing and batching

Equations 10 and 12 are fast, backprop-free implementations of exact PG and exact LOLA for normal-form games. Still, it would be prohibitively slow to apply these updates iteratively on single agent pairs in evolutionary-scale simulations. At each *evolution step*, all agents go through their revision protocol once. To make this process scalable, we batch learning updates across the entire population. More precisely, at each evolution step, we shuffle the entire population and randomly pair all agents two-by-two. Since all interactions are pairwise in evolutionary simulations, this enables batching updates across agent pairs. To optimize the process even further, we mirror all pairs to perform the entire population update in one single batched operation. This batching procedure is illustrated in Figure 3. We found batching populations in this manner to be extremely efficient. Combining this procedure with the analytical PG and LOLA implementations described in previous Sections, we are able to simulate populations of 200,000 learning agents for thousands of evolution steps in a matter of seconds on a consumer-grade GPU⁶. To ensure reproducibility and foster future work in this direction, we open-source our implementation⁷.

4 Experiments

We use the normal-form matrices presented in Section 2.1 as our three tested interactions: *Stag hunt* (SH), *Hawk-Dove* (HD) and *Rock-Paper-Scissors* (RPS). Each individual agent has its own persistent learning rule: either PG (gradient descent on Equation 10) or LOLA (gradient descent on Equation 12). The learning rate has no relevant impact on the dynamics presented in this paper other than varying their speed, thus we use unit learning rates everywhere.



Figure 4: Duration of a full evolution step (lower is better)

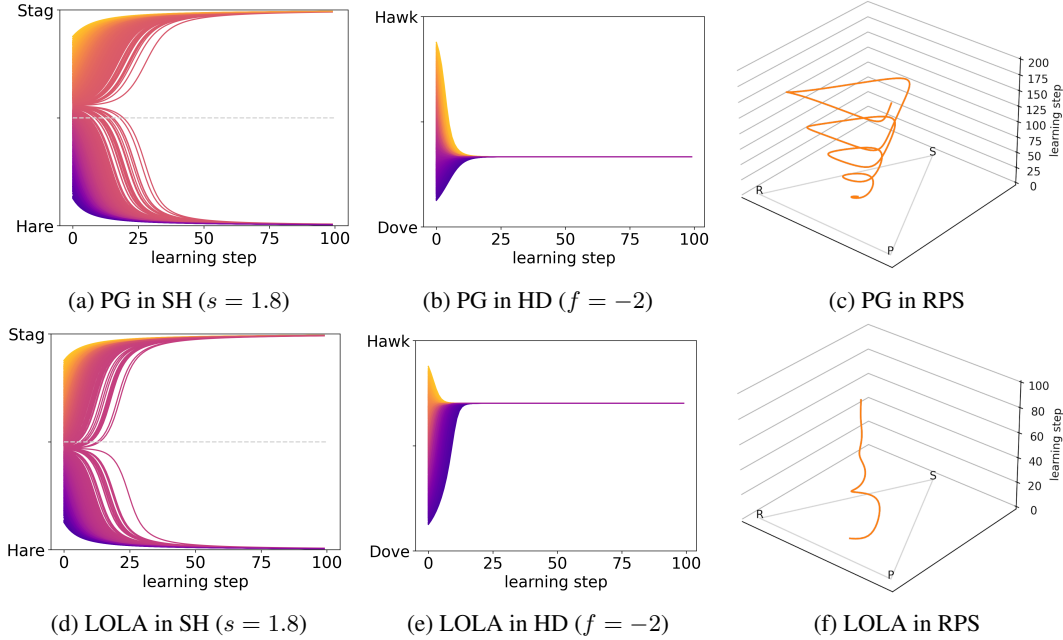


Figure 5: Self-play. In SH and HD, the color marks the initial policy.

4.1 Scalability

Figure 4 reports the computational performance of our approach (“matrix batched LOLA”), compared to two baselines. The “autograd iterative LOLA” baseline reproduces how LOLA updates are usually performed in classical MARL scenarios: using PyTorch’s autograd to compute the LOLA gradient, and updating the policies of all agents in an iterative fashion. This baseline is clearly not a viable implementation for evolutionary-scale simulations and is only provided for illustration. On the other hand, our “autograd batched LOLA” baseline is of more interest for future work. While the “matrix batched” approach is faster, it is limited to single-shot multi-armed bandits⁸. In particular, the “matrix batched” approach does not allow episodic interactions. Therefore, we have implemented the approach of Section 3.4 along with autograd, which yields a potentially more general implementation. We present the performance of this alternative approach as “autograd batched LOLA”. Clearly, the main reason why we can perform these large-scale simulations at all is that the pairwise structure of EGT simulations enables batching interactions and policy updates. Since our matrix-based implementation is faster for normal-form games, we use it in the following.

⁶All experiments in this paper are conducted with an i7-12700H CPU, an RTX 3080 Ti GPU, and 64G of RAM.

⁷<URL hidden for blind peer review>

⁸Section 3 is possible because the value function has a straightforward formulation in 2-agent multi-armed bandits.

4.2 Expected results

The methodology proposed in Section 3 is limited to very simple, single-shot bandit interactions between random pairs of learning agents. Remember how, under the population-policy equivalence described in Section 2.3, a population of pure-strategy agents can equivalently be seen as a stochastic policy over types. In our setting, agents have full stochastic policies assigning non-zero probabilities to all available actions, but are still simple stateless multi-armed bandits. Uniformly sampling a random pair of agents from such a population and then sampling from their policies is equivalent in expectation to sampling two actions from the average policy of all population members. In other words, at least in expectation, we anticipate that the population dynamic will show some resemblance to self-play over the population’s average policy (represented in Figure 5).

4.3 Empirical results

Stag Hunt. Figures 5a and 5d show how a single self-play agent learns against itself in Stag Hunt, via naive Policy Gradient and LOLA, respectively. The vertical axis represents its policy, and the horizontal axis represents time expressed in learning steps. Policies are color-coded by their initial configuration, with yellow policies starting close to the deterministic Stag policy and purple policies starting close to the deterministic Hare policy. Notably, PG tends to converge to the individualistic Nash equilibrium (i.e., deterministic Hare) for most initial configurations, whereas LOLA tends to converge to the pro-social Nash equilibrium (i.e., deterministic Stag). Notice that the forks are on different sides of the uniform random policy (middle tick), which is important because we will initialize all our population experiments with a Gaussian distribution around this neutral policy. As explained in Section 4.2, we expect the final population dynamic to follow a similar pattern. Figures 1a and 1d display the results of our first evolutionary simulation, featuring a population of 200,000 learning agents in Stag Hunt. Dark shades of blue represent high concentrations of agents. An *evolution step* corresponds to one learning step performed per agent in the population. In Figure 1a, the population is exclusively formed of naive learners, and quickly converges to the individualistic policy, as predicted by Figure 5a. In Figure 1d, the population is formed of only LOLA agents. Contrary to the naive population, a population of non-stationarity aware learners such as LOLA evolves to unanimously adopt the superior pro-social equilibrium (i.e., deterministic Stag). This effect is modulated by the payoff of the pro-social strategy, as measured in Figure 2a. Appendix E further shows that, when enough opponent-learning aware learners are present in a mixed population, they become able to pull naive learners toward the pro-social strategy.

Hawk-Dove. Figures 5b and 5e show that, in Hawk-Dove, self-play converges to definite policies regardless of where training starts from. Naive learning (Figure 5b) converges to the mixed Nash equilibrium. However, LOLA (Figure 5e) converges to another, inferior policy, where it selects Hawk 70% of the time (which yields a smaller payoff for both players, and is not a Nash equilibrium). A similar behavior has been described as “arrogance” in Letcher et al. (2018), where both LOLA learners make wrong assumptions about the response of their opponent and thus pull away from the equilibrium. From these observations, one could imagine that all learning agents in the population would converge to these policies, similar to what we observed for SH, but this is not at all what happens in practice. In HD, whether naive learning (Figure 1b) or LOLA (Figure 1e) is used as the learning rule of the entire population, it evolves into a mix of Hawks and Doves, most of them with close-to-deterministic policies, and in both cases with an average population policy that corresponds to the mixed Nash equilibrium. Notice that the convergence to deterministic strategies is however much slower than what we observed for SH, and it is in fact not clear whether this will eventually happen entirely, even after 10,000 steps (Appendix F). Nonetheless, what can be observed from Figures 1b and 7a is that LOLA learners converge faster to deterministic policies during early steps (shades of blue) but the population takes longer to stabilize (dotted black). In additional experiments, we mixed LOLA and PG learners to see whether LOLA learners would be more inclined toward the Hawk strategy (as suggested by Figure 5e). However, these experiments invalidated this hypothesis: half LOLA learners and half PG learners were present amongst both final sub-populations of Hawk-inclined and Dove-inclined individuals.

Rock-Paper-Scissors. Our final population experiment takes place in the 3-action Rock-Paper-Scissors environment, used in EGT to explain the coexistence of competitively unbalanced species (Allesina & Levine, 2011). Interestingly, we will see that, while populations of naive learners agree with this explanation of sustainable diversity, populations of LOLA learners yield quite the opposite result. Similarly to previous Sections, Figures 5c and 5f show how 2-agent self-play behaves for both PG and LOLA. Only 1 policy is displayed for readability (other initial conditions yield similar effects). The triangle on the bottom of each plot represents the policy, and the vertical axis represents the number of learning steps. It can be seen that PG slowly spirals outward from the mixed Nash equilibrium (due to performing straight policy updates with a non-zero learning rate following a circular vector field), whereas LOLA quickly spirals inward until it reaches the mixed Nash equilibrium. Figures 1c and 1f present the results of our evolutionary simulations in the RPS game. The color-code follows the same principle as our previous population plots, with the bottom triangle being the policy, and evolution steps being the vertical dimension. Similarly to Figure 1b, Figure 1c shows that populations of naive learners evolve into 3 equally distributed groups of close-to-deterministic agents always playing Rock, Paper and Scissors respectively. The reason why this happens is of interest, and is more clearly understood from Figure 9 in Appendix G. In short, this dynamic results from the loss of plasticity modeled by the softmax function⁹. At the beginning of evolution, naive learners are erratically moving around as they encounter all types of strategies. However, after some time, agents get “trapped” near the border of the policy simplex, where gradients toward the opposite action are near-zero. After a long time, three groups of near-deterministic agents emerge, and a small number of them continuously escape toward the strategy that counters the majority group, which eventually creates a new majority, and so on, yielding a cyclic evolution pattern. In other words, diversity emerges from populations of naive learners in the RPS model. On the other hand, Figure 1f tells the opposite story about populations of LOLA agents, which instead quickly and unanimously converge to the mixed Nash equilibrium of this game.

4.4 Remark

From our simulation results, it looks like the mean policy averaged over the entire population always converges near a Nash equilibrium of the game, even when the learning rule itself does not converge to this equilibrium in the conventional 2-agent MARL setting (see Figures 5e and 5c). This property is however merely a consequence of the uniform random opponent matching scheme that we chose to implement in this paper. For instance, let us consider an extreme opposite scheme, where all pairs would instead interact persistently. All individual pairs would then converge to the pure Nash equilibrium in the Hawk-Dove game (that is, exactly one deterministic Hawk and one deterministic Dove per pair): this would average to a uniform random policy, as opposed to Figures 1b and 1e. In reality, partner selection is more structured (Anastassacos et al., 2020) and can lead to different outcomes, which we plan to explore in the future.

5 Conclusions

We have presented a methodology enabling large-scale evolutionary simulations of independent learning agents, for both naive (Policy Gradient) and advanced non-stationarity-aware (LOLA) learning rules. We have demonstrated the scalability of our approach by performing very-large-scale evolutionary simulations of 200,000 independent learning agents, interacting in the classic games of Stag Hunt, Hawk-Dove and Rock-Paper-Scissors. Our work essentially explores the effect of Multi-Agent Reinforcement Learning on the usual model of evolution adopted in Evolutionary Game Theory, and demonstrates compelling social dynamics originating from both naive and non-stationarity-aware learners. While the approach presented in this paper is specifically designed for normal-form matrix games (i.e., stateless 2-player multi-armed bandits), exploring stateful episodic scenarios with a gradient-based approach might be possible, and is an avenue for future work.

⁹A model very similar to the “cost of motion” described by Mertikopoulos & Sandholm (2018).

References

- Stefano Allesina and Jonathan M Levine. A competitive network theory of species diversity. *Proceedings of the National Academy of Sciences*, 108(14):5638–5642, 2011.
- Nicolas Anastassacos, Stephen Hailes, and Mirco Musolesi. Partner selection for the emergence of cooperation in multi-agent systems using reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:7047–7054, 04 2020. DOI: 10.1609/aaai.v34i05.6190.
- Jose Apesteguia, Steffen Huck, and Jörg Oechssler. Imitation—theory and experimental evidence. *Journal of Economic Theory*, 136(1):217–235, 2007.
- Robert Axelrod and William D Hamilton. The evolution of cooperation. *science*, 211(4489):1390–1396, 1981.
- Alan JS Beavan, Maria Rosa Domingo-Sananes, and James O McInerney. Contingency, repeatability, and predictability in the evolution of a prokaryotic pangenome. *Proceedings of the National Academy of Sciences*, 121(1):e2304934120, 2024.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębniak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Daan Bloembergen, Karl Tuyls, Daniel Hennes, and Michael Kaisers. Evolutionary dynamics of multi-agent learning: A survey. *Journal of Artificial Intelligence Research*, 53:659–697, 2015.
- Richard Dawkins. The selfish gene, 1976.
- Jakob Foerster, Gregory Farquhar, Maruan Al-Shedivat, Tim Rocktäschel, Eric Xing, and Shimon Whiteson. DiCE: The infinitely differentiable Monte Carlo estimator. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1529–1538. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/foerster18a.html>.
- Jakob N Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness. *arXiv preprint arXiv:1709.04326*, 2017.
- The Anh Han, Flavio L Pinheiro, Simon T Powers, Julián García, and Matthijs van Veelen. No strategy can win in the repeated prisoner’s dilemma: Linking game theory and computer simulations. *Frontiers in Robotics and AI* | www.frontiersin.org, 1:102, 2018. DOI: 10.3389/frobt.2018.00102. URL www.frontiersin.org.
- Daniel Hennes, Dustin Morrill, Shayegan Omidshafiei, Rémi Munos, Julien Perolat, Marc Lanctot, Audrunas Gruslys, Jean-Baptiste Lespiau, Paavo Parmas, Edgar Duéñez-Guzmán, et al. Neural replicator dynamics: Multiagent learning via hedging policy gradients. In *Proceedings of the 19th international conference on autonomous agents and multiagent systems*, pp. 492–501, 2020.
- Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. Population based training of neural networks. *arXiv preprint arXiv:1711.09846*, 2017.
- Alistair Letcher, Jakob Foerster, David Balduzzi, Tim Rocktäschel, and Shimon Whiteson. Stable opponent shaping in differentiable games. *arXiv preprint arXiv:1811.08469*, 2018.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments, 2020. URL <https://arxiv.org/abs/1706.02275>.
- Chris Lu, Timon Willi, C. S. D. Witt, and Jakob Nicolaus Foerster. Model-free opponent shaping. In *International Conference on Machine Learning*, 2022. URL <https://api.semanticscholar.org/CorpusID:248505991>.

- Andrei Lupu and Doina Precup. Gifting in multi-agent reinforcement learning. In *Proceedings of the 19th International Conference on autonomous agents and multiagent systems*, pp. 789–797, 2020.
- Panayotis Mertikopoulos and William H. Sandholm. Riemannian game dynamics, 2018. URL <https://arxiv.org/abs/1603.09173>.
- Martin A Nowak. Five rules for the evolution of cooperation. *science*, 314(5805):1560–1563, 2006.
- Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland, Jean-Baptiste Lespiau, Wojciech M Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos. α -rank: Multi-agent evaluation by evolution. *Scientific reports*, 9(1):9937, 2019. DOI: 10.1038/s41598-019-45619-9. URL <https://doi.org/10.1038/s41598-019-45619-9>.
- Georgios Papoudakis, Filippos Christianos, Arrasy Rahman, and Stefano V. Albrecht. Dealing with non-stationarity in multi-agent deep reinforcement learning, 2019. URL <https://arxiv.org/abs/1906.04737>.
- William H Sandholm. *Population games and evolutionary dynamics*. MIT press, 2010.
- Francisco C Santos and Jorge M Pacheco. Scale-free networks provide a unifying framework for the emergence of cooperation. *Physical review letters*, 95(9):098104, 2005.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- J Maynard Smith and George R Price. The logic of animal conflict. *Nature*, 246(5427):15–18, 1973.
- Richard S Sutton. Reinforcement learning: an introduction. *A Bradford Book*, 2018.
- Karl Tuyls and Ann Nowé. Evolutionary game theory and multi-agent reinforcement learning. *The Knowledge Engineering Review*, 20(1):63–90, 2005.
- Jörgen W Weibull. *Evolutionary game theory*. MIT press, 1997.
- Timon Willi, Johannes Treutlein, Alistair Letcher, and Jakob Nicolaus Foerster. Cola: Consistent learning with opponent-learning awareness. *ArXiv*, abs/2203.04098, 2022. URL <https://api.semanticscholar.org/CorpusID:247315224>.
- Haoxiang Xia, Huili Wang, and Zhaoguo Xuan. Opinion dynamics: A multidisciplinary review and perspective on future research. *International Journal of Knowledge and Systems Science (IJKSS)*, 2(4):72–91, 2011.
- Jiachen Yang, Ang Li, Mehrdad Farajtabar, Peter Sunehag, Edward Hughes, and Hongyuan Zha. Learning to incentivize other learning agents. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS ’20, Red Hook, NY, USA, 2020a. Curran Associates Inc. ISBN 9781713829546.
- Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field multi-agent reinforcement learning, 2020b. URL <https://arxiv.org/abs/1802.05438>.
- Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative, multi-agent games, 2022. URL <https://arxiv.org/abs/2103.01955>.

Supplementary Materials

The following content was not necessarily subject to peer review.

A Related literature

This paper is a first step toward describing the fast-paced, social evolution that stems from individuals actively optimizing their own fitness in real-world societies and economies. While the idea is probably not novel (Tuyls & Nowé, 2005), no scientific progress has been made in developing a theory around it so far. In fact, Evolutionary Game Theory (EGT) usually avoids MARL entirely, as it introduces a substantial amount of theoretical complexity when modeling population dynamics and does not relate to the stochasticity-driven model of genetic evolution accepted so far (a model however recently challenged by the work of Beavan et al. (2024)). Before our work, MARL was complex to simulate at evolutionary scale and very few papers have attempted anything similar. Thus, the aim of this Section is to motivate our line of work by providing a brief overview of the most closely related fields of research and positioning ourselves with respect to their literature.

Evolutionary Game Theory. Originally rooted in biology, EGT extends classical Game Theory to the study of evolving populations. Smith & Price (1973) laid its mathematical foundations and introduced the concept of Evolutionarily Stable Strategies, formalizing how different types are maintained in genetic evolution. Axelrod & Hamilton (1981) famously described how cooperation naturally arises in the framework of EGT. More recently, Nowak (2006) proposed five mechanisms fostering the evolution of cooperation: kin selection, direct reciprocity, indirect reciprocity, network reciprocity, and - an idea that has a long history of controversy amongst evolutionary biologists - group selection. Finally, Mertikopoulos & Sandholm (2018) extended the mathematical framework available for studying EGT, by formalizing population dynamics (such as the Replicator Dynamic) under the more general class of Riemannian game dynamics.

Opinion dynamics. In parallel to its original motivation in biology, EGT has gathered interest from different fields, especially economy and social dynamics. In an influential book promoting the “selfish” gene-centered view of evolution, the ethologist Dawkins (1976) speculated about “memes”, an alleged generalization of genetics to cultural dynamics. Amongst other assumptions, his idea was based on the belief that *replication of the fittest*, the gradient-free revision protocol through which populations of non-learning agents are thought to evolve in nature (Apesteguia et al., 2007), extends beyond biology and, in particular, to opinions and strategies. Because this belief seems relevant to situations where individuals learn via imitation, it has later inspired a large body of work under the name “imitation dynamics” (Sandholm, 2010; Xia et al., 2011). Mentioning this competing line of thought is interesting, as our work studies the consequences of fundamentally different assumptions. Namely, whereas genetic evolution arguably stems from replication dynamics and random mutations (which may also be the case for social strategies/opinions to some extent), the dynamics that we specifically study in this paper instead stem from continual learning. We do not consider any replication process, but only a gradient-informed, learning-based “mutation” process.

Population-Based Training (PBT) (Jaderberg et al., 2017) is a line of efficient bio-inspired MARL approaches illustrating the potential complementarity of both views. In PBT, two processes coexist: an inner training loop lets a group of agents learn via MARL rules, while an outer genetic loop selects and replicates the fittest few resulting policies.

Opponent shaping. Advanced MARL rules are able to take advantage of the non-stationarity stemming from the learning processes of other agents. For instance, an agent can shape its opponent by learning to share its own rewards (Lupu & Precup, 2020; Yang et al., 2020a). However, this involves learning to reward opponents when they exhibit favorable behaviors, which exacerbates non-stationarity as credit assignment gets harder. Opponent-learning awareness methods such as LOLA overcome this issue by instead differentiating through the learning step of other agents (Foerster et al., 2017). An extensive line of research recently emerged in this direction, attempting to

fix diverse issues in the original LOLA formulation (Foerster et al., 2018; Letcher et al., 2018; Lu et al., 2022; Willi et al., 2022). In this paper, we find a way to simulate LOLA in large populations.

Simulation. Due to the high theoretical complexity of MARL-based population dynamics, our paper focuses on empirical simulation. Even when using simple replication-based revision protocols, resorting to simulation is common in EGT. For example, Santos & Pacheco (2005) famously simulated how structured locality fosters cooperation. Han et al. (2018) confirmed various theoretical predictions regarding cycles and their effect in games by simulating the evolution of non-learning strategies in iterated prisoner’s dilemmas. Interestingly, EGT simulations have also been used by Omidshafiei et al. (2019), who proposed to simulate imitation of the fittest to rank policies within a population of readily-trained agents. Finally, one work of particular relevance to ours has been conducted by Yang et al. (2020b), who used mean-field theory in a classical MARL scenario to reduce the environment dimensionality. In their proposed approach, naive learners approximate neighboring agents as one single, "mean-field" opponent. This essentially transforms n -agent MARL into pairwise MARL, thus reducing the environment complexity from the point of view of individual agents.

Parallels between EGT and single-agent RL. The population-policy equivalence described by Bloembergen et al. (2015) yields interesting parallels between single-agent Reinforcement Learning and the Replicator Dynamic. In particular, the resemblance of Policy Gradient with the Replicator Dynamic noted in Section 3.2 was further studied by Hennes et al. (2020). From this observation, they derived a single-agent algorithm that bypasses the loss of plasticity introduced by the softmax architecture of Equation 6. However, beyond the fact that their line of work uses concepts from RL and EGT, it is mostly unrelated to ours and we cite it here to clear a confusion made by early readers of our work: they are interested in finding high-performance single-agent RL algorithms, whereas we are interested in characterizing the population dynamics that stem from MARL revision protocols when individuals actively learn in massively multi-agent, evolutionary settings.

B Policy gradient

We derive an analytical formulation of the PG update in symmetric normal-form games:

$$\begin{aligned}
 \nabla_{\theta_1} v_1 &= \nabla_{\theta_1} P_1^\top A P_2 \\
 &= (\text{diag}(P_1) - P_1 P_1^\top) A P_2 \\
 &= \text{diag}(P_1) A P_2 - P_1 v_1 \\
 &= P_1 \odot (A P_2 - v_1 \mathbf{1}) \\
 &= P_1 \odot (A P_2 - \mathbf{1} P_1^\top A P_2) \\
 &= P_1 \odot (Q_1 - \mathbf{1} v_1)
 \end{aligned}$$

C LOLA

We now derive an analytical formulation of the LOLA update in symmetric normal-form games, similar to what we found for PG in Section 3.2. We are missing three terms from Equation 5:

- $\nabla_{\theta_2} v_1$
- $\nabla_{\theta_1} \nabla_{\theta_2} v_1$
- $\nabla_{\theta_1} \nabla_{\theta_2} v_2$

To compute the first term, we note that $v_1 = P_1^\top A P_2$ is a scalar and thus can also be written $v_1 = P_2^\top A^\top P_1$. We can then compute this term similarly to PG:

$$\begin{aligned}
\nabla_{\theta_2} v_1 &= \nabla_{\theta_2} P_2^\top A^\top P_1 \\
&= (\text{diag}(P_2) - P_2 P_2^\top) A^\top P_1 \\
&= P_2 \odot (A^\top P_1 - v_1 \mathbf{1}) \\
&= P_2 \odot (A^\top P_1 - \mathbf{1} P_2^\top A^\top P_1) \\
&= P_2 \odot (I - \mathbf{1} P_2^\top) A^\top P_1
\end{aligned} \tag{13}$$

Computing the two remaining terms is also possible.

Let us start with $\nabla_{\theta_1} \nabla_{\theta_2} v_2$:

$$\nabla_{\theta_1} \nabla_{\theta_2} v_2 = \nabla_{\theta_1} \nabla_{\theta_2} P_2^\top A P_1 \tag{14}$$

$$\begin{aligned}
&= \nabla_{\theta_1} (\text{diag}(P_2) - P_2 P_2^\top) A P_1 \\
&= (\text{diag}(P_2) - P_2 P_2^\top) A (\text{diag}(P_1) - P_1 P_1^\top)
\end{aligned} \tag{15}$$

While it would already possible to implement this formulation, we further derive a more efficient implementation in Appendix D:

$$\nabla_{\theta_1} \nabla_{\theta_2} v_2 = T \odot (I - \mathbf{1} P_2^\top) A (I - P_1 \mathbf{1}^\top) \tag{16}$$

where $T := P_2 P_1^\top$.

Computing $\nabla_{\theta_1} \nabla_{\theta_2} v_1$ is fairly straightforward:

$$\begin{aligned}
\nabla_{\theta_1} \nabla_{\theta_2} v_1 &= \nabla_{\theta_1} \nabla_{\theta_2} P_2^\top A^\top P_1 \\
&= \nabla_{\theta_1} \nabla_{\theta_2} P_2^\top B P_1 & (B := A^\top) \\
&= T \odot (I - \mathbf{1} P_2^\top) B (I - P_1 \mathbf{1}^\top) & (\text{c.f. 14,16}) \\
&= T \odot (I - \mathbf{1} P_2^\top) A^\top (I - P_1 \mathbf{1}^\top)
\end{aligned} \tag{17}$$

$$= T \odot (I - \mathbf{1} P_2^\top) A^\top (I - P_1 \mathbf{1}^\top) \tag{18}$$

Substituting Equations 10, 11, 13, 16 and 18 in Equation 5 yields the following analytical formulation of the LOLA gradient in the symmetric normal-form game defined by matrix A :

$$\begin{aligned}
\nabla_{\theta_1} v_1(\theta_1, \theta_2 + \Delta \theta_2) &\approx P_1 \odot X_1 A P_2 \\
&\quad + \eta (T^\top \odot X_1 A X_2^\top) (P_2 \odot X_2 A P_1) \\
&\quad + \eta (T^\top \odot X_1 A^\top X_2^\top) (P_2 \odot X_2 A^\top P_1)
\end{aligned} \tag{19}$$

where $X_1 := I - \mathbf{1} P_1^\top$, $X_2 := I - \mathbf{1} P_2^\top$, and $T := P_2 P_1^\top$.

D Second-order policy gradients

In this Section, we show that:

$$\nabla_{\theta_1} \nabla_{\theta_2} v_2 = T \odot (I - \mathbf{1} P_2^\top) A (I - P_1 \mathbf{1}^\top)$$

where $T := P_2 P_1^\top$ is agent 2's transition matrix.

Proof.

$$\begin{aligned}
\nabla_{\theta_1} \nabla_{\theta_2} v_2 &= \text{diag}(P_2) - P_2 P_2^\top A (\text{diag}(P_1) - P_1 P_1^\top) \\
&= \text{diag}(P_2) A \text{diag}(P_1) - \text{diag}(P_2) A P_1 P_1^\top - P_2 P_2^\top A \text{diag}(P_1) + P_2 P_2^\top A P_1 P_1^\top
\end{aligned} \tag{20}$$

Note that, for $X, Y \in \mathbb{R}^n$:

$$XY^\top = X\mathbf{1}^\top \odot \mathbf{1}Y^\top$$

since:

$$\begin{pmatrix} x_1y_1 & x_1y_2 & \dots & x_1y_n \\ x_2y_1 & x_2y_2 & \dots & x_2y_n \\ \vdots & \vdots & & \vdots \\ x_ny_1 & x_ny_2 & \dots & x_ny_n \end{pmatrix} = \begin{pmatrix} x_1 & x_1 & \dots & x_1 \\ x_2 & x_2 & \dots & x_2 \\ \vdots & \vdots & & \vdots \\ x_n & x_n & \dots & x_n \end{pmatrix} \odot \begin{pmatrix} y_1 & y_2 & \dots & y_n \\ y_1 & y_2 & \dots & y_n \\ \vdots & \vdots & & \vdots \\ y_1 & y_2 & \dots & y_n \end{pmatrix}$$

Also, for $M \in \mathbb{R}^{n,n}$ note that:

$$\text{diag}(X)M = X\mathbf{1}^\top \odot M$$

since:

$$\begin{pmatrix} x_1m_{1,1} & x_1m_{1,2} & \dots & x_1m_{1,n} \\ x_2m_{2,1} & x_2m_{2,2} & \dots & x_2m_{2,n} \\ \vdots & \vdots & & \vdots \\ x_nm_{n,1} & x_nm_{n,2} & \dots & x_nm_{n,n} \end{pmatrix} = \begin{pmatrix} x_1 & x_1 & \dots & x_1 \\ x_2 & x_2 & \dots & x_2 \\ \vdots & \vdots & & \vdots \\ x_n & x_n & \dots & x_n \end{pmatrix} \odot \begin{pmatrix} m_{1,1} & m_{1,2} & \dots & m_{1,n} \\ m_{2,1} & m_{2,2} & \dots & m_{2,n} \\ \vdots & \vdots & & \vdots \\ m_{n,1} & m_{n,2} & \dots & m_{n,n} \end{pmatrix}$$

And similarly:

$$M\text{diag}(X) = (\text{diag}(X)M^\top)^\top = (X\mathbf{1}^\top \odot M^\top)^\top = M \odot \mathbf{1}X^\top$$

So, taking a closer look at each term in Equation 20:

$$\text{diag}(P_2)A \text{diag}(P_1) = T \odot A$$

$$\begin{aligned} \text{diag}(P_2)AP_1P_1^\top &= P_2\mathbf{1}^\top \odot AP_1P_1^\top \\ &= P_2\mathbf{1}^\top \odot AP_1\mathbf{1}^\top \odot \mathbf{1}P_1^\top \\ &= P_2\mathbf{1}^\top \odot \mathbf{1}P_1^\top \odot AP_1\mathbf{1}^\top \\ &= T \odot AP_1\mathbf{1}^\top \end{aligned}$$

$$\begin{aligned} P_2P_2^\top A \text{diag}(P_1) &= P_2P_2^\top A \odot \mathbf{1}P_1^\top \\ &= P_2\mathbf{1}^\top \odot \mathbf{1}P_2^\top A \odot \mathbf{1}P_1^\top \\ &= P_2\mathbf{1}^\top \odot \mathbf{1}P_1^\top \odot \mathbf{1}P_2^\top A \\ &= T \odot \mathbf{1}P_2^\top A \end{aligned}$$

$$P_2P_2^\top AP_1P_1^\top = P_2v_2P_1^\top = T \odot v_2\mathbf{1}\mathbf{1}^\top$$

This enables writing the LOLA second-order gradient as:

$$\nabla_{\theta_1} \nabla_{\theta_2} v_2 = T \odot (A - AP_1\mathbf{1}^\top - \mathbf{1}P_2^\top A + v_2\mathbf{1}\mathbf{1}^\top)$$

The term between parentheses can be factorized:

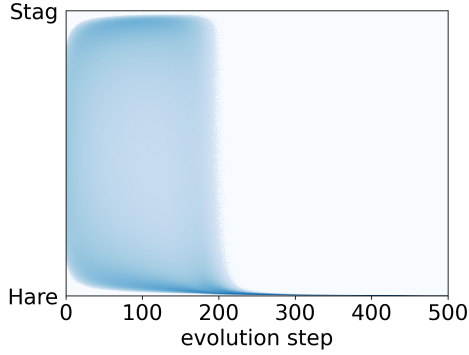
$$\begin{aligned}
 A - AP_1\mathbf{1}^\top - \mathbf{1}P_2^\top A + v_2\mathbf{1}\mathbf{1}^\top &= A - AP_1\mathbf{1}^\top - \mathbf{1}P_2^\top A + \mathbf{1}v_2\mathbf{1}^\top \\
 &= A - AP_1\mathbf{1}^\top - \mathbf{1}P_2^\top A + \mathbf{1}P_2^\top AP_1\mathbf{1}^\top \\
 &= (I - \mathbf{1}P_2^\top)A - (I - \mathbf{1}P_2^\top)AP_1\mathbf{1}^\top \\
 &= (I - \mathbf{1}P_2^\top)(A - AP_1\mathbf{1}^\top) \\
 &= (I - \mathbf{1}P_2^\top)A(I - P_1\mathbf{1}^\top)
 \end{aligned}$$

So:

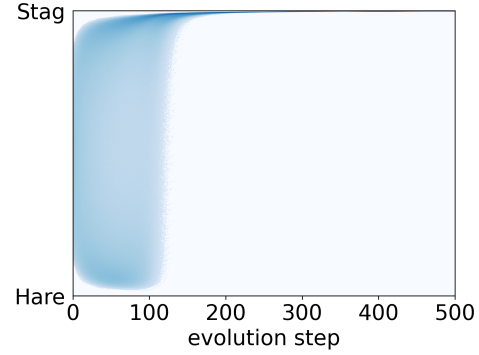
$$\nabla_{\theta_1} \nabla_{\theta_2} v_2 = T \odot (I - \mathbf{1}P_2^\top)A(I - P_1\mathbf{1}^\top)$$

□

E Stag Hunt



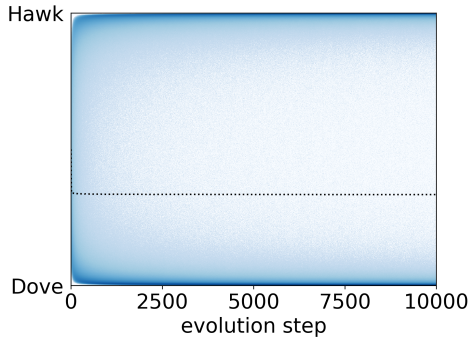
(a) 85% LOLA



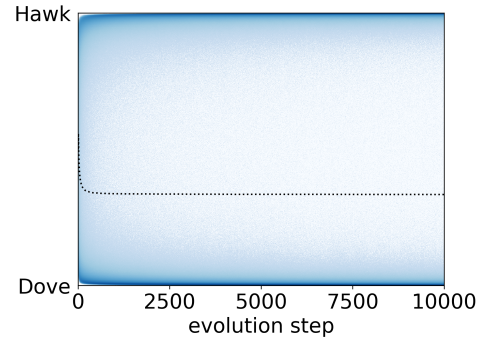
(b) 86% LOLA

Figure 6: Mixed PG and LOLA in Stag Hunt ($s = 1.8$). When more than 86% of the population is made of LOLA agents, opponent-aware learners bring the entire population to the pro-social equilibrium (NB: the higher s is, the lower this threshold becomes; it reaches 0% when $s = 2$).

F Hawk Dove



(a) PG late evolution



(b) LOLA late evolution

Figure 7: Late evolution in Hawk-Dove ($f = -2$)

G Rock-Paper-Scissors

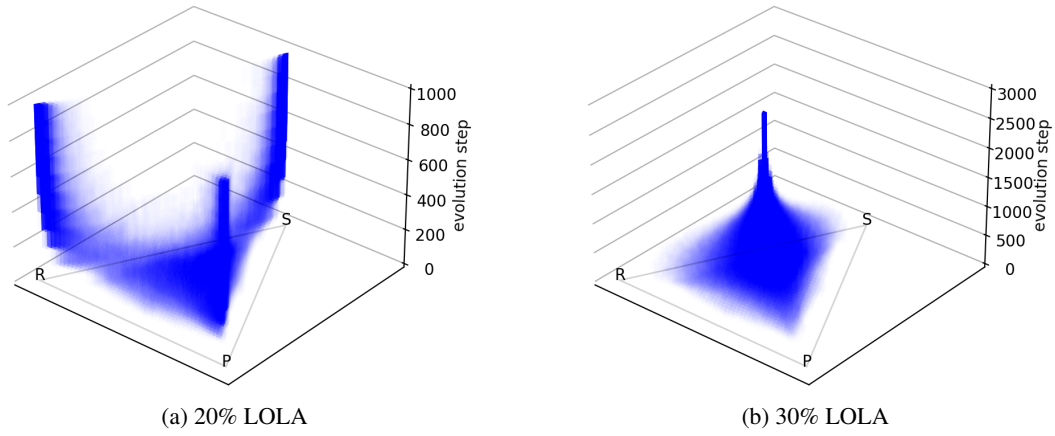


Figure 8: Mixed PG and LOLA in Rock-Paper-Scissors

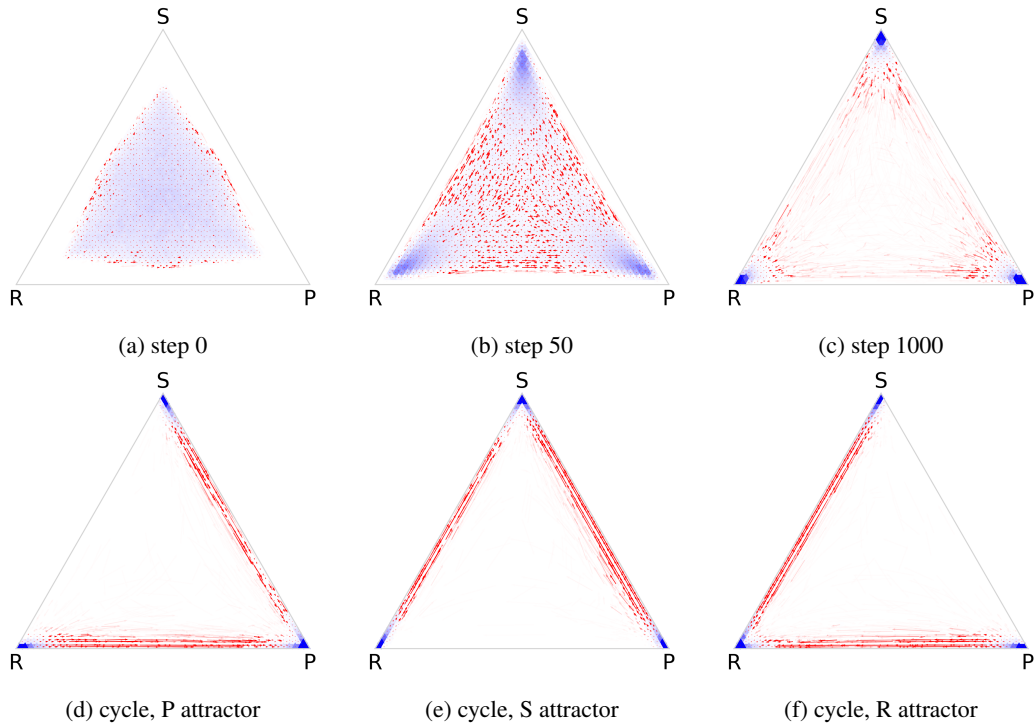


Figure 9: Naive learning in RPS, late evolution. Shades of blue indicate the concentration of individuals, while red arrows indicate their average measured movement. After about 4000 evolution steps, random drift slightly unbalances the 3 groups of near-deterministic individuals, which generates a cyclic attractor pattern in the population. In the RPS model, naive learning sustains diversity.

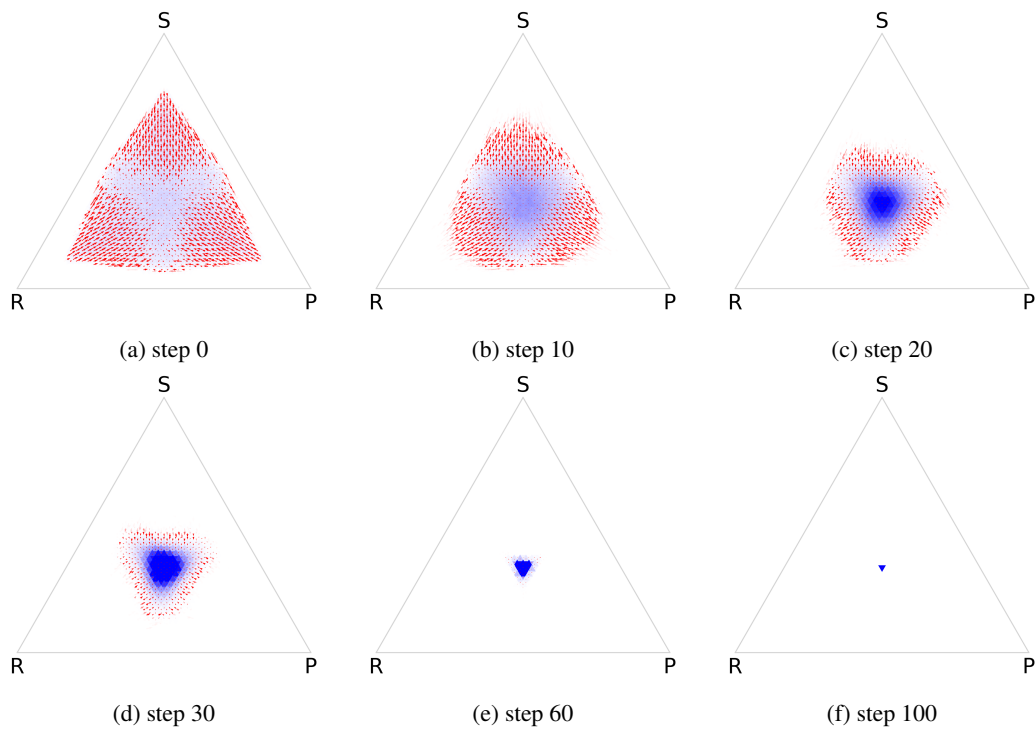


Figure 10: LOLA in RPS. Opponent-learning-awareness quickly brings the entire population to unanimously play the Nash equilibrium of this game (even when 70% of the population is naive, as shown in Figure 8b). In the RPS model, LOLA hinders diversity.