

HYBRID LEARNING IN SIGNALING GAMES

JEFFREY A. BARRETT, CALVIN T. COCHRAN, NAOKI FUJIWARA,
SIMON HUTTEGGER

ABSTRACT. Lewis-Skyrms signaling games [13, 16] have been studied under a variety of low-rationality learning dynamics [3, 5, 11, 12, 10]. Reinforcement dynamics are stable but slow and prone to evolving suboptimal signaling conventions. Trial-and-error learning is fast and reliable at finding perfect signaling conventions but unstable in the context of noise or agent error. Here we consider a low-rationality hybrid of reinforcement and trial-and-error learning that exhibits the virtues of both reinforcement and trial-and-error learning. This hybrid is reliable, stable, and exceptionally fast.

1. INTRODUCTION

Lewis-Skyrms signaling games illustrate how it is possible for agents with limited dispositional resources to evolve successful signaling systems as they interact with the world and each other.¹ The simplest sort of signaling game consists of one sender and one receiver. In an $n \times n \times n$ signaling game, there are n states of nature, n signals the sender might use, and n actions the receiver may perform on receiving a signal. Each action is appropriate to precisely one state of nature. On a play of the game, the sender observes a randomly selected state of nature then sends a signal. The receiver, who cannot see the state of nature, observes the signal then performs an action that either matches the state of nature and is successful or does not and is unsuccessful. The agents update their dispositions to signal (conditional on the state of nature) and to act (conditional on the signal) based on the success or failure of the receiver's action.

To be successful the sender and receiver must evolve a simple signaling system by updating their dispositions on repeated plays of the game. This is a subtle learning task since they must simultaneously establish interrelated conventions and learn to use these conventions for successful action in the context of many degrees of freedom. How the agents update their dispositions on repeated plays of the game is given by their learning dynamics. A learning dynamics where the agents are

Date: February 26, 2016.

¹See David Lewis's [13] characterization of signaling games. See Barrett [2] for an example of the evolution of a simple grammar in a two-sender signaling game, Skyrms [16] for an overview signaling games, and Barrett and Skyrms [4] for a discussion of how signaling games may themselves evolve.

expected to establish optimal signaling conventions in an $n \times n \times n$ signaling game is strong, especially if n is significantly greater than 2.²

In order to provide a compelling story regarding how linguistic conventions might evolve in nature, one would like to show how it is possible for agents to evolve a successful signaling system by means of a low-rationality learning dynamics that is both simple and generic. The more sophisticated or ad hoc the learning dynamics required, the less compelling the explanation. Simple, generic learning dynamics are also well-suited to applications in automated decision-making and artificial intelligence more generally. Hence, we will focus here on such dynamics.

Signaling games have been studied under a variety of low-rationality learning dynamics.³ Most of these involve either reinforcement or trial-and-error. On reinforcement learning, the sender and receiver gradually tune their conditional dispositions on the basis of their success and failure in action. Such reinforcements act as an evolving memory of what has worked well in the past. This gradually evolving memory also provides the agents with stable dispositions. If an agent makes a mistake observing the state of nature, sending or receiving a signal, or performing an action or if the signal is flipped by channel noise on a play of the game, that play may not lead to successful action. But such one-shot errors typically do little to change the dispositions the agents have forged by gradual reinforcement over time. This stability, however, comes at a significant cost. Reinforcement learning is slow, and its sluggishness frequently leaves agents stranded playing suboptimal signaling strategies, especially for large n .

In contrast, trial-and-error learning often allows agents to evolve a perfect signaling system quickly. On this sort of dynamics, an agent randomly selects a conditional strategy, then only shifts to a new strategy if their present strategy fails. Since they shift on failure, the agents do not get stuck playing suboptimal strategies.⁴ But the cost of such flexibility is instability. If the sender and receiver are playing optimal signaling strategies and an agent makes a mistake or a signal is flipped by channel noise, then they will fail on that play of the game and, consequently, shift away from the optimal coordinated conventions they have evolved. They then have to find their way back to optimal play by randomly shifting on subsequent failures.

In the present paper we consider the learning dynamics win-stay/lose-randomize with reinforcement (WS/LRwR), a low-rationality hybrid of reinforcement and trial-and-error learning that preserves the virtues of both. It is reliable, stable,

²For all its other virtues, as we will see, this is not the case for simple reinforcement learning.

³See [3, 5, 11, 12, 10] for examples.

⁴See [15] and [12] for discussions of convergence to optimal signaling for probe and adjust learning.

and exceptionally fast. In particular, it easily outperforms both simple reinforcement learning and win-stay/lose-randomize learning in the difficult task of evolving a successful signaling system.

2. REINFORCEMENT LEARNING

On simple reinforcement learning (SR) in an $n \times n \times n$ signaling game, the sender has one urn for each possible state of nature, with each urn initially containing a single ball of each possible signal type; and the receiver has one urn for each possible signal type, with each urn initially containing a single ball of each possible act type.⁵ On a play of the game, the sender sees the state of nature, draws a ball at random from the corresponding sender urn, then sends the signal indicated on that ball. The receiver sees the signal, draws a ball at random from the corresponding receiver urn, then performs the action indicated on that ball. If the action matches the state, it is successful, and each agent returns the ball she drew to the urn from which she drew it then adds a new ball to that urn of the same type; otherwise, each agent simply returns the ball she drew to the urn from which she drew it.

One of the virtues of reinforcement learning is that it is exceptionally simple. As with the other low-rationality dynamics considered here, it might be implemented in the dispositions of a simple physical system.⁶

Initially, such agents will only be successful by blind luck, but as they reinforce on success, they may evolve a signaling system that coordinates the receiver's action with the current state of nature. For the $2 \times 2 \times 2$ signaling game with unbiased nature, one can show that simple reinforcement learning will lead to a perfect signaling system with probability one.⁷ But for a $n \times n \times n$ game where n is greater than 2 or if nature is biased, the agents will often (indeed for larger n or significant bias, typically) not converge to an optimal signaling system.⁸ When they fail, they end up in a suboptimal partial pooling equilibria playing mixed conditional strategies.

A few examples of simple reinforcement learning in $n \times n \times n$ signaling games where $n > 2$ will be useful for the purpose of comparison. On simulations of the $n=3$ game, the system fails to do better than a 0.8 cumulative success rate on just 0.096 of 10^3 runs with 10^6 plays per run. For the $n=4$ game, the failure rate is higher at 0.219. And for the $n=8$ game, the failure rate is much higher at 0.594.

⁵See Herrnstein [7] for an early characterization of simple reinforcement learning. More sophisticated forms of reinforcement learning have also been studied. Some of these model human behavior well in some learning contexts [14] [6] [3] [5].

⁶See [4] for a description of how such systems then might evolve a signaling game by ritualization.

⁷See Argiento, Pemantle, Skyrms, and Volkov [1].

⁸See [3] and [8] for detailed discussions.

In those runs where the game fails to evolve perfect signaling, the agents end up playing a suboptimal combination of mixed conditional strategies.⁹

While simple reinforcement learning is relatively slow and often fails to evolve perfect signaling for larger n , whatever degree of success the composite system does attain is stable since it is these historically successful dispositions that are most strongly reinforced.

3. WIN-STAY/LOSE-RANDOMIZE

Win-stay/lose-randomize learning (WS/LR) is a form of trial-and-error learning. It can be applied to complete strategies or to individual actions.¹⁰ Here we will consider WS/LR applied to individual conditional actions.

In an $n \times n \times n$ signaling game on WS/LR the sender starts by randomly assigning a signal to each of the n states of nature, and the receiver starts by randomly assigning an act to each of the n signals.¹¹ On each play of the game, the sender observes the state of nature, then sends the signal currently assigned to that state. The receiver observes the signal, then performs the action currently assigned to that signal. If the act is successful, the agents keep their assignments of signals to states and acts to signals (win-stay); otherwise, the sender randomly selects a signal with uniform probabilities and assigns it to the current state and the receiver randomly selects an act with uniform probabilities and assigns it to the current signal (lose-randomize).¹²

The first two lines of table 1 show the mean and median number of plays for convergence to perfect signaling for an $n \times n \times n$ signaling game with WS/LR learning for each n from 2 to 8.¹³ While the differences between the two learning dynamics prevent a simple comparison, WS/LR is roughly three orders of magnitude faster

⁹See [3, 16, 9, 10] for further details regarding the behavior of simple reinforcement learning.

¹⁰See [12] for the former sort of trial-and-error learning. Win-stay/lose-randomize is closely related to win-stay/lose-shift and probe-and-adjust learning. Win-stay/lose-shift requires the agent to move to a new strategy if the current strategy fails. One virtue of win-stay/lose-randomize, is that the agent need not track what strategy led to failure when a new strategy is selected. In contrast, a virtue of probe-and-adjust is that one does not need a notion of what it is to win on the play of a strategy. Rather, one just needs to know whether one does *better* when one probes another option. While the present paper is concerned with WS/LS, as one might imagine, a hybrid dynamics of probe-and-adjust and reinforcement also has virtues over probe-and-adjust alone.

¹¹Note that the initial random assignment of signals to states may not use all of the signals and that the initial assignment of acts to signals may not use all of the acts.

¹²Note that the currently used signal or act may serve as the new signal or act since there is no restriction on the assignment when it is randomized on failure. Note also that this dynamics may assign the same signal to different states or the same act to different signals. But if a perfect signaling system evolves, and it will on this dynamics, such assignments will be temporary.

¹³The simulations of WS/LR and WS/LRwR were run in JAVA using the Eclipse integrated development environment. The game was run 10^4 times for each n .

than SR for the games considered here. Further, unlike SR, WS/LR is not susceptible to suboptimal pooling equilibria. On the simulations we discuss here WS/LR was always observed to converge to perfect signaling in finite times.¹⁴

	2x2x2	3x3x3	4x4x4	5x5x5	6x6x6	7x7x7	8x8x8
WS/LR mean	11.2	90	674	6,590	83,912	1,362,729	25,605,008
WS/LR median	7	62	463	4,583	58,001	940,789	17,539,708
WS/LRwR mean	9.77 [0.9985]	50 [0.9968]	165 [0.9949]	545 [0.9953]	3,322 [0.9964]	40,021 [0.9958]	565,594 [0.9950]
WS/LRwR median	6	32	88	188	353	598	973

TABLE 1. Speeds for WS/LR and WS/LRwR (noise-free)

It is the forgetfulness of WS/LR that prevents the agents from getting stuck playing suboptimal strategies, but it is the same forgetfulness that makes the dynamics extremely unstable. If perfect signaling has evolved under WS/LR, then a single mistake by either player or a single signal flipped by channel noise may kick the agents out of equilibrium. And, as we will see later, returning to equilibrium after being kicked out is a difficult task for WS/LR.

4. WIN-STAY/LOSE-RANDOMIZE WITH REINFORCEMENT

Ideally, one would like to have a low-rationality dynamics with the stability of SR and the speed of WS/LR. Win-stay/lose-randomize with reinforcement (WS/LRwR) is a hybrid dynamics that provides both virtues.

In an $n \times n \times n$ signaling game on WS/LRwR learning, the sender starts by randomly assigning a signal to each of the n states of nature, and the receiver starts by randomly assigning an act to each of the n possible signals. On a play of the game, the sender observes the state or nature, then sends the signal currently assigned to that state. The receiver observes the signal, then performs the action currently assigned to that signal. If the act is successful, the agents keep their assignments of signals to states and acts to signals and tally the success for both the current conditional signal and current conditional act (win-stay and reinforce).¹⁵

¹⁴Indeed, one can prove that WS/LR will converge to an optimal signaling system with probability one. See [12] and [15] for proofs of convergence in finite times for closely associated learning dynamics.

¹⁵This is the mechanism borrowed from SR. Rather than determine the probability of a particular action, however, reinforcement on this dynamics determines the probability of each conditional action being selected as a new strategy if the currently strategy fails.

Otherwise, the sender randomly selects a new signal with probabilities equal to the past success rate for each of the possible signals conditional on the current state, then assigns that signal to the current state; and the receiver randomly selects an act with probabilities equal to the past success rate for each of the possible acts conditional on the current signal, then assigns that act to the current signal (lose-randomize using probabilities determined by past reinforcements).¹⁶

The third and fourth lines of table 1 show the mean and median number of plays to convergence for an $n \times n \times n$ signaling game with WS/LRwR learning for each n . While the thought was to capture the stability of SR with the hybrid dynamics, an issue we will consider in the next section, WS/LRwR is also much faster than WS/LR, especially for larger n . The mean time to convergence is two orders of magnitude faster, and the median is five orders of magnitude faster for $n = 8$. The reinforcements allow WS/LRwR to hold steady those parts of the signaling system that are working while finding conventions for the parts where successful conventions have not yet been formed.

Before considering the relative stability of the two dynamics, it is important to note that there is one sense in which WS/LR is better behaved than WS/LRwR. While every run of WS/LR was successful, WS/LRwR occasionally failed to find a signaling system even when run for twice as long as the longest time it took for WS/LR to find a signaling system.¹⁷ The numbers in square brackets on the second row of table 1 indicate the proportion of runs where the agents converged to perfect signaling on WS/LRwR for each n . The upshot is that WS/LRwR gains a burst in speed on most runs, but at the cost of occasionally taking longer than WS/LR in finding an optimal signaling system.¹⁸

5. NOISE

While the increase in speed of WS/LRwR over WS/LR is welcome, the motivation behind the hybrid dynamics was the instability of WS/LR when there is possibility of agent error or noise. To compare the two dynamics in this regard, we considered their stability in the context of random channel noise where, with probability 0.01, the receiver observes a random signal rather than the signal sent by the sender.¹⁹

¹⁶As with WS/LR, the initial random assignment may not use all of the signals or acts, and the current signal or act may serve as the new randomly selected signal or act on WS/LRwR.

¹⁷We used a cutoff of twice the longest time it took for WS/LR to find an optimal signaling system for each n when we ran WS/LRwR. The means and medians are calculated on the runs that converged before hitting the cut-off.

¹⁸Note that whether WS/LRwR is sure to eventually find an optimal signaling system is currently an open question.

¹⁹The probability of each random signal is $1/n$ for each game. Note that it is possible that the randomly selected signal matches the signal sent.

The mean and median number of plays to reach optimal signaling for both dynamics is given in table 2. Again, the numbers in square brackets indicate the proportion of runs where the agents found optimal signaling on WS/LRwR for each n . Tables 1 and 2 are comparable.

	2x2x2	3x3x3	4x4x4	5x5x5	6x6x6	7x7x7	8x8x8
WS/LR mean	11.4	92	693	6,957	92,014	1,512,384	29,400,236
WS/LR median	8	65	482	4,852	64,006	1,037,499	19,941,227
WS/LRwR mean	10.0 [0.9985]	50 [0.9956]	175 [0.9938]	511 [0.9934]	2,101 [0.9895]	32,328 [0.9874]	409,742 [0.9823]
WS/LRwR median	7	33	91	192	369	620	971

TABLE 2. Speeds for WS/LR and WS/LRwR (0.01 noise)

In the noise-free case, both WS/LR and WS/LRwR are stable since the agents are always successful after finding an optimal signaling equilibrium. But channel noise makes optimal play unstable for both learning dynamics. Once they find an optimal signaling system, if any signal is randomly flipped by channel noise, the agents will fail on that play. Since WS/LR has no memory, the agents will immediately randomize their conditional strategies with uniform probabilities for each. Agents on WS/LRwR, however, typically continue to play their current conditional strategies insofar as these strategies have been historically successful and hence have been significantly reinforced.

	2x2x2	3x3x3	4x4x4	5x5x5	6x6x6	7x7x7	8x8x8
WS/LR mean	265	169	144	132	122	120	116
WS/LR median	183	117	101	90	85	83	80
WS/LRwR mean	90 [0.008]	243 [0.2967]	450 [0.5491]	640 [0.7526]	799 [0.8764]	879 [0.9391]	936 [0.9747]
WS/LRwR median	37	114	199	255	301	307	335

TABLE 3. Mean plays to break optimal signaling (0.01 noise)

Table 3 indicates the mean and median number of plays that it takes for the 0.01 channel noise to break the optimal signaling system for those runs which are observed to leave the first evolved signaling system.²⁰ Note that WS/LRwR is an order of magnitude more stable than WS/LR for larger n . Indeed, if a single signal is randomly flipped, WS/LR nearly always breaks and WS/LRwR nearly always preserves optimal play. Further, when WS/LR breaks, with no memory of what has worked well in the past, it typically unravels the entire system of evolved conventions. The first two lines of table 4 give the mean and median number of plays per channel error it takes for WS/LR to find an optimal signaling system when it is kicked out.²¹ In short, it typically takes as long as it took to find optimal signaling in the first place. Finally, when it finds an optimal signaling system again, the new system is typically different from the one that initially evolved.²² WS/LR, then, is extremely unstable in the context of error or noise.

The third and fourth lines of table 4 give the mean and median number of plays for WS/LRwR to return to an optimal signaling system per channel error. Measured by the mean, WS/LRwR is three orders of magnitude more stable than WS/LR for $n = 8$. Measured by the median, it is seven orders of magnitude more stable. Further, as the median indicates, WS/LRwR typically either does not leave optimal signaling at all or just bounces back to its initial optimal signaling system when it gets kicked out by error. Indeed, on simulation, the system is usually observed to find the same equilibrium it had initially evolved. Optimal signaling conventions do sometimes unravel under WS/LRwR. But, while this is typical for WS/LR, it is rare for WS/LRwR. The reinforcements under WS/LRwR provide a memory that typically guides the system back to optimal play by biasing the random shifts on failure to favor conditional actions that have worked well in the past.

Significantly, the longer the system spends at optimal play, the stronger the reinforcements that guide it back when it leaves. One can, consequently, expect agent error and system noise to become less effective at ever kicking the system out of optimality the longer the game is played. One can also expect a yet surer path back to optimality if error does kick the system out of optimal play. And finally, since the system typically does not stray too far from the initial optimal signaling system it evolved, one can expect a high degree of success while it is returning. In other words, given the memory provided by past reinforcements on success, one can expect WS/LRwR typically to produce ever more stable successful dispositions.

²⁰Many runs under WS/LRwR hit the run-length cap never leaving the signaling system they initially evolve. For $n = 2$, for example, 0.92 of the runs are never observed to leave the initially evolved signaling system. For $n = 4$, the proportion of runs that never leave is somewhat lower at 0.45.

²¹All of the data in this table is per channel error. The medians are medians of the mean number of plays to return to convergence per channel error.

²²For the $n = 6$ game, for example, the agents almost always evolve a new signaling system.

	2x2x2	3x3x3	4x4x4	5x5x5	6x6x6	7x7x7	8x8x8
WS/LR mean	7.46	64	556	5,836	82,227	1,370,470	27,000,015
WS/LR median	4	38	341	3,837	53,915	901,404	24,584,681
WS/LRwR mean	0.10	1.35	6.84	20.50	169	6,702	73,516
WS/LRwR median	0.00	0.00	0.14	0.60	1.00	1.50	2.00

TABLE 4. Plays to return to optimal signaling (0.01 noise)

6. DISCUSSION

Since the agents must establish and learn conventions in the context of many degrees of freedom, signaling games pose a difficult learning problem. Simple reinforcement learning (SR) is successful in finding optimal conventions only in the context of the simplest games. The agents almost always find an optimal equilibrium for the $n=2$ game with unbiased nature, but they are highly unlikely to do so if n is significantly greater than 2 or if nature exhibits a strong bias in states. And even when the agents are successful in evolving a perfect signaling system under SR, they do so very slowly. On the other hand, the dispositions the agents evolve by reinforcement are stable since occasional agent error or system noise do little to disturb their past reinforcements.

Win-stay/lose-randomize (WS/LR) learning is much faster than SR, it always eventually finds an optimal signaling strategy, and it sticks with the optimal strategy it finds in the noise-free case. But WS/LR is wildly unstable. A single mistake or noisy signal typically knocks the system out of a coordinated system of optimal signaling conventions. When it does, it takes as long to find new signaling conventions as it took to find the first. And, when they are found, they are typically different.

The problem is that WS/LR is perfectly forgetful. All that matters is what happened in the last play of the game. This prevents past reinforcements on marginally successful strategies from standing in the way of the agents finding a perfect signaling system, but it also means that the agents have no record of what worked well in the past to guide them if an error occurs on a particular play. Hence, their dispositions under the dynamics are maximally unstable.

Win-stay/lose-randomize with reinforcement (WS/LRwR) combines the virtues of SR and WS/LR. While suboptimal partial pooling equilibria are possible, they

are rare for both noise-free and noisy systems. It is typically many orders of magnitude faster than WS/LR in finding a perfect signaling system, and the stability it inherits from SR helps it keep the system it evolves. The trial-and-error feature of the dynamics allows for the free exploration of alternative strategies and the reinforcements act as a memory of what has worked well in the past. The memory typically keeps the system close to optimal play and provides increasing stability over time for successful dispositions.

WS/LRwR is just trial-and-error learning where the likelihood of a conditional action being tried on failure is determined by how well that action has worked in the past. As such, it is a generic, low-rationality dynamics that is easily implemented in an artificial system and that one might also expect to find exhibited in the behavior of natural agents.

REFERENCES

- [1] Argiento, Raffaele, Robin Pemantle, Brian Skyrms and Stas Volkov (2009) “Learning to Signal: Analysis of a Micro-Level Reinforcement Model,” *Stochastic Processes and Their Applications* 119(2): 373–390.
- [2] Barrett, J. A. (2007a) “Dynamic Partitioning and the Conventionality of Kinds,” *Philosophy of Science* 74: 527–546.
- [3] Barrett, J. A. (2006) “Numerical Simulations of the Lewis Signaling Game: Learning Strategies, Pooling Equilibria, and the Evolution of Grammar,” *Institute for Mathematical Behavioral Sciences Paper* 54. <http://repositories.cdlib.org/imbs/54>.
- [4] Barrett, J. A. and B. Skyrms (2015) “Self-Assembling Games,” *The British Journal for the Philosophy of Science*. Published online 13 September 2015. doi: 10.1093/bjps/axv043
- [5] Barrett, J. A. and Kevin Zollman (2009) “The Role of Forgetting in the Evolution and Learning of Language,” *Journal of Experimental and Theoretical Artificial Intelligence* 21(4): 293–309.
- [6] Erev, I. and A. E. Roth (1998) “Predicting How People Play Games: Reinforcement Learning in Experimental Games with Unique, Mixed Strategy Equilibria” *American Economic Review* 88: 848–81.
- [7] Herrnstein, R. J. (1970) “On the Law of Effect,” *Journal of the Experimental Analysis of Behavior* 13: 243–266.
- [8] Huttegger, Simon (2007) “Evolutionary Explanations of Indicatives and Imperatives. *Erkenntnis* 66, 2007, 409–436.
- [9] Hofbauer, Josef and Simon Huttegger (2008) “Feasibility of Communication in Binary Signaling Games,” *Journal of Theoretical Biology* 254(4): 843–849.
- [10] Simon M. Huttegger, Brian Skyrms, Rory Smead, Kevin J. S. Zollman (2010) “Evolutionary Dynamics of Lewis Signaling Games: Signaling Systems vs. Partial Pooling” *Synthese* 172(1): 177–191.
- [11] Huttegger, Simon, Brian Skyrms, Pierre Tarrès, and Elliott Wagner (2014) “Some Dynamics of Signaling Games,” *Proceedings of the National Academy of Sciences* 111(S3): 10873–10880.
- [12] Simon M. Huttegger, Brian Skyrms, and Kevin J. S. Zollman (2014) “Probe and Adjust in Information Transfer Games” *Erkenntnis* 79 (S4):1–19 (2013)
- [13] Lewis, David (1969) *Convention*. Cambridge, MA: Harvard University Press.
- [14] Roth, A. E. and I. Erev (1995) “Learning in Extensive Form Games: Experimental Data and Simple Dynamical Models in the Immediate Term,” *Games and Economic Behavior* 8:164–212.
- [15] Skyrms, Brian (2014) “Learning to Signal with Two Kinds of Trial and Error” forthcoming in *Foundations and Methods from Mathematics to Neuroscience: Essays Inspired by Patrick Suppes* Colleen E. Crangle, Adolfo Garcia de la Sienra, and Helen E. Longino (eds): CSLI Publications.
- [16] Skyrms, Brian (2010) *Signals Evolution, Learning, & Information*. New York: Oxford University Press.