

Complex dynamics in learning complicated games

Tobias Galla^{a,1} and J. Doyne Farmer^{b,c,1}

^aTheoretical Physics, School of Physics and Astronomy, University of Manchester, Manchester M13 9PL, United Kingdom; ^bInstitute for New Economic Thinking at the Oxford Martin School and Mathematical Institute, University of Oxford, Oxford OX1 3LP, United Kingdom; and ^cSanta Fe Institute, Santa Fe, NM 87501

Edited by Kenneth Wachtler, University of California, Berkeley, CA, and approved November 27, 2012 (received for review June 29, 2011)

Game theory is the standard tool used to model strategic interactions in evolutionary biology and social science. Traditionally, game theory studies the equilibria of simple games. However, is this useful if the game is complicated, and if not, what is? We define a complicated game as one with many possible moves, and therefore many possible payoffs conditional on those moves. We investigate two-person games in which the players learn based on a type of reinforcement learning called experience-weighted attraction (EWA). By generating games at random, we characterize the learning dynamics under EWA and show that there are three clearly separated regimes: (i) convergence to a unique fixed point, (ii) a huge multiplicity of stable fixed points, and (iii) chaotic behavior. In case (iii), the dimension of the chaotic attractors can be very high, implying that the learning dynamics are effectively random. In the chaotic regime, the total payoffs fluctuate intermittently, showing bursts of rapid change punctuated by periods of quiescence, with heavy tails similar to what is observed in fluid turbulence and financial markets. Our results suggest that, at least for some learning algorithms, there is a large parameter regime for which complicated strategic interactions generate inherently unpredictable behavior that is best described in the language of dynamical systems theory.

high-dimensional chaos | statistical mechanics

The majority of results in game theory concern simple games with a few players and a few possible moves, characterizing them in terms of their equilibria (1, 2). The applicability of this approach is not clear when the game becomes more complicated, for example due to more players or a larger strategy space, which can cause an explosion in the number of possible equilibria (3–6). This is further complicated if the players are not rational and must learn their strategies (7–11). In a few special cases, it has been observed that the strategies display complex dynamics and fail to converge to equilibrium solutions (12–14). Are such games special, or is this typical behavior? More generally, under what circumstances should we expect that games have a multiplicity of solutions, or that they become so hard to learn that their dynamics fail to converge? What kind of behavior should we expect and how should we characterize the solutions?

We do not answer these questions in full generality here, but we are able shed some light on them by investigating randomly constructed games using a specific family of learning algorithms. This is inspired by the work of Opper and Diederich (5, 6), who investigated random games with replicator dynamics and by Berg, Weigt, and McLennan (3, 4), who showed that as one deviates from the zero sum case the number of Nash equilibria grows exponentially.

As an example of what we mean by a complicated vs. a simple game, compare tic-tac-toe and chess. Tic-tac-toe is a simple game with only 765 possible positions and 26,830 distinct sequences of moves. Young children easily discover the Nash equilibrium, which results in a draw, and once their friends discover this too the game becomes uninteresting. In contrast, chess is a complicated game with roughly 10^{47} possible positions and 10^{123} possible sequences of moves; despite a huge effort, the Nash equilibrium (corresponding to an ideal game) remains unknown. Equilibrium concepts of game theory are not useful in describing complicated games such as chess or go (which has an even larger game tree with roughly 10^{360} possible sequences of moves). Another

example is investing in financial markets, which is a nonzero sum game where players can choose between thousands of assets and a rich set of possible trading strategies.

Here, we investigate a type of reinforcement learning that is extensively used both in practical applications in machine learning and to explain social experiments. We study complicated games that are constrained by the average correlation between the payoffs of the two players, but are otherwise random. Depending on the payoff correlation and the learning memory parameter, we find the asymptotic behavior of the strategy dynamics has clearly separated regimes. In regime (i), the strategies converge to unique fixed points, in regime (ii) they may converge but the number of possible fixed points is huge, and in regime (iii), no matter how long the players learn, the strategies never converge to a fixed strategy. Instead they continually vary as each player responds to past conditions and attempts to do better than the other player. The trajectories in some parts of regime (iii) display high-dimensional chaos, suggesting that for most intents and purposes the behavior is essentially random.

I. Games and Learning Algorithm

To address the questions raised above, we study two-player games. For convenience, call the two players Alice and Bob. At each time step t player $\mu \in \{\text{Alice} = A, \text{Bob} = B\}$ chooses between one of N possible moves, picking the i th move with frequency $x_i^\mu(t)$, where $i = 1, \dots, N$. The frequency vector $\mathbf{x}^\mu(t) = (x_1^\mu, \dots, x_N^\mu)$ is the strategy of player μ . If Alice plays i and Bob plays j , Alice receives payoff Π_{ij}^A and Bob receives payoff Π_{ij}^B .

We assume that the players learn their strategies \mathbf{x}^μ via a form of reinforcement learning called “experience-weighted attraction.” This has been extensively studied by experimental economists who have shown that it provides a reasonable approximation for how real people learn in games (7–9). Actions that have proved to be successful in the past are played more frequently and moves that have been less successful are played less frequently. To be more specific, the probability of a given move is as follows:

$$x_i^\mu(t) = \frac{e^{\beta Q_i^\mu(t)}}{\sum_k e^{\beta Q_k^\mu(t)}}, \quad [1]$$

where Q_i^μ is called the “attraction” for player i to strategy μ . In the special case of experience-weighted attraction that we mainly focus on here, Alice’s attractions are updated according to the following:

$$Q_i^A(t+1) = (1 - \alpha)Q_i^A(t) + \sum_j \Pi_{ij}^A x_j^B(t), \quad [2]$$

and similarly for Bob with A and B interchanged.

Author contributions: T.G. and J.D.F. designed research; T.G. performed research; T.G. and J.D.F. analyzed data; and T.G. and J.D.F. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence may be addressed. E-mail: tobias.galla@manchester.ac.uk or jdf@santafe.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1109672110/-DCSupplemental.

Under this update rule, Alice knows her own payoffs and also knows the frequency x_i^B with which Bob makes each of his possible moves (and similarly for Bob). This approximates the situation in which the players vary their strategies slowly in comparison with the timescale on which they play the game, so that Alice can collect good statistics about Bob before updating her strategy. In the machine learning literature, the practice of infrequent parameter updating is called “batch learning.” We prefer to focus on batch learning because in the large batch limit the dynamics for updating the strategies x^t of the two players are deterministic, which allows us to compute many properties of the games analytically. The alternative is “online learning,” in which each player updates strategies after every move. Because the moves are randomly chosen, this implies the learning dynamics have a stochastic component. In *SI Appendix, section V*, we present numerical simulations that show that, with a few exceptions, batch learning and online learning behave similarly, and the main results of the paper hold for both cases.

The key parameters that characterize the learning strategy are α and β . The parameter β is called the intensity of choice; when β is large, a small historical advantage for a given move causes that move to be very probable, and when $\beta = 0$, all moves are equally likely. The parameter α specifies the memory in the learning; when $\alpha = 1$, there is no memory of previous learning steps, and when $\alpha = 0$, all learning steps are remembered and are given equal weight, regardless of how far in the past. The case $\alpha = 0$ corresponds to the much-studied replicator dynamics used to describe evolutionary processes in population biology (15–17), where x_i^μ is the concentration of species i in population μ (*SI Appendix, section II*).

II. Why Investigate Random Games?

Our goal here is to characterize the typical behavior one expects a priori when the strategy for playing a game is learned under reinforcement learning. We first describe what we do, and then explain why we think it is useful.

We construct members of an ensemble of games at random subject to constraints. The constraints we use are that the payoffs have zero mean and a given positive variance, and the payoff to Alice has a given correlation to the payoff to Bob. In mathematical terms, this means that $E[\Pi_{ij}^A] = 0$, $E[(\Pi_{ij}^A)^2] = 1$, and $E[\Pi_{ij}^A \Pi_{ij}^B] = \Gamma$, where $E[x]$ denotes the average of x . Conforming to previous work involving random dynamical systems (3–6, 18), we draw the payoff matrices Π_{ij}^t from a multivariate normal distribution. This was previously presented as an arbitrary choice. In fact, the principle of maximum entropy, which is the foundation of statistical mechanics (19, 20) and has many practical applications in signal processing and machine learning (21), dictates that this is the natural choice, as it is the one that maximizes entropy subject to the above constraints.

The variable Γ is a crucial parameter that measures the correlation in the payoffs of the two players. When $\Gamma = -1$ the game is zero sum, i.e., the amount Alice wins is equal to the amount Bob loses, whereas when $\Gamma = 0$ their payoffs are uncorrelated, and when $\Gamma = 1$ their payoffs are identical. Thus, Γ can be regarded as a competition parameter, where smaller values of Γ indicate more competition.

What can we learn by studying randomly generated games? This depends on whether or not the ensemble of randomly generated games that we have constructed has characteristics that are representative of the “real” games that naturally occur in biology and social science. Can real games be regarded as typical members of the ensemble of random games satisfying the above constraints? If so, then we can obviously learn a great deal by studying the properties of the ensemble. For example, in *SI Appendix*, see our discussion of 2×2 games under replicator dynamics (*SI Appendix, section II*).

However, what if it turns out that randomly generated games as we generate them here are not representative of real games? In this case, our approach is still valuable as a null hypothesis that can be used to sharpen understanding of what makes real games special. In this case, we would seek alternative constraints leading to alternative ensembles more representative of real games.

For example, if the variance of the payoffs of real games were typically unbounded (which we doubt), this would suggest that the payoffs should be drawn instead from a heavy-tailed Levy distribution, which is the maximum entropy solution in this case. Or there might be additional constraints characterizing biological and social systems that have so far not been identified, which would modify the payoff distribution. Understanding such constraints would obviously be very illuminating, in that it would require discovering new general principles about the generic nature of strategic interactions in real contexts.

An example of work in a similar spirit is the 1972 paper of Robert May (18), which analyzed the generic stability properties of differential equations modeling predator–prey interactions with random coupling coefficients. This work challenged the conventional wisdom that more complex ecosystems are necessarily more stable than simple ones by showing that for a particular ensemble of random equations the opposite is true. The question of whether complex ecosystems are more or less stable remains controversial. In any case, May’s paper has played a vital role by focusing the debate and forcing ecologists to think carefully about the generic properties of ecological interactions. Our intent is similar. Even if it turns out that we are wrong, explaining why we are wrong will hopefully stimulate game theorists to think more carefully about the generic properties of real games.

The characterization of turbulence in fluid mechanics provides a success story analogous to the one we are seeking. The Reynolds number is the ratio of inertial forces to viscous forces in a fluid flow. Based on a very simple analysis, the Reynolds number provides an a priori estimate of whether a fluid flow is likely to be turbulent, and if so, how turbulent it will be. This estimate is crude and inaccurate, but nonetheless very useful. Similarly, our long-term goal is to predict the qualitative nature of the long-term dynamics of a complicated game based on a very cursory analysis of its parameters; we have begun that project by studying a particular ensemble of random games under reinforcement learning.

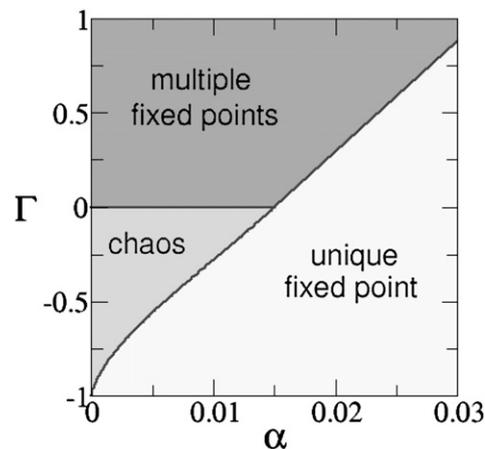


Fig. 1. Schematic illustration of the qualitative nature of the asymptotic learning dynamics in the parameter space for complicated random games. $\beta = 0.07$.

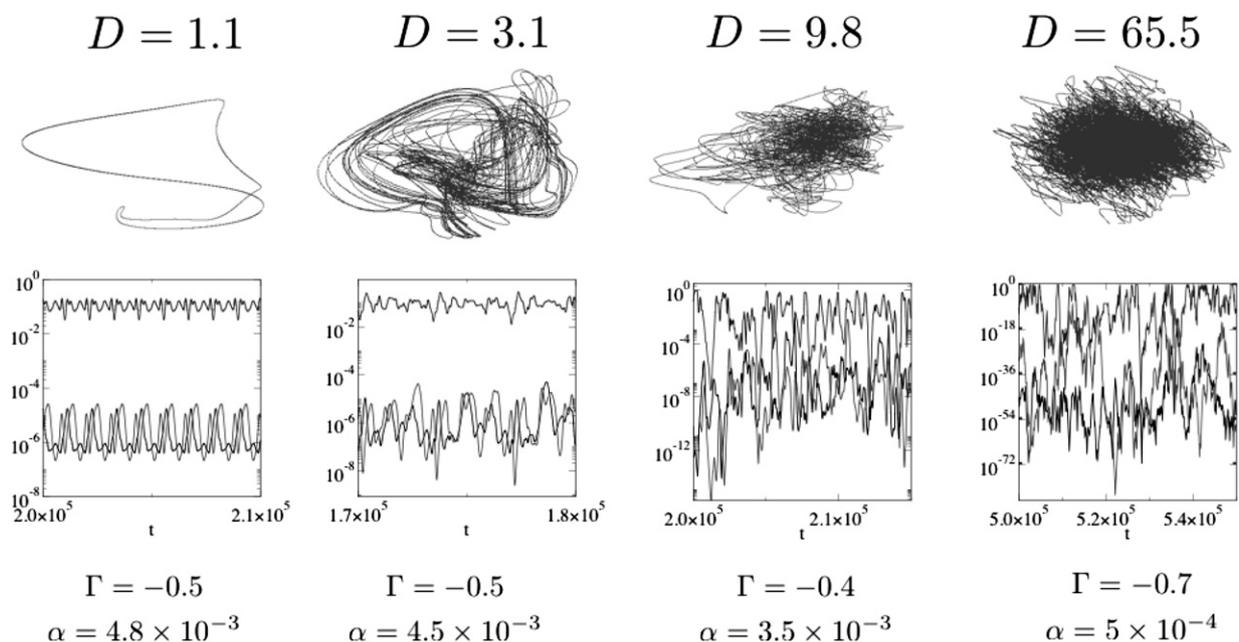


Fig. 2. An illustration of complex learning dynamics showing trajectories of the strategy $x^*(t)$ for four different sets of parameters (varying from *Left to Right*). For each parameter set, we show a phase plot on top and a time series below. For the phase plots, we project the 98-dimensional dynamics onto two dimensions, and we show three representative time series, corresponding to a particular choice of the player and actions. For clarity, we use logarithmic scale. The estimates for the attractor dimensions are obtained as explained in *SI Appendix, section IV*. As the dimension of the attractor increases, so does the range of x^*_i . For the highest dimensional case, a given move has occasional bursts where it is highly probable, and long periods where it is extremely improbable (as low as 10^{-72}). $\beta = 0.07$ in all panels.

III. Phase Diagram

We simulate randomly constructed games with $n = 50$ possible moves, corresponding to a 98-dimensional state space (there are two 50-dimensional strategy vectors and two probability constraints). Roughly speaking, we observe three regimes, as illustrated in the phase diagram of Fig. 1.

- i*) “Unique fixed point”: For large values of α and small values of Γ , roughly the lower right triangle in the phase diagram, the learning dynamics $x^*(t)$ converge to a unique stable fixed point.
- ii*) “Multiple fixed points”: In the remainder of the phase diagram with $\Gamma > 0$, we observe multiple fixed points. In this case, the multiplicity of the fixed points is often high, e.g., more than 100.
- iii*) “Chaos or limit cycle”: For Γ negative and α small, roughly the lower left triangle in the phase diagram, the learning dynamics tend to converge to limit cycles or chaotic attractors, although occasionally we observe other behaviors discussed below.

IV. Characterizing Multiple Attractors and Chaos

For regime (*i*), in which we observe unique fixed points, there is not much to say about the asymptotic behavior of the learning trajectories $x^*(t)$: A fixed point is a fixed point. For the other two regimes, however, it is interesting to investigate how the number of stable fixed points varies with parameters in regime (*ii*) and how the dimensionality of the attractors varies with parameters in regime (*iii*).

Regime (*ii*), Multiple Fixed Points. In regime (*ii*), we observe multiple stable fixed points depending on initial conditions. A crude survey of the behavior in this region is given in *SI Appendix, section IVC*. The number of fixed points can be quite high. The most complicated behavior is observed for small values of α and values of Γ near 1. In this case, we typically observe more than 100 fixed points. We want to emphasize that the problem of counting the number of fixed points is difficult, and we cannot be

sure that our answers are exact. However, we believe that our results are good enough to accurately indicate the variation between different parts of the parameter space.

Regime (*iii*), Limit Cycles and Chaos. We give several examples of the observed learning dynamics at different parameter values in regime (*iii*) in Fig. 2. These include a limit cycle and chaotic attractors of varying dimensionality. There can also be long transients in which the trajectory follows a complicated orbit for a long time and then suddenly collapses into a fixed point. In general, the behavior observed depends on the random draws of the payoff matrices Π_{ij}^u —the outcomes vary from realization to realization, but if the system is self-averaging then in the limit of large N one expects the behavior to become more and more consistent.

To characterize the local stability properties of the attractors, we numerically compute the Lyapunov exponents λ_i , $i = 1, \dots, 2N - 2$, which quantify the rate of expansion or contraction of nearby points in the state space. The Lyapunov exponents also determine the Lyapunov dimension D , which measures the number of degrees of freedom of the motion on the attractor.

Simulating games at many different parameter values gives the stability diagram shown in Fig. 3. Roughly speaking, we find that the dynamics are stable when Γ is strongly negative (nearly zero-sum games) and α is large (short memory), i.e., in the lower right of the diagram.* They are unstable when Γ is not too close to zero-sum behavior (imperfectly correlated payoffs) and α is small (long memory). Interestingly, for reasons that we do not understand, the highest dimensional behavior is observed when the payoffs are moderately anticorrelated ($\Gamma \sim -0.6$) and when players have long memory ($\alpha \sim 0$). In this regime, we also encounter numerical problems due to the extreme nonlinearities

*Note that the fixed point reached in the stable regime is only a Nash equilibrium at $\Gamma = 0$ and in the limit $\alpha \rightarrow 0$. When $\alpha > 0$, the players are effectively assuming their opponent’s behavior is nonstationary and that more recent moves are more useful than moves in the distant past.

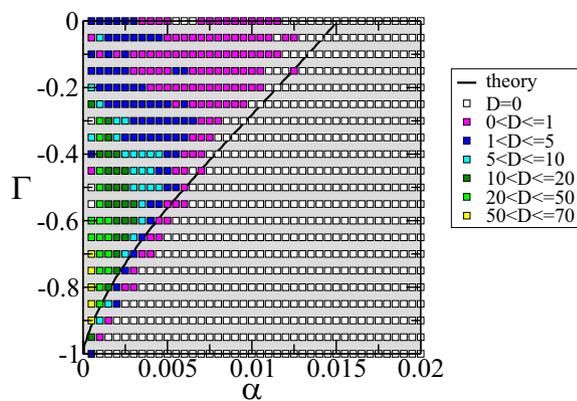


Fig. 3. Stability diagram showing where stable vs. chaotic learning is likely when $\Gamma < 0$. The colored squares represent the typical dimension of the attractor observed in our simulations (averaged over 10 or more independent payoff matrices at each grid point). White indicates fixed points, pink indicates limit cycles, and yellow indicates high-dimensional chaotic attractors. The solid line is the stability boundary estimated using path-integral methods. $\beta = 0.07$.

associated with fixed points near the edges of the simplex, and in some cases the dimension is hard to compute reliably.

A good approximation of the boundary between the stable and unstable regions of the parameter space can be computed analytically using techniques from statistical physics. We use path-integral methods from the theory of disordered systems (5, 22) to compute the stability in the limit of infinite payoff matrices, $N \rightarrow \infty$. We do this in a continuous-time limit where, for fixed Γ , stability then depends only on the ratio α/β (*SI Appendix, section III*). The results of doing this are illustrated by the solid black line in Fig. 3, which gives a good approximation for the stability boundary.

We have not yet extensively studied the behavior as N is varied. If $D > 0$ at small N , then we expect the dimension D to increase with N . At this stage, we have been unable to tell whether D reaches a finite limit or grows without bound as $N \rightarrow \infty$.

V. Clustered Volatility and Heavy Tails

An interesting property of this system is the time dependence of the received payoffs. As shown in Fig. 4, when the dynamics are chaotic the total payoff to all of the players varies, with intermittent bursts of large fluctuations punctuated by relative quiescence. This is observed, although to varying degrees, throughout the chaotic part of the parameter space. There is a strong

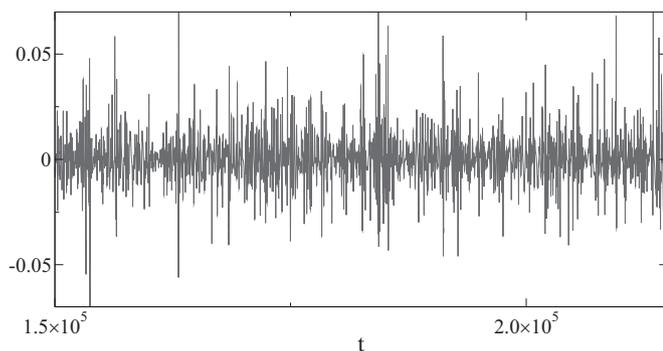


Fig. 4. Chaotic dynamics display clustered volatility. We plot the difference of payoffs on successive time steps for parameter values corresponding to the third panel from the left in Fig. 2. The amplitude of the fluctuations increases with the dimension of the attractor.

resemblance to the clustered volatility observed in financial markets, which in turn resembles the intermittency of fluid turbulence (12, 23). Financial markets are in fact multiplayer games; our preliminary studies of multiplayer games suggest that these effects become even stronger, and that most parameters lead to chaotic learning dynamics.

Similarly, we also typically observe heavy tails in the distribution of the fluctuations, as shown in *SI Appendix, section IVD*. The tails that we observe here are exponential, and not power law; we believe this is due to the fact that α sets a fixed timescale. We conjecture that if this timescale is allowed to vary, or in situations with more players with a broad distribution of learning timescales, power law tails may be observed.

The fact that this behavior occurs generically and more or less automatically suggests that these properties, which have received a great deal of attention in studies of financial markets, may occur simply because they are generic properties of the learning dynamics of complicated games.[†] Understanding and more thoroughly characterizing this will have to wait for a broader study including a variety of different learning algorithms.

VI. Why Is Dimensionality Relevant?

The fact that the learning dynamics do not converge to an equilibrium under a particular learning algorithm, such as reinforcement learning, does not in general imply that convergent learning might not happen with another algorithm. High dimensionality is relevant because it indicates that the learning dynamics are effectively random. As we explain below, this means that there is no obvious alteration in the learning algorithm of a given player that will improve that player's performance.

In contrast, if the learning dynamics settle into a limit cycle or a low dimensional attractor, Alice could collect data and make better predictions about Bob's likely decisions using the method of analogs (24) or refinements based on local approximation (25). Assuming Bob does not alter his strategy, and assuming the dynamics are sufficiently stable under variations in parameters, Alice could use her predictions to make small alterations in her strategy so as not to alter the combined learning dynamics of the two players too much, thereby perturbing the system onto a nearby attractor and improving her average payoff.

If the dimension of the chaotic attractor is too high, however, the curse of dimensionality makes this impossible with any reasonable amount of data (25). Thus, high dimensionality indicates that the chaotic attractors we observe here are rather complicated Nash equilibria, in the sense that there are no easily learnable strategies nearby.

This raises the question of whether there are games where learning will fail to cause the strategies to converge to a fixed point under any inductive learning algorithm. The observation of high-dimensional chaos adds weight to the suggestion that there are. The existence of such high-dimensional chaos is reminiscent of the ergodic conjecture of statistical mechanics, which loosely speaking says that for many nonlinear dynamical systems almost all of the trajectories display high-dimensional chaos, and consequently can only be characterized by statistical averages. See also refs. 14, 26, and 27.

VII. Concluding Discussion

Our work here indicates that for games drawn from the ensemble that we describe here, it is possible to characterize the learning dynamics under experience-weighted attraction a priori. We have shown that a key property of such games is their

[†]In contrast to financial markets, for the behavior we observe here, the distribution of heavy tails decays exponentially (as opposed to following a power law). We hypothesize that this is because the players in financial markets use a variety of different timescales α .

